

<http://www.jyu.fi/library/tutkielmat/27/>

Jukka Salmikuukka

**AIKASARJOJEN PERUSRAKENNEMALLEISTA JA NIIDEN
SOVELTAMINEN JYVÄSKYLÄN KAUKOLÄMMÖN KULU-
TUKSEN ANALYSOINTIIN JA ENNUSTAMISEEN**

Tilastotieteen
pro gradu -tutkielma
9.5.1997

Jyväskylän Yliopisto
Tilastotieteen laitos
Informaatioteknologian maisteriohjelmat
Tilastotoimen menetelmät

Tilastotoimen menetelmien maisteriohjelma

(<http://www.stat.jyu.fi/>)

(Päivitetty 16.5.1997)

Tilastotoimen menetelmien maisteriohjelman tavoite on kouluttaa opiskelija tilastotiedon keruun, jatkojalostuksen ja käytön asiantuntijaksi nykyaikaisessa tilastojärjestelmäympäristössä, jossa datat ovat survey-, koeasetelma- tai rekisteriperusteisia. Koulutus järjestetään yhteistyönä Jyväskylän yliopiston tilastotieteen laitoksen kanssa. Tilastotieteelliseltä kannalta kyse on survey-menetelmiin, biostatistiikkaan ja ekonometriaan erikoistuneiden tilastoasiantuntijoiden koulutuksesta. Ohjelman kesto päätoimisella opiskelijalla on kaksi lukuvuotta.

Tilastotoimen ytimenä ovat tilastotieteen, erityisesti tilastotoimen teorian ja tietojenkäsittelyn kurssit. Ohjelmaan voidaan sisällyttää myös koulutustavoitteita tukevia opintojaksoja lähitieteistä (esim. talous-, yhteiskunta- ja viestintätieteet) sekä yritystoimintaan perehdyttäviä kursseja. Tarjottava opetus koostuu osittain tilastotieteen laitoksen opetusohjelmasta, erityisesti suunnitelluista tilastotoimen teorian kursseista ja ostopalveluina hankituista muista erikoiskursseista. Opettajista pääosa on kotimaasta. Vierailijoina on myös ulkomaalaisia asiantuntijoita, joten opetus on osin englanninkielistä.

Ohjelman tärkeä osa on yliopiston ulkopuolisessa yhteistyötoimipaikassa suoritettu pro gradu -tutkielma ja siihen liittyvä harjoittelu. Yhteistyötoimipaikat ovat virallisesta tilastotoimesta, suuryrityksistä tai tutkimuslaitoksista. Ne osallistuvat osaltaan harjoittelun kustannuksiin. Perusajatus on, että pro gradu -tutkielma tai osa siitä toisi tutkimustulostensa osalta lisäarvoa yhteistyötoimipaikalle. Tätä tukenee se, että harjoittelu jaetaan kahteen osaan, joista ensimmäinen eli orientoiva osa on ensimmäisen opintovuoden lopussa. Sen jälkeen opiskelija palaa yliopisto-opiskeluihin yhden lukukauden ajaksi syventääkseen tietojensa siinä tilastotieteen osa-alueessa, joka on tarpeen aiotussa pro gradu -tutkielmassa. Maisteriohjelman viimeinen lukukausi muodostaa yhtenäisen tutkimusjakson, jonka aikana tehdään pro gradu -tutkielma.

Master's Programme in Statistical Systems

(<http://www.stat.jyu.fi/>)

(Updated 16.5.1997)

The main target of the Program is to educate high qualified professionals in the collection analysis, managing and dissemination of large data sets. The Program has been built mainly on the regular curriculum of the Department of Statistics in the University of Jyväskylä. In point of view of statistical sciences, the Program in Statistical Systems concentrate in the knowledge of survey methodology, biometry and econometry. Full-time students can graduate in two Academic years. The core courses cover the advanced theory in mathematical statistics, of statical systems expecially the survey methodology and some courses in information technology. Still the program is flexible that a moderate amount of different kinds of studies in the neighboring sciences can be included as studies in economics, social sciences and communication. Most of the instructors come from Finland but visiting professors from outboard are eventual using English.

An important part of the Program crows from the cooperational research work created between the Department of Statistics and the research units located outside the University of Jyväskylä. Those units are research and development departments in the bodies of official statistics, big business firms and research institutes. They share the costs of practices. Basic idea is that the MS.c. Thesis wil be written from the topics given by the cooperational research units and so the results could be of some contributed value to them. As an appropriate mode of policy may be here thus that the time of practice skills for the Thesis. This is proceeded in 6 final months of the Program, in collaboration with the same cooperating research unit.

Sisällysluettelo

	sivu
Merkinnot ja lyhenteet	iii
1. Johdanto	1
1.1 Työn tarkoitus	1
1.2 Aikasarjan luonne	1
1.3 Aikasarjateorian vaiheista	2
1.4 Yleistä	3
2. Rakenneyhtälömallit	4
2.1 Johdanto rakenneyhtälömalleihin	4
2.2 Perusrakennemalli ja sen laajennuksia	5
2.2.1 Mallin rakenne	5
2.2.2 Perusrakennemallin laajennuksia	7
2.3 Aikasarjan stationarisointi	8
2.3.1 Erilaisia malleja	8
2.3.2 Aikasarjan stationarisointi	11
3. Tila-aika-malli ja Kalmanin suodin	14
3.1 Tila-aika-malli	14
3.1.1 Mallin yleinen muoto	14
3.1.2 Mallin rakenteesta	15
3.2 Kalmanin suodin	16
3.2.1 Kalmanin suotimen tehtävä	16
3.2.2 Kalmanin suotimen yleinen muoto	16
3.2.3 Kalmanin suotimen ominaisuuksista	17
3.3 Aikainvarianttisuus	18
3.4 Rakenneyhtälömallin yhteys uskottavuusfunktioon	19
4. Estimointi	22
4.1 Perusrakennemalli ja tila-aikamuoto	22
4.2 Parametrien estimointi	23
4.2.1 Hyperparametrien estimointi	23
4.2.2 Uskottavuusfunktion maksimointi	24
4.2.3 Deterministiset komponentit	26
4.3 Huomautuksia	26
5. Mallin spesifiointi, testaaminen ja diagnostiikka	27
5.1 Mallin spesifiointi	27
5.2 Testaaminen ja diagnostiikka	28
5.2.1 Testaamisen perusteet	28
5.2.2 Mallin diagnostiset testit	29
6. Ennustaminen	36
6.1 Perusrakennemallilla ennustaminen	36
6.2 Huomautuksia	37

7. Perusrakennemallin käytöstä	38
7.1 Ennustaminen	38
7.2 Aikasarjatasoituksista	38
7.3 Trendipuhdistus	39
7.4 Kausipuhdistus	39
7.5 Komponenttikohtainen tarkastelu	40
7.6 Huomautuksia	40
8. Selittävät muuttujat	41
8.1 Johdanto	41
8.1.1 Eksogeenisuus	41
8.1.2 Erilaisia selitysmalleja	42
8.2 Selittävät muuttujat ja parametrien estimointi	44
8.3 Selittävien muuttujien valinta ja mallin diagnostiikka	46
8.4 Interventiot	47
9. Empiirinen ongelma ja tutkimusaineisto	50
9.1 Empiirinen ongelma	50
9.1.1 Tutkimuksen tavoitteet	50
9.1.2 Kaukolämpö	50
9.1.3 Kaukolämpöterminologiaa	51
9.1.4 Kaukolämmön toimintaperiaate	51
9.2 Tutkimusaineisto	52
9.2.1 Aineistosta yleisesti	52
9.2.2 Analysoitavat muuttujat	52
9.2.3 Aineistoon tehdyistä korjauksista	54
10. Analyysitulokset	55
10.1 Kaukolämmön kokonaiskulutusennuste vuodelle 1996	55
10.1.1 Yksinkertainen perusrakennemalli	55
10.1.2 Regressiomalli	59
10.1.3 Laajennettu perusrakennemalli	61
10.1.4 Mallien vertailua	63
10.2 Kulutusennusteet vuodelle 1997	65
11. Yhteenveto	69
Liite 1. Havaintoaineisto	71
Lähteet	74

Merkinnät ja lyhenteet

Merkintöjä

- a) Matriisit ja vektorit merkitään **lihavoituna**, matriisit isolla, vektorit pienellä kirjaimella.
- b) Parametrit ja tilat merkitään kreikkalaisin kirjaimin, parametrien lineaariestimaattorit lähimmällä vastaavalla latinan kirjaimella. Tilde (\sim) kreikkalaisen kirjaimen päällä merkitsee, että kyseessä on parametrin SU-estimaattori tai tilan MMSE-estimaattori. Muulloin käytetään sen sijaan estimaatin symbolina kirjaimen yläpuolella olevaa hattua (\wedge).
- c) Ajan hetkeä ilmentävä alaindeksi t osoittaa että on käytetty informaatio hetkeen t saakka (t mukaan lukien). Estimaattorin tai estimaatin ehdollinen merkintä, esimerkiksi $t|\tau$, kertoo puolestaan, että estimaatti perustuu hetkellä τ käytettävissä olleeseen informaatioon. Tässä tapauksessa τ voi olla suurempi tai pienempi kuin t .
- d) Merkintä 'log' merkitsee aina luonnollista logaritmia.

Lyhenteitä

AIC	Akaiken informaatiokriteeri (Akaike information criterion)
akf	autokorrelaatiofunktio (autocorrelationfunction)
ARIMA	Autoregressiivinen integroitu liukuvan keskiarvon (Auto-Regressive Integrated Moving Average) malli
BIC	Bayesin informaatiokriteeri (Bayes information criterion)
CUSUM	kumulatiivinen summa (cumulative sum)
D-W	Durbin-Watsonin testisuure
KS	Kalmanin suodin (Kalman filter)
MMSE	keskineliöpoikkeaman minimoiva estimaattori, -estimaatti (minimum mean square error estimator, -estimate)
MSE	keskineliöpoikkeama (mean square error)
NID(*,**)	riippumattomasti normaalijakautunut odotusarvolla * ja varianssilla **
p.e.v.	ennustevirheen varianssi (prediction error variance)
RMSE	keskineliövirhe (root mean square error)
WLS	painotettu pns-menetelmä (weighted least squares)

Aikasarjojen perusrakennemalleista ja niiden soveltaminen Jyväskylän kaukolämmön kulutuksen analysointiin ja ennustamiseen

Tilastotieteen pro gradu -tutkielma
9.5.1997

Tekijä: Jukka Salmikuukka
Ohjaaja: prof. Esko Leskinen
Julkaisija: Jyväskylän yliopisto
Sivumäärä: 76 sivua

Abstrakti

Tässä tutkielmassa tarkastellaan aikasarjojen rakenneyhtälömalleja perusrakennemallien osalta. Perusrakennemallien teoria esitetään vaihe vaiheelta ja tämän jälkeen teoriaa sovelletaan empiiriseen ongelmaan, joka on kaukolämmön kulutuksen ennustaminen.

Perusrakennemallien keskeisenä ajatuksena on esittää tutkittava ilmiö erilaisten rakennekomponenttien avulla. Keskeisenä erona aikaisempiin menetelmiin verrattuna on komponenttien stokastisointi. Perusrakennemallien teoria linkitetään klassiseen estimointiteoriaan tila-aika -malliin perustuvan Kalmanin suotimen avulla. Niin kutsuttua ennustevirrehajotelmaa hyödyntäen mallille saadaan muodostettua uskottavuusfunktio, jonka myötä klassinen testiteoria on käytettävissä. Näin ollen voidaan täsmällisesti määrittää estimoitavan mallin komponenttien merkitsevyytasot. Perusrakennemallien teoriakehikon joustavuudesta johtuen malleihin voidaan sisällyttää myös lisäinformaatiota erityyppisten selittävien muuttujien avulla.

Empiirisessä osassa analysoidaan Jyväskylän kaukolämmön kulutusta vuosina 1989 - 1996 kuvaavaa kuukausitason aikasarjaa. Sarjasta muodostetaan erilaisia malleja, joista osassa hyödynnetään lisäinformaationa yleistä kiinteistöjen lämmitystarvetta kuvaavaa astepäivälukujen aikasarjaa. Aluksi mallit muodostetaan hyödyntäen havaintoaineistoa ainoastaan vuoden 1995 loppuun saakka. Muodostetuilla malleilla tuotetaan sitten ennusteet vuodelle 1996. Saatuja ennusteita verrataan sitten toteutuneeseen kulutukseen ja voidaan arvioida erilaisten mallien ominaisuuksia. Tämän jälkeen hyödynnetään koko aineistoa ja tuotetaan aito ennuste vuoden 1997 kaukolämmön kulutukselle.

Havaitaan, että perusrakennemallien avulla pystytään selittämään jopa 90% kulutuksen kokonaisvaihtelusta. Parhaaseen tulokseen päästään lisäinformaatiota hyödyntämällä ja vuoden 1997 kaukolämmön kulutukselle saadaan hienoista kasvua ennakoiva kulutusennuste, joka perustuu oletukselle pitkän aikavälin keskiarvoja noudattavista ulkolämpötiloista.

Summary

This study presents the theory and practice of basic structural time-series models. First, the theory of basic structural models is shown and after that follows the practical application of the theory. The empirical application deals with forecasting district heating consumption in the city of Jyväskylä.

Basic structural models are formulated in terms of structural components. The main difference compared to former component methods is the fact, that the components in basic structural models are stochastic. The theory is linked to the classical estimation theory by Kalman filter, which is a statistical algorithm based to the state-space model. By using so called *prediction error decomposition* the likelihood function can be produced. After that the classical testing theory can be used to estimate the significance of the components. Because of the flexibility of the basic structural model theory, models can be expanded with variables containing auxiliary information.

In the empirical application the task is to analyze the consumption of energy transmitted by district heating system in Jyväskylä. The data contains monthly observations from January 1989 to December 1996. Also auxiliary information will be used in some models. This is the monthly count of degree days ($^{\circ}\text{Cd}$). First the models will be made based to the information from the January 1989 to December 1995 and forecasting period is 12 months forward. Then it is possible to compare the forecasts to the real values and the best model can be found. After that the model based to the full information is made and best possible forecasts can be made for the year 1997.

As a result it will be noticed, that by basic structural models almost 90% of the consumption's total variation can be explained. The model with auxiliary information gives the most reliable forecasts. To the year 1997 the model gives a forecast which indicates a slight growth in energy consumption assumed that $^{\circ}\text{Cd}$ follows its long-time average.

1. Johdanto

1.1 Työn tarkoitus

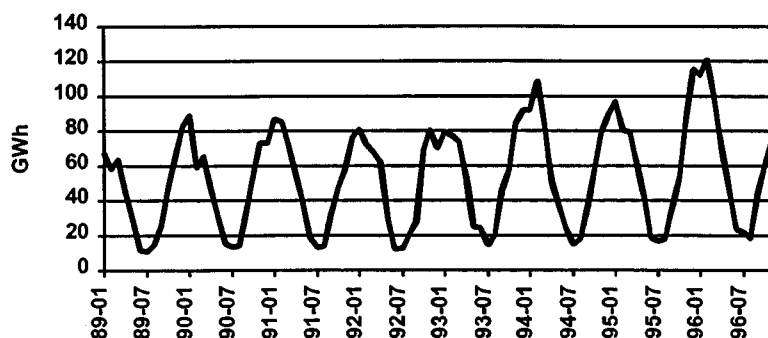
Tässä tilastotieteen pro gradu-työssä esitellään aikasarjojen rakenneyhtälömallien teoriaa (Harvey, 1989) ja sovelletaan sitä kaukolämmön kulutuskehityksen arviointiin. Työ jakautuu kahtia siten, että ensin esitetään teoreettinen perusta, jota sitten sovelletaan tutkimusaineistoon. Empiirisessä osassa analysoidaan ja ennustetaan Jyväskylän Energia Oy:n kaukolämmön kulutusta soveltaen rakenneyhtälömalleja. Analysointi ja ennustaminen suoritetaan käyttäen apuna **STAMP-tietokoneohjelmistoa** (STAMP 5.0 *Structural Time series Analyser, Modeller and Predictor*, Koopman et al., 1995).

Työssä rajoitutaan lähinnä niin kutsuttujen **perusrakennemallien** (*basic structural model*) tarkasteluun. Aluksi tarkastellaan aikasarjan käsitettä, aikasarja-analyysin kehittymistä ja perusrakennemallin peruseriaatteita (luvut 1 ja 2). Luvussa kolme käydään läpi rakenneyhtälömalleille keskeisen **Kalmanin suodimen** (*Kalman filter*) toiminta, joka perustuu nk. **tila-aika-malliin** (*state-space-model*). Kalmanin suodin mahdollistaa mallien estimoinnin, ennustamisen ja analysoinnin (luvut 4, 5 ja 6). Myös erilaisia perusrakennemallien käyttötapoja käsitellään (luku 7). Rakenneyhtälömallien teoria mahdollistaa myös selittävien muuttujien lisäämisen malliin, tätä tarkastellaan luvussa 8.

Luvussa 9 esitellään empiirinen ongelma, joka käsittelee Jyväskylän Energia Oy:n Jyväskylään myymän kaukolämmön kulutuksen ennustamista vuodelle 1997. Samalla tarkastellaan kaukolämmön käsitteitä ja toimintaperiaatteita ja esitellään tutkimusaineisto, jonka jälkeen suoritetaan analysointi ja esitetään analyysitulokset (luku 10). Lopuksi (luku 11) arvioidaan sekä saatuja empiirisiä tuloksia, että teoriakehikkoa kokonaisuudessaankin.

1.2 Aikasarjojen luonne

Aikasarja on kokoelma tasaisin väliajoin eri ajan hetkinä mitattuja samaa ilmiötä mittaavia havaintoja y_t , $t = 1, \dots, T$. Esimerkiksi Jyväskylän kuukausittainen kaukolämmön kulutus gigawattitunteina (kuvio 1.1) sisältää kaupungin kaukolämpöverkkoon kuuluvien kiinteistöjen yhteenlasketun lämmönkulutuksen tarkasteluajavälillä.



Kuvio 1.1 Kaukolämmön kulutus Jyväskylässä 1/1989 - 12/1996

Kun aikasarjaa tarkastellaan lähemmin, siitä voidaan havaita erilaisia systemaattisia piirteitä. Esimerkkisarjassa ensimmäinen huomiota kiinnittävä seikka näyttäisi olevan lämmönkulutuksen systemaattinen vaihtelu vuodenaikojen mukaan. Tätä piirrettä kutsutaan kausivaihteluksi. Osa vaihtelusta johtuu siitä, että lämmitettäviä rakennuksia tulee lisää ja jo lämmitettävien talojen kunnan huononeminen lisäävät systemaattisesti vuosittaista kokonaislämmönkulutusta, eli kyseessä on siis kasvava trendi. Huomiota kiinnittävä piirre on varmasti myös se, että saman kuukauden havainnot poikkeavat toisistaan eri tarkasteluvuosina. Tämän eron aiheuttaa osittain sääolojen vaihtelu eri vuosina ja osa eroista tulee satunnaisten tekijöiden summana. Lukuisat pienet lämmönkulutukseen vaikuttavat tekijät yhdessä vaihtelevien sääolojen kanssa tuovat oman vaikutuksensa, joka on niin kutsuttu satunnaiskomponentti. Näin siis saadaan kaukolämmön kulutusta kuvaavaksi malliksi additiivinen malli:

$$\text{havaittu lämmönkulutus} = \text{trendi} + \text{kausivaihtelu} + \text{satunnaiskomponentti}. \quad (1.1)$$

Joissakin tapauksissa mallin voidaan ajatella olevan kerrannainen:

$$\text{havaittu lämmönkulutus} = \text{trendi} \times \text{kausivaihtelu} \times \text{satunnaiskomponentti}. \quad (1.2)$$

Kun näin on, voidaan malli muuntaa muodon (1.1) kaltaiseksi suorittamalla aikasarjan logaritmisointi.

Esittämällä aikasarja erillisten komponenttien avulla saadaan kiinnostuksen kohteena olevasta sarjasta informaatiota, jota ei välttämättä aikasarjaa silmämääräisesti tarkastelelemalla havaita. Luonnollisesti aikasarjaa voi silmämääräisestäikin 'analysoida'. Pelkästään tällä tavalla ei usein kuitenkaan pystytä saamaan riittävän tarkkaa tietoa tutkittavasta ilmiöstä. Tämän vuoksi aikasarjan tilastollinen analysointi on tarpeen, jotta pystytään täsmällisemmin selvittämään, kuinka ilmiö käyttäytyy.

1.3 Aikasarjateorian vaiheista

Aikasarjojen historia ulottaa tietyvästi ainakin 1800-luvun alkupuolelle, jolloin muuan William Playfair, jollei ensimmäisenä niin ensimmäisten joukossa esitti aikasarjoja nykyisessä, modernissa muodossa (Kendall, 1973, s 2). Aikasarjoja on analysoitu jollakin tapaa jo Playfairin ajoista lähtien, joskin aikasarja-analyttisen tutkimuksen alkuna voidaan pitää vuotta 1927, jolloin auringon pilkkujen tutkija Udny Yule loi pohjan stokastisen prosessin käsitteelle (Kendall, 1973, s 4). Aikasarjojen analysoinnin keskeisenä tavoitteena on tuottaa ennusteita kohteena olevasta ilmiöstä. Yksinkertaisimmillaan ennustaminen tapahtuu sovittamalla jokin ajan funktio aineistoon ja ekstrapoloimalla tulevaisuuteen. Ensimmäisenä *ad hoc* ennustamisproseduurina voidaan mainita **eksponentiaalisesti painotetun liukuvan keskiarvon menetelmä** (*EWMA, exponentially weighted moving average*), jonka Harveyn mukaan kehittivät Holt (1957) ja Winters (1960). He toivat kulmakerroin-komponentin ennustamisfunktioon ja hyväksyivät kausivaihtelutekijöiden olemassaolon. Toisenlaisen lähestymistavan asiaan omaksui Brown (1963), joka vei ennustamisproseduurin regressioympäristöön ja kehitti k. painotetun pns-menetelmän (*discounted least squares*). Harveyn mukaan Muth esitti 1960 ensimmäisenä täsmällisen perustelun EWMA:lle käyttäen apuna **satunnaiskulkua plus kohina -mallia** (*random walk plus noise model*). Theil ja Wage (1964) ja Nerlove ja Wage (1964) täydensivät mallia siten, että siihen sisältyi kulmakerrointermi.

1960-Luvulla kehitys pysähtyikin tähän, pääosin kehittymättömän laskentakoneiston takia. Tilastotieteilijöiltä jäi myös joksikin aikaa huomioimatta Schweppen (1965) esittämä **Kalmanin suotimeen perustuva uskottavuusfunktion ratkaisumenetelmä**, jonka soveltaminen olisi tosin ollut melkoisen vaikea juuri laskentakapasiteetin puutteen vuoksi (Harvey 1989, s 22 - 23).

Merkittävimmästä aikasarja-analyttisestä työstä 1960-luvulla vastasivat Box ja Jenkins (1970, 1976), jotka formuloivat **ARIMA-malliteorian** Holtin ja Wintersin sekä Muthin esityksiä laajentaen. Rakenneyhtälömallien osalta Harrison ja Stevens (1971, 1976) jatkoivat työtä ja saavuttivatkin merkittävää edistystä Kalmanin suotimen osalta. Vaikka ARIMA-malliajattelu dominoikin aikasarja-analyttistä ajattelua koko 1970-luvun ja 80-luvun alun, Kalmanin suotimen mahdollisuudet ekonometristen ja tilastollisten ongelmien ratkaisussa havaittiin 1970-luvulla mm Rosenbergin (1973), Englen (1978), Harveyn ja Phillipsin (1979) sekä Garbaden (1977) toimesta ja 1980-luvun lopulla päästiin mm Harveyn (1989) toimesta kokonaisvaltaiseen tilastolliseen esitykseen rakenneyhtälömallien osalta.

1.4 Yleistä

Aikasarjojen analysoinnissa voidaan siis erottaa kaksi eri lähestymistapaa, toisaalta nk. Boxin ja Jenkinsin esittämä ARIMA-malleihin perustuva lähestymistapa (Box & Jenkins, 1970, 1976), ja toisaalta tässä työssä tarkasteltavana oleva rakenneyhtälömalleihin perustuva ajattelutapa (esim. Harvey, 1989). Nämä kaksi eivät ole kuitenkaan toisensa poissulkevia, vaan ARIMA-mallit voidaan esittää myös rakenneyhtälömallien avulla. ARIMA-mallien lähestymistapaan verrattuna rakenneyhtälömalliajattelussa on kuitenkin joitakin merkittäviä poikkeuksia.

ARIMA-mallin rakentaminen perustuu oletukselle, jonka mukaan haluttu aikasarja saadaan stationarisoitua differensoimalla riittävän monta kertaa. Tämän jälkeen stationarisoitu sarja esitetään autoregressiivisten (AR-) sekä liukuvan keskiarvon (MA-) komponenttien avulla (ARMA-malli). Jo oletus, että differensoimalla aina päästään stationaariseen tulokseen, on kuitenkin kyseenalainen. Lisäksi, vaikka stationarisointi tällä tavoin onnistuisikin, oikean ARMA-mallin valinta voi olla hankalaa (Harvey, 1989 s 80-81).

Rakenneyhtälömallilla pyritään kuvaamaan aikasarjan keskeiset ominaisuudet, kuten esimerkiksi trendi tai sykli toisistaan erillisinä, stokastisina komponentteina. Juuri komponenttien stokastisointi on selkeä laajennus verrattuna ARIMA-malleihin, joissa jäännöstermiä lukuunottamatta muut komponentit ovat deterministisiä. Stationaarisuuteen differensoimalla pyrkiminen ei ole rakenneyhtälömalleissa niin keskeisessä roolissa kuin ARIMA-mallien tapauksessa. Mallin etuihin kuuluu vielä se, että malliin voidaan sisällyttää selittäviä muuttujia regressioanalyysin tapaan.

Rakenneyhtälömallien taustalla oleva teoriakehikko antaa erittäin monipuoliset mahdollisuudet erilaisiin vertailuihin ja jatkoanalyysiin, kuten myöhemmin saamme huomata.

2. Rakenneyhtälömallit

2.1 Johdanto rakenneyhtälömalleihin

Tässä työssä tarkastellaan yhden muuttujan rakenneyhtälömalleja. Aluksi määritellään niin kutsuttu **perusrakennemalli** (*basic structural model*), joka on myös laajennettavissa joiltakin osin. Ennen tätä on kuitenkin syytä tehdä joitakin apumäärittelyjä.

Viiveoperaattori L on keskeisessä roolissa aikasarjojen matemaattisessa formuloinnissa. Se voidaan määritellä muodossa

$$Ly_t = y_{t-1}, \quad (2.1)$$

joka yleistyy muotoon

$$L^\tau y_t = y_{t-\tau}. \quad (2.1b)$$

Viiveoperaattori L :n avulla voidaan määritellä **differenssioperaattori** Δ . Ensimmäinen differenssi differenssioperaattorilla esitettynä on

$$\Delta = 1 - L. \quad (2.2)$$

Differensoimalla ensimmäinen differenssi uudelleen saadaan toisen asteen differenssi

$$\Delta^2 = (1 - L)^2 = 1 - 2L + L^2, \quad (2.3a)$$

ja vastaavalla tavalla differenssioperaattorin yleinen muoto

$$\Delta^d = (1 - L)^d \quad (2.3b)$$

on d :nnen asteen differenssi. Differenssi viiveellä τ esitetään puolestaan muodossa

$$\Delta_\tau = 1 - L^\tau, \quad \tau = 1, 2, \dots \quad (2.4a)$$

eli

$$\Delta_\tau y_t = y_t - y_{t-\tau}. \quad (2.4b)$$

Kausivaihteludifferenssi on differenssioperaattorin erikoistapaus, siinä differensoinnin viive τ on sama kuin kausivaihtelujaksojen s lukumäärä.

Astetta τ oleva **summaoperaattori** on muotoa

$$S_\tau(L) = 1 + L + L^2 + \dots + L^{\tau-1}, \quad \tau = 1, 2, \dots, \quad (2.5)$$

ja näin ollen havaintojen y_t ($t = \tau, \tau+1, \dots, T$) summa voidaan esittää muodossa

$$S_\tau(L)y_t = \sum_{j=0}^{\tau-1} y_{t-j}. \quad (2.6)$$

Kausivaihtelun summausoperaattoriksi kaavan (2.6) operaattori muodostuu silloin, kun summattava jakso vastaa kausivaihtelun pituutta, eli $\tau = s$.

Trigonometrinen operaattori on

$$\gamma(L) = 1 - (2 \cos \lambda)L + L^2, \quad (2.7)$$

jossa λ on jokin ennalta määrätty radiaaneina mitattu taajuus väliltä $0 < \lambda < \pi$. $\gamma(L)$:n juurten tulo on yksi ja vaihekulma λ .

Viimeiseksi määritellään vielä summausoperaattorin ja trigonometrinen funktioiden välinen yhteys. Olkoon $\lambda_c = 2\pi j/\tau$, kun $j = 1, \dots, [\tau/2]$. Edellisessä

$$[\tau/2] = \begin{cases} \tau/2 & \text{kun } \tau \text{ parillinen} \\ (\tau-1)/2 & \text{kun } \tau \text{ pariton} \end{cases} \quad (2.8)$$

Kun τ on pariton, trigonometrinen operaattori on

$$\gamma(L) = 1 - (2 \cos \lambda)L + L^2, \quad j = 1, \dots, [\tau/2]. \quad (2.9a)$$

Kun τ on parillinen ja $j = 1, \dots, \tau/2 - 1$, operaattori määritellään kuten (2.9), ja kun $j = \tau/2$ operaattori on muotoa

$$\gamma_{\frac{\tau}{2}}(L) = 1 + L. \quad (2.9b)$$

Summausoperaattori voidaan kirjoittaa edellisten avulla muotoon

$$S_\tau(L) = \prod_{j=1}^{[\tau/2]} \gamma_j(L). \quad (2.10)$$

2.2 Perusrakennemalli ja sen laajennuksia

2.2.1 Mallin rakenne

Yhden muuttujan malleissa tarkastellaan tilannetta, jossa selitettävänä on yksi yksittäinen aikasarja ja merkitään tätä y_t :llä. Määritellään kaavan (1.2.1) mukainen **perusrakennemalli**

$$y_t = \mu_t + \gamma_t + \varepsilon_t, \quad \varepsilon_t \sim \text{NID}(0, \sigma_\varepsilon^2). \quad (2.11)$$

Trendikomponentilla μ_t kuvataan sarjan pitkän aikavälin kehitystä, jotka voidaan ekstrapoloida tulevaisuuteen. Trendi määritellään muodossa

$$\mu_t = \mu_{t-1} + \beta_{t-1} + \eta_t, \quad \eta_t \sim \text{NID}(0, \sigma_\eta^2). \quad (2.12)$$

Trendikomponentti rakentuu kahdesta osasta, trendikomponentin edellisen ajankohdan tasoparametrista μ_{t-1} ja trendin kulmakerroinkomponentista β_{t-1} , joka voidaan määrittellä muodossa

$$\beta_t = \beta_{t-1} + \zeta_t, \quad \zeta_t \sim \text{NID}(0, \sigma_\zeta^2). \quad (2.13)$$

Kausivaihtelukomponentti voidaan esittää muodossa

$$\gamma_t = -\sum_{j=1}^{s-1} \gamma_{t-j} + \omega_t, \quad \omega_t \sim \text{NID}(0, \sigma_\omega^2), \quad (2.14)$$

jolle yleensä asetetaan rajoitus

$$\sum_{j=0}^{s-1} \gamma_{t-j} = 0, \quad (2.15)$$

eli

$$S_s(L)\gamma_t = 0. \quad (2.16)$$

Vaihtoehdoisen esitystavan ovat esittäneet Harrison ja Stevens (1976, s 217-218). Nyt oletetaan, että kullakin kausivaihteluterminillä on oma satunnaisprosessinsa, eli

$$\sum_{j=1}^s \gamma_{jt} = 0 = \sum_{j=1}^s \omega_{jt}. \quad (2.17)$$

Lausekkeen (2.15) stokastinen muoto perustuu kiinteään kausivaihtelun mallittamiseen perinteisillä dummy-muuttujilla. Harvey'n mukaan vaihtoehtoinen tapa on suorittaa mallitus taajuusalueen avulla (Hannan, Terrell ja Tuckwell, 1970). Deterministinen kausivaihtelu hetkellä t on tällöin muotoa

$$\gamma_t = \sum_{j=1}^{[s/2]} (\gamma_j \cos \lambda_j t + \gamma_j^* \sin \lambda_j t). \quad (2.18)$$

Kun s on parillinen ja $j = s/2$, sini-termi on nolla. Näin ollen trigonometrinen parametrien γ_j ja γ_j^* lukumäärä lausekkeessa (2.18) on aina $s - 1$, mikä on sama kuin dummy-muuttujiin perustuvan kausivaihtelumallituksen tuottamien kausivaihtelukertoimien lukumäärä. On myös helposti osoitettavissa, että yhtälö (2.16) on voimassa standardeja trigonometrisia määrittelyjä käytettäessä. Voidaan myös osoittaa, että tietyissä tapauksissa kausivaihtelun trigonometrinen mallitus ja mallitus dummy-muuttujan avulla johtavat täysin identtisiin kausivaihtelukertoimi-

hin (Harvey, 1989 s 42). Tämä edellyttää kuitenkin mm. sitä, että trigonometrinen kausivaihtelukomponentti on täyttä astetta. Käytännössä kausivaihtelurakenne muuttuu kuitenkin niin hitaasti, että usein kausivaihtelukomponentista voidaan jättää pois joitakin korkeimpia taajuuksia. Joissakin tapauksissa riittää yhden taajuusalueen parametri, esimerkiksi kuukausitaso-aineistossa kahdentoista kuukauden jakson kattava parametri $\gamma_1 = 2\pi/12$.

Trigonometrinen kausivaihtelukomponentti voidaan stokastisoida lisäämällä mukaan satunnaiskomponentit ω_{jt} ja ω_{jt}^* , molempien noudattaessa jakaumaa $NID(0, \sigma_j^2)$. Jakauman ei tarvitsisi välttämättä olla molemmille parametreille γ_j ja γ_j^* sama, mutta käytännössä varianssien yhdistämisestä saatavat hyödyt ovat Harveyn mukaan haittoja suuremmat. Stokastinen kausivaihtelukomponentti taajuusalueen avulla esitettynä on siten

$$\begin{aligned} \gamma_{jt} &= \gamma_{j,t-1} \cos \lambda_j + \gamma_{j,t-1}^* \sin \lambda_j + \omega_{jt} \\ \gamma_{jt}^* &= -\gamma_{j,t-1} \sin \lambda_j + \gamma_{j,t-1}^* \cos \lambda_j + \omega_{jt}^*, \end{aligned} \quad j = 1, \dots, [s/2], \quad (2.19)$$

jossa siis $\omega_{jt}, \omega_{jt}^* \sim NID(0, \sigma_j^2)$.

2.2.2 Perusrakennemallin laajennuksia

Sykli määritellään radiaaneina mitatun taajuusalueen avulla. Olkoon ψ_t **syklikomponentti**, joka deterministisessä muodossa kosini-käyrän avulla esitettynä on

$$\psi_t = A \cos(\lambda_c t - \theta), \quad t = 1, \dots, T, \quad (2.20)$$

jossa syklin amplitudi on A ja θ on vaihekulma. Syklin kesto on puolestaan $2\pi/\lambda_c$. Käyttökelpoisempi muotoilu on syklin esittäminen sini- ja kosini-käyrien yhdistelmänä, eli

$$\psi_t = \alpha \cos \lambda_c t + \beta \sin \lambda_c t, \quad (2.21)$$

jossa $(\alpha^2 + \beta^2)^{1/2}$ on amplitudi ja $\tan^{-1}(\beta/\alpha)$ on vaihekulma. Stokastisoidaan systeemi sallimalla parametrien α ja β vaihdella eri ajan hetkinä. Ennen stokastisointia varmistetaan syklin jatkuvuus kirjoittamalla syklimalli muotoon

$$\begin{bmatrix} \psi_t \\ \psi_t^* \end{bmatrix} = \begin{bmatrix} \cos \lambda_c & \sin \lambda_c \\ -\sin \lambda_c & \cos \lambda_c \end{bmatrix} \begin{bmatrix} \psi_{t-1} \\ \psi_{t-1}^* \end{bmatrix}, \quad t = 1, \dots, T \quad (2.22)$$

ja jossa $\psi_0 = \alpha$ ja $\psi_0^* = \beta$. Näin saadaan parametrit ψ_t ja ψ_t^* , joista ensimmäinen on syklin arvo hetkellä t ja jälkimmäistä käytetään konstruoidessa ensimmäistä. Stokastisoidaan (2.22) lisäämällä siihen kaksi **valkoisen kohinan** (*white noise*) komponenttia, joka johtaa muotoon

$$\begin{bmatrix} \psi_t \\ \psi_t^* \end{bmatrix} = \begin{bmatrix} \cos \lambda_c & \sin \lambda_c \\ -\sin \lambda_c & \cos \lambda_c \end{bmatrix} \begin{bmatrix} \psi_{t-1} \\ \psi_{t-1}^* \end{bmatrix} + \begin{bmatrix} \kappa_t \\ \kappa_t^* \end{bmatrix}. \quad (2.23)$$

Jotta malli olisi identifioitavissa, on oletettava, että kohinakomponentit ovat keskenään korreloimattomia, tai että niillä on sama varianssi. Käytännössä kuitenkin selvyyden vuoksi tehdään nämä molemmat oletukset. Syklistä voidaan muokata vieläkin joustavampi lisäämällä siihen **vaimentava tekijä** (*damping factor*) ρ , jolloin (2.23) laajenee:

$$\begin{bmatrix} \psi_t \\ \psi_t^* \end{bmatrix} = \rho \begin{bmatrix} \cos \lambda_c & \sin \lambda_c \\ -\sin \lambda_c & \cos \lambda_c \end{bmatrix} \begin{bmatrix} \psi_{t-1} \\ \psi_{t-1}^* \end{bmatrix} + \begin{bmatrix} \kappa_t \\ \kappa_t^* \end{bmatrix}, \quad (2.24)$$

jossa $0 \leq \rho \leq 1$. Näin malli (2.24) on vektoriarvoinen AR(1)-prosessi.

Syklin erikoistapauksena voidaan mainita tilanne, jossa $\lambda_c = 0$ tai π . Tällöinhän $\sin \lambda_c = 0$ ja yhtälössä oleva ψ_t^* on merkityksetön ja (2.24) supistuu muotoon

$$\psi_t = \rho \psi_{t-1} + \kappa_t, \quad (2.25a)$$

kun $\lambda_c = 0$ ja

$$\psi_t = -\rho \psi_{t-1} + \kappa_t, \quad (2.25b)$$

kun $\lambda_c = \pi$. Sykli supistuu silloin yksinkertaiseksi AR(1)-prosessiksi.

Perusrakennemallia voidaan laajentaa myös huomioimalla **päivittäisvaikutukset** (*daily effects*). Päivittäisvaikutusten mallittaminen tulee kyseeseen, kun käytettävissä on päivittäishavaintoaineisto, jossa halutaan huomioida erilaisten viikonpäivien erilainen vaikutus, esimerkiksi tilanne, jossa arkipäivät maanantaista perjantaihin ovat keskenään samanlaisia, mutta viikonloput poikkeavat näistä.

2.3. Aikasarjan stationarisointi

2.3.1 Erilaisia malleja

Luvussa 2.2.1 tarkasteltiin jo niin kutsuttua perusrakennemallia. Tässä luvussa tarkastellaan miten eri mallikomponentit yhdistetään kokonaiseksi rakenneyhtälömalliksi ja miten malli on esitettävissä yhtenä yksittäisenä yhtälönä.

Paikallinen tasomalli (*local level/random walk plus noise*) voidaan esittää muodossa

$$y_t = \mu_t + \varepsilon_t, \quad t = 1, \dots, T, \quad \varepsilon_t \sim \text{NID}(0, \sigma_\varepsilon^2) \quad (2.26a)$$

jossa trendikomponentti, μ_t , on yksinkertaisesti taso, joka vaihtelee ylös- ja alaspäin satunnaiskävelyn (*random walk*) mukaan.

$$\mu_t = \mu_{t-1} + \eta_t, \quad t = 1, 2, \dots, \quad \eta_t \sim \text{NID}(0, \sigma_\eta^2) \quad (2.26b)$$

Alkuarvoa ei tarvitse erikseen määrittellä, sillä voidaan olettaa, että prosessi on alkanut jollain hetkellä menneisyydessä. Malli voidaan myös kirjoittaa muodossa

$$\mu_t - \mu_{t-1} = \Delta\mu_t = \eta_t, \quad (2.26c)$$

josta on helposti ratkaistavissa

$$\mu_t = \frac{\eta_t}{\Delta}, \quad (2.26d)$$

joka sijoitetaan yhtälöön (2.26a), saadaan

$$y_t = \frac{\eta_t}{\Delta} + \varepsilon_t, \quad t = 1, \dots, T. \quad (2.27)$$

Lisättään nyt edelliseen malliin kulmakerrointermi, jolloin saadaan muodostettua **paikallinen lineaarinen trendimalli** (*local linear trend*), jossa tasomalliin on lisätty kulmakerroinkomponentti β_t . Malli on muotoa

$$\begin{aligned} \mu_t &= \mu_{t-1} + \beta_{t-1} + \eta_t \\ \beta_t &= \beta_{t-1} + \zeta_t, \end{aligned} \quad (2.28a)$$

josta kulmakerroinkomponentti voidaan esittää muodossa

$$\beta_{t-1} = \zeta_{t-1}/\Delta, \quad (2.28b)$$

jonka jälkeen ratkaistaan trendikomponentti vastaavasti, saadaan

$$y_t = \frac{\eta_t}{\Delta} + \frac{\zeta_{t-1}}{\Delta^2} + \varepsilon_t, \quad t = 1, \dots, T. \quad (2.29a)$$

Mikäli kulmakerroinkomponentin varianssi $\sigma_\zeta^2 = 0$, on kulmakerroin deterministinen, eli β on kiinteä estimoitava parametri. Yhtälö supistuu muotoon

$$y_t = \frac{\eta_t}{\Delta} + \beta + \varepsilon_t. \quad (2.29b)$$

Stokastinen syklisten komponenttien malli (*the cycle plus noise*) on muotoa

$$y_t = \mu_t + \psi_t + \varepsilon_t, \quad t = 1, \dots, T. \quad (2.30)$$

Jossa ψ_t on stokastinen prosessi (2.23). Jäännöstermin ε_t oletetaan olevan korreloimaton termien κ_t ja κ_t^* kanssa. Yhden yhtälön esitysmuotoon päästään käsittelemällä syklikomponenttia vektoriarvoisena AR(1)-prosessina, merkitään

$$\begin{bmatrix} \psi_t \\ \psi_t^* \end{bmatrix} = \begin{bmatrix} 1 - \rho \cos \lambda_c \cdot L & -\rho \sin \lambda_c \cdot L \\ \rho \sin \lambda_c \cdot L & 1 - \rho \cos \lambda_c \cdot L \end{bmatrix}^{-1} \begin{bmatrix} \kappa_t \\ \kappa_t^* \end{bmatrix}, \quad (2.31)$$

jossa L on viivästysoperaattori. Sijoittamalla edellinen yhtälöön (2.30) saadaan

$$y_t = \mu_t + \frac{(1 - \rho \cos \lambda_c L) \kappa_t + (\rho \sin \lambda_c L) \kappa_t^*}{1 - 2\rho \cos \lambda_c L + \rho^2 L^2} + \varepsilon_t, \quad t = 1, \dots, T \quad (2.32)$$

Sykliä voidaan käyttää mallissa monella tavalla. Ehkäpä kaksi keskeisintä ovat **trendi plus sykli**, jossa trendi ja sykli yhdistetään toisiinsa, ja **syklinen trendi**, jossa sykli vietään trendikomponentin sisään. Molemmissa tapauksissa syklikomponentin oletetaan olevan stationaarinen, eli ρ :n on ehdottomasti oltava pienempi kuin yksi.

Perusrakennemalli (2.12) - (2.15) saatetaan yhden yhtälön muotoon samalla tavoin kuin lineaarinen trendimallikin, joihin laajenuksena perusrakennemalli poikkeaa vain siinä, että mallissa on mukana myös γ_t , joka on joko trigonometrinen tai dummy-muuttujaan perustuva kausivaihtelukomponentti. Komponentin muokkaamisessa käytetään hyväksi summausoperaattoria $S(L)$. Summausoperaattorin avulla kausivaihtelu voidaan esittää muodossa

$$\gamma_t = \frac{\omega_t}{S_s(L)}, \quad (2.33)$$

joka sijoitetaan yhtälöön (2.10), saadaan perusrakennemalli muotoon

$$y_t = \frac{\eta_t}{\Delta} + \frac{\zeta_{t-1}}{\Delta^2} + \frac{\omega_t}{S_s(L)} + \varepsilon_t. \quad (2.34)$$

Autoregressiiviset mallit sisältävät autoregressiivisen komponentin, joka pystytään helposti yhdistämään rakenneyhtälömalliin. Autoregressiivinen prosessi voidaan liittää esimerkiksi satunnais/jäännöstermiin ε_t . Esimerkiksi ensimmäisen asteen autoregressiivinen komponentti paikallisen tasomallin satunnaiskomponenttiin liitettynä johtaa muotoon

$$y_t = \mu_t + v_t \quad (2.35a)$$

$$v_t = \phi_1 v_{t-1} + \varepsilon_t \quad (2.35b)$$

Yleisemmässä muodossa sama voidaan formuloida

$$y_t = \mu_t + \phi^{-1}(L) \varepsilon_t, \quad (2.36)$$

jossa

$$\phi(L) = 1 - \phi_1 L - \dots - \phi_p L^p. \quad (2.37)$$

Vaihtoehtoinen tapa autoregressiivisen mekanismin mukaantuomiselle on autoregressiivisen komponentin sijoittaminen suoraan y_t :hen. Tällöin ensimmäisen asteen autoregressiivisen komponentin tapauksessa (2.32):n sijasta saadaan

$$y_t = \phi_1 y_{t-1} + \mu_t + \varepsilon_t, \quad (2.38)$$

jonka esitys yleisessä muodossa on

$$y_t = \mu^\circ_t + \varphi^{-1}(L)\varepsilon_t, \quad (2.39a)$$

jossa

$$\mu^\circ_t = \varphi^{-1}(L)\mu_t. \quad (2.39b)$$

2.3.2 Aikasarjan stationarisointi

Rakenneyhtälömallin komponentit, kuten esimerkiksi trendi tai kausivaihtelu supistuvat deterministisiksi ajan funktioiksi, kun komponenttien varianssit ovat nolliä.

Stokastisuuden merkitys mallissa on keskeinen, koska se mahdollistaa sen, että komponentin arvo voi vaihdella jonkin kiinteän arvon ympärillä. Näiden stokastisten osien karakterisointi mahdollistaa mallin analysoinnin ja estimoinnin. Ensimmäinen askel stokastisten osien karakterisoinnissa on aikasarjan esittäminen stationaarisessa muodossa. Tavallisesti tämä tapahtuu differenssoinnin ja muiden sitä vastaavien toimien avulla. Stationarisoinnin jälkeen aikasarjaa voidaan analysoida mm. hyödyntämällä Kalmanin suodinta ja klassista testiteoriaa. Myös epästationaarisin tapauksiin soveltuvia ratkaisumalleja on kehitetty (de Jong & Chu-Chun-Lin, 1994).

Voidaan ajatella, että stokastisen prosessin tuottama aikasarja on yksi prosessin realisaatio. Prosessin stationaarisuus voidaan määritellä kahdella tavalla:

Prosessi on **heikosti stationaarinen** (*weakly stationary*), kun seuraavat ehdot ovat voimassa kaikille mahdollisille prosessin reaalisatioille ja kaikille havainnoille t .

$$E(y_t) = \mu, \quad (2.40)$$

$$E[(y_t - \mu)^2] = \text{Var}(y_t) = \gamma(0) \quad (2.41)$$

ja

$$E[(y_t - \mu)(y_{t-\tau} - \mu)] = \gamma(\tau), \quad \tau = 1, 2, \dots \quad (2.42)$$

Stationaarisen prosessin keskeinen piirre on se, että prosessin ominaisuudet eivät muutu ajan kuluessa. Toisin sanoen keskiarvo pysyy vakiona ja on riippumaton ajasta t . Myöskään prosessin varianssi tai autokovarianssi eivät saa riippua ajasta. Edellä mainitut ominaisuudet mahdollistavat sen, että keskiarvo, varianssi ja kovarianssit pystytään estimoimaan yksittäisen reaalisation perusteella:

$$\hat{\mu} = \bar{y} = T^{-1} \sum_{t=1}^T y_t, \quad (2.43)$$

$$\hat{\gamma}(0) = c(0) = T^{-1} \sum_{t=1}^T (y_t - \bar{y})^2 \quad (2.44)$$

ja

$$\hat{\gamma}(\tau) = c(\tau) = T^{-1} \sum_{t=\tau+1}^T (y_t - \bar{y})(y_{t-\tau} - \bar{y}). \quad (2.45)$$

Stokastisen stationaarisen prosessin dynaamiset ominaisuudet voidaan esittää autokovarianssi-funktion avulla, eli kuvaamalla $\gamma(\tau)$ kaikilla mahdollisilla τ :n positiivisilla arvoilla. Tarkastelua ei tarvitse laajentaa τ :n negatiivisiin arvoihin, koska $\gamma(\tau) = \gamma(-\tau)$. Kun autokovarianssit jaetaan prosessin varianssilla, saadaan muodostettua autokorrelaatiofunktio

$$\rho(\tau) = \frac{\gamma(\tau)}{\gamma(0)}, \tau = 1, 2, \dots \quad (2.46)$$

Autokorrelaatiofunktio saadaan kuvaamalla $\rho(\tau)$ τ :n suhteen. Otosautokorrelaatio määritellään vastaavasti

$$r(\tau) = \frac{c(\tau)}{c(0)}, \quad \tau = 1, 2, \dots \quad (2.47)$$

Stationaarisuus voidaan määritellä myös voimakkaamman määritelmän mukaan: Prosessi on **vahvasti stationaarinen** (*strictly stationary*), mikäli sen yhteystiheysfunktio ajasta riippumaton. Vahvasti stationaarinen prosessi on aina heikostikin stationaarinen, mutta päinvastoin väite ei ole aina tosi. Kuitenkin, mikäli heikosti stationaarinen prosessi on normaalisti jakautunut, on se stationaarinen myös vahvemman määritelmän mukaisesti.

Yhtälössä (2.31) esitettiin paikallinen lineaarinen trendimalli yhden yhtälön muodossa. Yhtälö saadaan stationaariseen muotoon kertomalla molemmat puolet kahteen kertaan differenssioperaattorilla Δ , eli

$$\Delta^2 y_t = \Delta \eta_t + \zeta_{t-1} + \Delta^2 \varepsilon_t. \quad (2.48)$$

Perusrakennemallin tapauksessa stationarisointi tapahtuu kertomalla molemmat puolet ensin differenssioperaattorilla Δ ja sitten vielä kausivaihtelun differenssioperaattorilla Δ_s . Tälle pätee $\Delta_s = \Delta S(L)$. Dummy-muuttujaisen kausivaihtelukomponentin sisältävälle perusrakennemallille (2.34) stationarisoitu muoto on

$$\Delta \Delta_s y_t = \Delta_s \eta_t + S(L) \zeta_{t-1} + \Delta^2 \omega_t + \Delta \Delta_s \varepsilon_t, \quad (2.49a)$$

ja taajuusalueeseen perustuvalla kausivaihtelukomponentilla varustetulle mallille stationarisoitu muoto saadaan myös summausoperaattoria (2.10) hyväksi käyttäen

$$S(L) \gamma_t = \left\{ \sum_{j=1}^{s/2-1} \bar{\gamma}_j(L) [(1 - \cos \lambda_j L) \omega_{jt} + (\sin \lambda_j L) \omega_{jt}^*] \right\} + \bar{\gamma}_{s/2}(L) \omega_{s/2,t}, \quad (2.49b)$$

jossa

$$\bar{\gamma}_j(L) = \frac{S(L)}{\gamma_j(L)} = \prod_{i \neq j} \gamma_i(L). \quad (2.49c)$$

Kaavan (2.49a) oikean puolen kolmas komponentti on korvattu Δ^2 :lla kerrotulla kaavan (2.49b) oikealla puolella. Kun aikasarja on stationaarissa muodossa, voidaan hyödyntää Kalmanin suodinta, jonka toimintaa tarkastellaan seuraavassa luvussa.

3. Tila-aika-malli ja Kalmanin suodin

Seuraavassa tarkastellaan aikasarjamallien **tila-aika-muotoa** (*state space form*, SSF), joka on keskeinen malli useiden erilaisten aikasarjamallien käsittelyssä. Kun malli on muunnettu tila-aika-muotoon, voidaan käyttää hyväksi Kalmanin suodinta, jonka kautta mahdollistuu myös mallin tuntemattomien parametrien estimointi suurimman uskottavuuden menetelmän avulla.

3.1 Tila-aika-malli

3.1.1 Mallin yleinen muoto

Tila-aika-muoto rakentuu kahdesta osasta, **mittausyhtälöstä** (*measurement equation*) ja **siirtymäyhtälöstä** (*transition equation*). Tila-aika-muodon yleinen esitysmuoto määritellään N :n aikasarjan tapauksessa, jossa siis y_t on $N \times 1$ vektori. Nyt voidaan kuitenkin rajoittaa tarkastelu yksinkertaiseen tilanteeseen, jossa $N = 1$. Tällöin mittausyhtälö voidaan esittää muodossa

$$y_t = z_t' \alpha_t + d_t + \varepsilon_t, \quad \text{Var}(\varepsilon_t) = h_t, \quad t = 1, \dots, T. \quad (3.1)$$

jossa z_t on vakioita tai parametreja sisältävä $1 \times m$ vektori, d_t on skalaariparametri ja ε_t on normaalisti jakautunut satunnaismuuttuja odotusarvolla 0 ja varianssilla h_t . α_t on satunnaismuuttujia sisältävä $m \times 1$ tilavektori, joka voidaan esittää siirtymäyhtälön avulla

$$\alpha_t = \mathbf{T}_t \alpha_{t-1} + \mathbf{c}_t + \mathbf{R}_t \eta_t, \quad t = 1, \dots, T, \quad (3.2)$$

jossa \mathbf{T}_t on $m \times m$ parametrimatriisi, \mathbf{c}_t $m \times 1$ vektori, \mathbf{R}_t $m \times g$ matriisi ja η_t $g \times 1$ jäännösvektori odotusarvolla 0 ja kovarianssimatriisilla \mathbf{Q}_t , eli

$$E(\eta_t) = 0 \quad \text{ja} \quad \text{Var}(\eta_t) = \mathbf{Q}_t. \quad (3.3)$$

Tila-aika-mallin määrittely on valmis, kun tehdään oletukset:

$$(a) E(\alpha_0) = \mathbf{a}_0 \quad \text{ja} \quad \text{Var}(\alpha_0) = \mathbf{P}_0, \quad (3.4)$$

(b) jäännösvektorit ε_t ja η_t ovat korreloimattomia, niin keskenään kuin alkutilan α_0 kanssa, eli

$$E(\varepsilon_t \eta_s') = 0 \quad \forall s, t = 1, \dots, T \quad (3.5)$$

ja

$$E(\varepsilon_t \alpha_0') = 0, \quad E(\varepsilon_t \alpha_0') = 0, \quad \text{kun } t = 1, \dots, T. \quad (3.6)$$

3.1.2 Mallin rakenteesta

Mittausyhtälön tekijöitä z_t , d_t ja h_t , samoin kuin siirtymäyhtälön tekijöitä T_t , c_t , R_t ja Q_t voidaan kutsua **systemikomponenteiksi**, ja mikäli muutoin ei määrätä oletetaan, että ne ovat deterministisiä. Systemikomponentitkin voivat toki vaihdella ajassa, mutta vaihtelu on ennalta määriteltyä. Edellä mainitun perusteella systemi on lineaarinen ja y_t voidaan kaikilla ajanhetkillä esittää sen hetkisten ja edeltävien jäännöstermien ε_t ja η_t sekä alkutilan α_0 avulla.

Jos systemikomponentit eivät muutu ajan kuluessa, mallin sanotaan olevan aikainvariantti tai aikahomogeeninen. Stationaariset mallit ovat erikoistapaus. Vaikka aikainvarianttien mallien luokka on paljon stationaaristen mallien luokkaa laajempi, on monilla epästationaarisilla aikainvariantteilla malleilla stationaarinen muoto, johon päästään muunnosten, kuten differenssoinnin avulla.

Tilavektorin α_t koostumus yksittäiselle tilastolliselle mallille riippuu mallin rakenteesta. Tilavektorin alkioille voi olla identifioitavissa oleva sisällöllinen tulkinta, esimerkiksi trendi tai kausivaihtelu, tai sitten ei. Tekniseltä kannalta katsottuna tilavektorin on tarkoitus sisältää kaikki olennainen informaatio aineistosta, mutta siten, että tilavektorin pituus on mahdollisimman pieni. Tila-aika-muotoa, joka minimoi tilavektorin pituuden sanotaan **minimoiduksi realisaatioksi** (*minimum realisation*). Eri esitysmuotojen ominaisuuksien ollessa muutoin tasavahvoja, minimoidun realisaation periaate on peruskriteeri esitysmuodon valinnalle. Periaate ei kuitenkaan merkitse sitä, että jokaisessa yksittäisessä tilanteessa löytyy yksikäsitteisesti paras muoto. Tilanteen, jossa yksikäsitteistä ratkaisua ei löydy, voidaan sanoa olevan pikemminkin sääntö kuin poikkeus (Harvey, 1989 s. 102).

Systemikomponentit z_t , h_t , T_t , R_t ja Q_t riippuvat tuntemattomista parametreista ja olennainen tavoite onkin näiden parametrien estimointi. Esimerkiksi perusrakennemallissa (2.11) myös mallin satunnaiskomponenttien varianssit σ_ε^2 , σ_η^2 ja σ_γ^2 ovat tuntemattomia. Nämä parametrit erotetaan $n \times 1$ vektoriin ψ . Näitä parametreja kutsutaan **hyperparametreiksi**. Hyperparametrit määrittelevät koko mallin stokastiset ominaisuudet, kun matriiseissa c_t ja d_t olevat parametrit vaikuttavat havaintoihin ainoastaan deterministisesti tila-aikamallin odotusarvon kautta. Jos c_t tai d_t on tuntemattomien parametrien lineaarinen funktio, näitä parametreja voidaan käsitellä tilamuuttujina.

Esimerkki 3.1 Oletetaan, että paikallisessa lineaarisessa trendimallissa (2.29b) μ_t on muotoa

$$y_t = \mu_t + \varepsilon_t = \frac{\eta_t}{\Delta} + \beta + \varepsilon_t, \quad (3.7)$$

eli trendin kulmakerroin komponentti β on deterministinen. Huolimatta tästä komponenttia voidaan käsitellä tilavektorin α_t alkiona määrittelemällä $\alpha_t = [\mu_t \ \beta_t]'$ ja kirjoittamalla tila-aika malli muodossa

$$y_t = [1 \ 0]\alpha_t + \varepsilon_t, \quad t = 1, \dots, T \quad (3.8)$$

$$\alpha_t = \begin{bmatrix} \mu_t \\ \beta_t \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mu_{t-1} \\ \beta_{t-1} \end{bmatrix} + \begin{bmatrix} \eta_t \\ 0 \end{bmatrix}. \quad (3.9)$$

3.2 Kalmanin suodin

3.2.1 Kalmanin suotimen tehtävä

Kun malli on saatettu tila-aika muotoon, voidaan malliin soveltaa erilaisia laskenta-algoritmeja, josta keskeisin on Kalmanin suodin (KS). KS on rekursiivinen proseduur, jolla pystytään tuottamaan tilavektorin optimaalinen estimaattori hetkellä t , kun käytössä on informaatio hetkelle t saakka. Systeemikomponentit sekä \mathbf{a}_0 ja \mathbf{P}_0 oletetaan tunnetuiksi.

Kalmanin suotimen keskeisenä roolina on toimia linkkinä klassiseen estimointi- ja testiteoriaan. Koska alkutilavektori ja tilavektorin alkiot ovat normaalijakautuneita, on mahdollista konstruoida uskottavuusfunktio ennustevirheen avulla. Näin päästään myös hyödyntämään klassista estimointi- ja testiteoriaa. Voidaan osoittaa, että normalisuusoletusten ollessa voimassa Kalmanin suodin tuottaa tilavektorille optimaalisen estimaattorin, joka on MMSE.

3.2.2 Kalmanin suotimen yleinen muoto

Tarkastellaan tila-aika-mallia (3.1) ja (3.2). Olkoon \mathbf{a}_{t-1} α_{t-1} :n optimaalinen estimaattori (MSE mielessä) perustuen havaintoihin $\mathbf{Y}_{t-1} = (y_1, y_2, \dots, y_{t-1})$. Olkoon \mathbf{P}_{t-1} estimointivirheen $\alpha_{t-1} - \mathbf{a}_{t-1}$ $m \times m$ kovarianssimatriisi, eli

$$\mathbf{P}_{t-1} = E[(\alpha_{t-1} - \mathbf{a}_{t-1})(\alpha_{t-1} - \mathbf{a}_{t-1})']. \quad (3.10)$$

Siirtymäyhtälöä käyttäen saadaan \mathbf{a}_{t-1} :n avulla optimaalinen estimaattori α_t :lle

$$\alpha_{t|t-1} = \mathbf{T}_t \mathbf{a}_{t-1} + \mathbf{c}_t. \quad (3.11a)$$

Estimointivirheen kovarianssimatriisi on puolestaan muotoa

$$\mathbf{P}_{t|t-1} = \mathbf{T}_t \mathbf{P}_{t-1} \mathbf{T}_t' + \mathbf{R}_t \mathbf{Q}_t \mathbf{R}_t', \quad t = 1, \dots, T \quad (3.11b)$$

Yhtälöitä (3.11a) ja (3.11b) kutsutaan **ennustamis-yhtälöiksi** (*prediction equations*).

Kun uusi havainto y_t saadaan aineistoon, α_t :n estimaattoria $\mathbf{a}_{t|t-1}$ voidaan päivittää. **Päivitysyhtälöt** (*updating equations*) ovat muotoa

$$\mathbf{a}_t = \mathbf{a}_{t|t-1} + \mathbf{P}_{t|t-1} \mathbf{z}_t' \mathbf{f}_t^{-1} (y_t - \mathbf{z}_t \mathbf{a}_{t|t-1} - \mathbf{d}_t), \quad (3.12a)$$

ja

$$\mathbf{P}_t = \mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1} \mathbf{z}_t' \mathbf{f}_t^{-1} \mathbf{z}_t \mathbf{P}_{t|t-1}, \quad (3.12b)$$

jossa \mathbf{f}_t on muotoa

$$\mathbf{f}_t = \mathbf{z}_t \mathbf{P}_{t|t-1} \mathbf{z}_t' + \mathbf{h}_t \quad (3.12c)$$

oleva skalaari. Yhtälöt (3.11) ja (3.12) yhdessä muodostavat Kalmanin suotimen. Yhtälöt voidaan yhdistää myös niin kutsutuksi **Riccati-yhtälöiksi** (*Riccati equations*):

$$\mathbf{a}_{t+1|t} = (\mathbf{T}_{t+1} - \mathbf{k}_t \mathbf{z}_t') \mathbf{a}_{t|t-1} + \mathbf{k}_t y_t + (\mathbf{c}_{t+1} - \mathbf{k}_t \mathbf{d}_t), \quad (3.13a)$$

jossa matriisi \mathbf{k}_t on muotoa

$$\mathbf{k}_t = \mathbf{T}_{t+1} \mathbf{P}_{t|t-1} \mathbf{z}_t' \mathbf{f}_t^{-1}, \quad t = 1, \dots, T. \quad (3.13b)$$

Kovarianssimatriisi \mathbf{P} on muotoa

$$\mathbf{P}_{t+1|t} = \mathbf{T}_{t+1} (\mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1} \mathbf{z}_t' \mathbf{f}_t^{-1} \mathbf{z}_t' \mathbf{P}_{t|t-1}) \mathbf{T}_{t+1}' + \mathbf{R}_{t+1} \mathbf{Q}_{t+1} \mathbf{R}_{t+1}', \quad t = 1, \dots, T \quad (3.13c)$$

jossa \mathbf{f} on määritelty kuten (3.12c). Eli Riccati-yhtälöt yhdistävät päivitys- ja ennustusyhtälöt.

Tilavektorin ja kovarianssimatriisin alkuarvoja voidaan merkitä \mathbf{a}_0 ja \mathbf{P}_0 . Alkuarvoista liikkeelle lähtemällä päivitetään tilavektoria havainto kerrallaan, kunnes kaikki havainnot on käyty lävitse. Näin edettynä tilavektorin estimaattori kerää mukaansa kaiken optimaalisten ennusteiden tekemiseen tarvittavan informaation. Laskennalliselta kannalta tarkasteltuna Kalmanin suodin ei ole aina välttämättä paras tapa edetä, vaan vaihtoehtoisen tavan Harveyn mukaan ovat esittäneet Andersson ja Moore (1979, luku 6). Tämän **informaatiosuotimen** (*information filter*) käyttö poikkeaa Kalmanin suotimesta vain siltä osin, että se käyttää laskennassa matriisin \mathbf{P}_t sijasta tämän käänteismatriisia \mathbf{P}_t^{-1} , joka myös informaatiomatriisin nimellä tunnetaan. Informaatiosuotimen edut eivät kuitenkaan tule juuri esille tässä työssä kuvattujen yhden aikasarjan perusrakennemallien kohdalla, vaan korostuvat tilanteissa, joissa analysoitavia sarjoja on useita samanaikaisesti ($\mathbf{y}_t = (y_{1t}, y_{2t}, \dots, y_{Nt})$) ja laskentaproseduurit vastaavasti raskaampia.

3.2.3 Kalmanin suotimen ominaisuuksista

Kalmanin suodin tuottaa siis α_t :lle ehdollisen odotusarvovektorin ja kovarianssimatriisin ehdollistettuna hetkellä t käytettävissä olevalla informaatiolla. Siis

$$\mathbf{a}_t = \mathbf{E}_t(\alpha_t) = \mathbf{E}(\alpha_t | \mathbf{y}_t), \quad (3.14a)$$

ja

$$\mathbf{P}_t = \mathbf{E}_t \{ [\alpha_t - \mathbf{E}_t(\alpha_t)] [\alpha_t - \mathbf{E}_t(\alpha_t)]' \}, \quad (3.14b)$$

jossa odotusarvo-operaattorin \mathbf{E} alla oleva t osoittaa odotusarvon olevan α_t :n ehdolliseen jakaumaan ajanhetkellä t perustuva odotusarvo. Ehdollinen odotusarvo on samalla α_t :n keskineliöpoikkeaman minimoiva estimaatti (MMSE).

Ehdollinen odotusarvo voidaan myös ymmärtää α_t :n estimaattoriksi. Estimaattorin ja estimaatin käsitteiden ero voidaan tulkita niin, että **estimaatti** on numeerinen arvo ja **estimaattori** puolestaan numeerisen arvon laskusääntö. Tässä yhteydessä tämä merkitsee sitä, että kun tarkastellaan estimaattia, ehdollinen odotusarvo on tiettyjen yksittäisten havaintojen numeerinen reaalisatio, kun puolestaan estimaattorin tapauksessa on kyse ilmauksesta, joka pätee mille tahansa havaintojoukolle. Näin ollen α_t :n ehdollinen odotusarvo onkin satunnaismuuttujien muodostama vektori. Kun ehdollinen odotusarvo otetaan perustuen kaikkiin käytettävissä oleviin havaintoihin, voidaan Harveyn mukaan osoittaa, että estimaattori minimoi keskineliöpoikkeaman (Andersson & Moore, 1979, s. 29 - 32).

Vaikka ehdollinen odotusarvo tässä tapauksessa onkin ymmärrettävissä satunnaisvektoriksi, ei silti ole sallittua puhua ehdollisen odotusarvovektorin kovarianssimatriisista. Sen sijaan, voidaan määrittellä harha estimointivirheen avulla ja näin voidaan myös määrittellä estimointivirheen kovarianssimatriisi, joka voidaan mieltää myös estimaattorin MSE-matriisiksi. On myös osoitettavissa, että ehdollisen odotusarvon estimointivirheen kovarianssimatriisi on samalla myös estimointivirheen yleinen kovarianssimatriisi ilman ehdollistamista. Näin ollen se voidaan jatkossa esittää ilman odotusarvo-operaattorin alla olevaa ehdollistamisen merkintää.

Edellä esitettyyn pohjautuen voidaan y_t :n ehdollinen odotusarvo ehdolla $t - 1$ esittää muodossa

$$\tilde{y}_{t|t-1} = \mathbf{z}_t \mathbf{a}_{t|t-1} + d_t. \quad (3.15)$$

Ennustevirheet voidaan määrittellä muodossa

$$v_t = y_t - y_{t|t-1} = \mathbf{z}_t (\alpha_t - \mathbf{a}_{t|t-1}) + \varepsilon_t, \quad t = 1, \dots, T. \quad (3.16)$$

Ennustevirhe edustaa siis viimeisen havainnon mukanaan tuomaa uutta informaatiota. Päivitysyhtälöissä (3.12a) - (3.12c) keskeisenä tehtävänä onkin tämän informaation lisääminen tilavektorissa aiempien havaintojen perusteella jo olevaan informaatioon. Ennustevirhe on riippumattomasti normaalijakautunut odotusarvolla 0 ja varianssilla f_t (3.12c), eli

$$v_t \sim \text{NID}(0, f_t). \quad (3.17)$$

Eri ajanhetkillä lasketut ennustevirheet oletetaan keskenään korreloimattomiksi

$$E(v_t v_s) = 0, \quad t \neq s \text{ ja } t, s = 1, \dots, T. \quad (3.18)$$

3.3 Aikainvarianttisuus

Useissa rakenneyhtälömallien sovelluksissa tila-aika-malli on aikainvariantti, toisin sanoen systeemikomponentit \mathbf{z}_t , d_t , h_t , \mathbf{T}_t , \mathbf{c}_t , \mathbf{R}_t ja \mathbf{Q}_t ovat kaikki ajasta riippumattomia, ja näin ollen ne voidaan esittää ilman aikaa kuvaavaa alaindeksiä t . Koska kaikki tässä yhteydessä kiinnostavat ominaisuudet sisältyvät malleihin, jotka sallivat matriisin \mathbf{c}_t ja skalaarin d_t vaihdella ajassa, joten tarkasteltavat mallit ovat muotoa

$$y_t = \mathbf{z}'_t \alpha_t + d_t + \varepsilon_t, \quad \text{Var}(\varepsilon_t) = h \quad (3.19a)$$

ja

$$\alpha_t = T\alpha_{t-1} + c_t + R\eta_t, \quad \text{Var}(\eta_t) = Q \quad (3.19b)$$

siten, että $E(\varepsilon_t \eta_s) = 0$ kaikille s ja t .

3.4 Rakenneyhtälömallin yhteys uskottavuusfunktioon

Klassinen suurimman uskottavuuden teoria perustuu tilanteeseen, jossa T kpl havaintoja ovat toisistaan riippumattomia ja samoin jakautuneita. Yhteystiheysfunktio on tällöin muotoa

$$L(\psi; y) = \prod_{t=1}^T f(y_t), \quad (3.20)$$

jossa $f(y_t)$ on t :nnen havainnon tiheysfunktio. Klassinen teoria ei kuitenkaan toimi aikasarjamallien kohdalla, koska perättäiset havainnot riippuvat toisistaan. Siksi (3.20):n sijasta käytetään ehdollista tiheysfunktioita

$$L(\psi; y) = \prod_{t=1}^T f(y_t | y_{t-1}), \quad (3.21)$$

jossa $f(y_t | y_{t-1})$ on havainnon y_t ehdollinen tiheysfunktio ehdolla $y_{t-1} = \{y_{t-1}, y_{t-2}, \dots, y_1\}$, eli hetkellä $t - 1$ käytettävissä olevan informaation perusteella. Jos tila-aika mallin (3.1) ja (3.2) jäännöstermit ja alkutilavektori α_0 ovat normaalisti jakautuneet, on y_t :n ehdollinen jakauma ehdolla Y_{t-1} myös normaali. Jakauman odotusarvo ja varianssi saadaan Kalmanin suotimen avulla. Kuten luvussa 3.2 johdettiin, odotusarvo on $a_{t|t-1}$ ja kovarianssimatriisi $P_{t|t-1}$. Mittausyhtälö (3.1) voidaan kirjoittaa muodossa

$$y_t = z_t a_{t|t-1} + z_t (\alpha_{t|t-1} - a_{t|t-1}) + d_t + \varepsilon_t, \quad (3.22)$$

joka on normaalisti jakautunut ja jonka ehdollinen odotusarvo on

$$E(y_t) = \tilde{y}_{t|t-1} = z_t a_{t|t-1} + d_t, \quad (3.23)$$

ja varianssi f_t on kuten (3.12c). Yhden muuttujan rakenneyhtälömalli tavallisesti uudelleenparametrisoidaan siten, että

$$\psi = [\psi_* \quad \sigma_*^2], \quad (3.24)$$

jossa ψ_* on vektori sisältäen $n - 1$ parametria ja σ_*^2 on mallin jonkin jäännöstermin varianssi. Jäännöstermien varianssit voidaan esittää muodossa

$$\text{Var}(\varepsilon_t) = \sigma_*^2 h_t \quad (3.25a)$$

ja

$$\text{Var}(\eta_t) = \sigma_*^2 \mathbf{Q}_t. \quad (3.25b)$$

Lisäksi edellä esitetyissä oletetaan, että h_t ja \mathbf{Q}_t riippuvat ψ_* :stä, mutta ei σ_*^2 :sta. Esittämällä varianssitermit σ_*^2 :n avulla termi saadaan ”vedettyä ulos” uskottavuusfunktioista jonka seurauksena uskottavuusfunktio tiivistyy. Näin ollen uskottavuusfunktio voidaan kirjoittaa **ennustevirrehajotelman** (*prediction error decomposition*) avulla muodossa

$$\log L(\psi; \mathbf{y}) = -\frac{T}{2} \log 2\pi - \frac{T}{2} \log \sigma_*^2 - \frac{1}{2} \sum_{t=1}^T \log f_t - \frac{1}{2\sigma_*^2} \sum_{t=1}^T \frac{v_t^2}{f_t}. \quad (3.26)$$

Koska y_t :n ehdollinen odotusarvo on samalla MMSE, voidaan v_t tulkita ennustevirhevektoriiksi,

$$v_t = y_t - \tilde{y}_{t|t-1}, \quad t = 1, \dots, T. \quad (3.27)$$

Derivoidaan (3.26) σ_*^2 :n suhteen, saadaan

$$\tilde{\sigma}_*^2(\psi_*) = \frac{1}{T} \sum_{t=1}^T \frac{v_t^2}{f_t}, \quad (3.28)$$

joka on σ_*^2 :n su-estimaattori ehdolla ψ_* . Sijoitetaan (3.27) uskottavuusfunktioon (3.26), saadaan **logaritmisoitu ja tiivistetty uskottavuusfunktio** (*concentrated log-likelihood function*)

$$\log L_c(\psi; \mathbf{y}) = -\frac{T}{2} (\log 2\pi - 1) - \frac{1}{2} \sum_{t=1}^T \log f_t - \frac{T}{2} \log \tilde{\sigma}_*^2(\psi_*), \quad (3.29)$$

joka maksimoidaan sitten vektorin ψ_* :n suhteen.

Jos alkutilan α_0 jakauma tunnetaan eksaktisti, toisin sanoen käytettävissä on täydellinen informaatio alkutilan komponenteista, Kalmanin suotimella saadaan johdettua ennustevirrehajotelman kautta täsmällinen uskottavuusfunktio havainnoille y . Täydellistä informaatiota ei kuitenkaan useinkaan ole käytettävissä. Tämän ongelma on johtanut kahdenlaisiin johtopäätöksiin: Kalmanin suodin on soveltuva tekniikka vain, 1) jos omaksutaan Bayesiläinen lähestymistapa, jossa α_0 :n jakauma on aina täsmällisesti määritelty. 2) jos otoskoko on niin suuri, että alkutilan täsmällisellä määrittelyllä ei ole merkitystä. Harveyn mukaan esitetyt johtopäätökset ovat perusteettomia, sillä yhden aikasarjan tapauksessa Kalmanin suodin saadaan käyttöön α_t :n ehdottoman jakaumaoletuksen perusteella, kun α_t on stationaarinen. Epästatio-naarisessa tapauksessa jakauma pystytään johtamaan tilavektorille α_d . Tällöin muodostetaan yhteistiheysfunktio havainnoille y_{d+1}, \dots, y_T ehdollistettuna havainnoilla y_1, \dots, y_d . Jos havainnot y_1, \dots, y_d oletetaan kiinteiksi, yhteistiheysfunktio on

$$\log L(\psi; \mathbf{y}) = -\frac{(T-d)}{2} (\log 2\pi - 1) - \frac{1}{2} \sum_{t=d+1}^T \log f_t - \frac{1}{2} \sum_{t=d+1}^T \frac{v_t^2}{f_t} \quad (3.30)$$

Tuntemattoman alkutilavektorin määrittämisen ongelmaa epästationaarisessa tapauksessa ovat käsitelleet mm. de Jong (1988), de Jong ja Chu-Chun-Lin (1994) ja Shephard (1993) ja heidän toimestaan on esitetty myös erilaisia menetelmiä alkutilan määrittämiseksi hetkelle t_0 .

Joskus saatetaan olla tilanteessa, jossa jokin α_0 :n alkioista onkin kiinteä. Tällaiset alkiot estimoidaan osana suurimman uskottavuuden estimointiproseduuria siten, että uskottavuusfunktiota tiivistetään niin, että kiinteät parametrit jäävät sen ulkopuolelle. Käytännössä estimointiproseduuria voidaan Harveyn mukaan suorittaa käyttäen **yleistettyä pns-menetelmää** (*generalised least squares*) tai Rosenbergin (1973) algoritmia.

Tähän saakka on oletettu, että havaintoaineisto on täydellinen, mutta mikäli aineistossa on puuttuvia havaintoja, joudutaan suorittamaan korjaustoimenpiteitä. Mikäli puuttuvien havaintojen lukumäärä on vähäinen, voidaan puuttuvat havainnot korvata nolllalla ja laajentaa mallia dummy-muuttujilla, siten että jokaista puuttunutta havaintoa kohden on oma dummy-muuttuja. Mikäli puuttuvia havaintoja on runsaasti ei dummy-muuttujiin perustuvaa puuttuvien havaintojen käsittelyä ole syytä tehdä, vaan tällöin tulisi soveltaa muita puuttuvan tiedon käsittelymenetelmiä. Vaihtoehtoinen tapa on esimerkiksi se, että Kalmanin suodinta käytettäessä ei päivitetäkkään tilavektoria α_t puuttuvan havainnon kohdalla, vaan mennään suoraan seuraavaan havaintoon.

4. Estimointi

Seuraavassa tarkastellaan ensin kuinka perusrakennemalli saatetaan tila-aika-muotoon ja kuinka tätä kautta mahdollistunut tuntemattomien parametrien estimointi tapahtuu.

4.1 Perusrakennemalli ja tila-aika-muoto

Aikainvariantti tila-aika-malli lineaariselle perusrakennemallille on muotoa

$$y_t = \mathbf{z}'\alpha_t + \varepsilon_t, \quad \text{Var}(\varepsilon_t) = h \quad (4.1a)$$

$$\alpha_t = \mathbf{T}\alpha_{t-1} + \mathbf{R}\eta_t, \quad \text{Var}(\eta_t) = \mathbf{Q} \quad (4.1b)$$

Kaikilla perusrakennemalleilla on olemassa aikainvariantti tila-aika-muoto. Sen sijaan malleissa, joissa mukaan on liitetty päivittäiskomponentti, matriisi \mathbf{z} on ajasta riippuva. Kun malli saadaan tila-aika-muotoon voidaan ottaa Kalmanin suodin käyttöön ja päästään tiivistetyn uskottavuusfunktion kautta klassisen estimointiteorian piiriin.

Additiivinen perusrakennemalli on muotoa

$$y_t = \mu_t + \gamma_t + \varepsilon_t, \quad t = 1, \dots, T, \quad (4.2)$$

jonka kausivaihtelukomponentti esitetään dummy-muuttujien avulla, eli

$$\sum_{j=0}^{s-1} \gamma_{t-j} = \omega_t, \quad \omega \sim \text{NID}(0, \sigma_\omega^2). \quad (4.3)$$

Olkoon $s = 4$, eli kyseessä esimerkiksi neljännesvuosiaineisto. Tällöin perusrakennemallin tila-aika-muoto on

$$y_t = [1 \ 0 \ 1 \ 0 \ 0]\alpha_t + \varepsilon_t \quad (4.4a)$$

$$\alpha_t = \begin{bmatrix} \mu_t \\ \beta_t \\ \dots \\ \gamma_t \\ \gamma_{t-1} \\ \gamma_{t-2} \end{bmatrix} = \begin{bmatrix} 1 & 1 & \vdots & & \mathbf{0} \\ 0 & 1 & \vdots & & \\ \dots & \dots & \dots & \dots & \dots \\ \vdots & & & -1 & -1 & -1 \\ \mathbf{0} & & \vdots & 1 & 0 & 0 \\ \vdots & & & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \mu_{t-1} \\ \beta_{t-1} \\ \dots \\ \gamma_{t-1} \\ \gamma_{t-2} \\ \gamma_{t-3} \end{bmatrix} + \begin{bmatrix} \eta_t \\ \beta_t \\ \dots \\ \omega_t \\ 0 \end{bmatrix} \quad (4.4b)$$

Tilavektorin kolmas komponentti edustaa kausivaihtelua hetkellä t , ja noudattaa yhtälössä (2.15) asetettua rajoitusta. Taajuusalueen avulla esitetyllä kausivaihtelukomponentilla (2.19) varustetun perusrakennemallin tila-aika-muoto on puolestaan

$$y_t = [1 \ 0 \ 1 \ 0 \ 1]\alpha_t + \varepsilon_t \quad (4.5a)$$

$$\alpha_t = \begin{bmatrix} \mu_t \\ \beta_t \\ \dots \\ \gamma_t \\ \gamma_t^* \\ \dots \\ \gamma_{2t} \end{bmatrix} = \begin{bmatrix} 1 & 1 & \vdots & \mathbf{0} & \dots & \mathbf{0} \\ 0 & 1 & \vdots & & \dots & \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \mathbf{0} & & \vdots & \mathbf{C}_1 & \vdots & \mathbf{0} \\ & & \vdots & & \vdots & \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \mathbf{0}' & & \vdots & \mathbf{0}' & \vdots & (-1)^t \end{bmatrix} \begin{bmatrix} \mu_{t-1} \\ \beta_{t-1} \\ \dots \\ \gamma_{1,t-1} \\ \gamma_{1,t-1}^* \\ \dots \\ \gamma_{2,t-1} \end{bmatrix} + \begin{bmatrix} \eta_t \\ \zeta_t \\ \dots \\ \omega_{1t} \\ \omega_{1t}^* \\ \dots \\ \omega_{2t} \end{bmatrix} \quad (4.5b)$$

jossa

$$\mathbf{C}_1 = \begin{bmatrix} \cos(\pi/2) & \sin(\pi/2) \\ -\sin(\pi/2) & \cos(\pi/2) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}. \quad (4.5c)$$

Tilavektorin kolmas ja neljäs termi vastaavat kausivaihtelusta hetkellä t . Kuukausiaineistolle siirtymämatriisiin \mathbf{T} rakenne säilyy samanlaisena, ainoastaan kausivaihteluun vaikuttavassa osassa on \mathbf{C}_1 :n sijaan \mathbf{C}_j , jossa $j = 1, \dots, 5$. Näin tilavektorin kausivaihtelutekijä hetkelle t ja siirtymämatriisissa olevat viisi kausivaihtelutekijää täyttävät yhdessä yhtälön (2.19) määrittelyyn sisältyneen ehdon $j = 1, \dots, [s/2]$, eli kuukausiaineistolle saadaan kuusi kausivaihtelutekijää.

4.2. Parametrien estimointi

4.2.1 Hyperparametrien estimointi

Oletetaan siis, että estimoitavat satunnaismuuttujat ε , η , ζ ja ω ovat normaalisti jakautuneita ja rajataan tarkastelu tilanteisiin, joissa tilavektorissa on enintään $d \leq m$ epästationaarista elementtiä. Varianssitermit h ja \mathbf{Q} voidaan määrittellä yhtälöiden (3.25a) ja (3.25b) osoittamalla tavalla, jossa σ_*^2 on positiivinen skalaari, joksi kiinnitetään joko σ_ε^2 tai σ_η^2 . Su-estimointi tehdään tiivistetylle uskottavuusfunktiolle maksimoimalla yhtälö ψ_* :n suhteen. ψ_* sisältää estimoitavat varianssiparametrit suhteutettuna kiinnitettyyn σ_*^2 :een. Perusrakennemallille alku-peräinen, täydellinen parametrivektori on

$$\psi = (\sigma_\varepsilon^2, \sigma_\eta^2, \sigma_\zeta^2, \sigma_\omega^2).$$

Valitaan σ_*^2 :ksi σ_η^2 , jolloin saadaan parametrivektori

$$\psi_* = (h, q_\zeta, q_\omega),$$

$$\text{jossa } h = \sigma_\varepsilon^2 / \sigma_\eta^2, \quad q_\zeta = \sigma_\zeta^2 / \sigma_\eta^2 \quad \text{ja} \quad q_\omega = \sigma_\omega^2 / \sigma_\eta^2.$$

Varianssitermejä voidaan estimoida autokovarianssi- ja autokorrelaatiofunktioita (2.46) - (2.48) hyväksikäyttäen. Yksinkertaiselle mallille kuten esimerkiksi paikallinen tasomalli (2.28) estimaattorit ovat helposti muodostettavia,

$$\hat{\sigma}_\eta^2 = c(0) + 2c(1) \quad (4.6)$$

ja

$$\hat{\sigma}_\varepsilon^2 = -c(1). \quad (4.7)$$

Monimutkaisemmille malleille estimaattorien konstruointi on monimutkaisempaa. Autokovarianssifunktioiden avulla muodostetut estimaattorit perustuvat momenttimenetelmään, jolla muodostetut estimaatit eivät ole asymptoottisesti tehokkaita.

4.2.2 Uskottavuusfunktion maksimointi

Ennustevirrehajotelman kautta konstruoitu uskottavuusfunktio maksimoidaan siis hyperparametrien suhteen. Maksimointi voidaan tehdä analyttisesti tai käyttäen erilaisia numeerisia optimointimenetelmiä, kuten esimerkiksi Gill-Murray-Pitfield- tai EM-algoritmia.

Kun alkutilavektorista ei ole saatavilla tarvittavaa informaatiota, muodostetaan ehdollistettu uskottavuusfunktio (3.21) havainnoille $y_{d+1}, y_{d+2}, \dots, y_T$, ehdollistettuna havainnoilla y_1, \dots, y_d , mikä voidaan esittää muodossa

$$\log L(\psi; \mathbf{y}) = -\frac{(T-d)}{2} \log 2\pi - \frac{1}{2} \sum_{t=d+1}^T \log f_t - \frac{1}{2} \sum_{t=d+1}^T \frac{v_t^2}{f_t}. \quad (4.8)$$

jossa d on tilavektorin epästationaaristen alkioden lukumäärä. Jos uskottavuusfunktiota tiivistetään, eli kiinnitetään σ_*^2 parametrivektorista ψ , uskottavuusfunktio voidaan maksimoida minimoimalla lauseke

$$S^\circ(\psi_*) = S(\psi_*) \left(\prod_{t=d+1}^T f_t \right)^{1/(T-d)}, \quad (4.9)$$

jossa

$$S(\psi_*) = \sum_{t=d+1}^T \frac{v_t^2}{f_t}. \quad (4.9b)$$

σ_*^2 :n su-estimaattori on silloin

$$\sigma_*^2(\psi_*) = \frac{S(\psi_*)}{(T-d)}. \quad (4.9c)$$

Tiivistämällä uskottavuusfunktiota kiinnittämällä parametri σ_*^2 pienenee optimoitavien parametrien lukumäärä ja laskennallinen proseduuri tehostuu samalla, kun tulosten luotettavuus paranee. Ainoa käytännön ongelma on σ_*^2 :n valinta. Jos valinta osuu termiin, joka on lähellä

nollaa, on seurauksena laskennallisia ongelmia, kun suhteelliset varianssit muodostuvat suuriksi.

Watsonin ja Englen (1983) mukaan voidaan myös EM-algoritmia käyttää tuntemattomien parametrien ratkaisemiseen. Oletetaan aikainvariantti malli (3.19), jossa σ_*^2 asetetaan ykköseksi. \mathbf{Q} oletetaan rajoittamattomaksi ja \mathbf{a}_0 ja \mathbf{P}_0 tunnetuiksi. Jos tilavektorin komponentit on havaittu ajanhetkillä $t = 0, \dots, T$, uskottavuusfunktio y_t :ille ja α_t :ille on

$$\begin{aligned} \log L(y, \alpha) = & -\frac{T}{2} \log 2\pi - \frac{T}{2} \log h - \frac{1}{2h} \sum_{t=1}^T (y_t - \mathbf{z}'\alpha_t)^2 \\ & - \frac{Tn}{2} \log 2\pi - \frac{T}{2} \log |\mathbf{Q}| - \frac{1}{2} \sum_{t=1}^T (\alpha_t - \mathbf{T}\alpha_{t-1})' \mathbf{Q}^{-1} (\alpha_t - \mathbf{T}\alpha_{t-1}) \\ & - \frac{n}{2} \log 2\pi - \frac{1}{2} \log |\mathbf{P}_0| - \frac{1}{2} (\alpha_0 - \mathbf{a}_0)' \mathbf{P}_0^{-1} (\alpha_0 - \mathbf{a}_0). \end{aligned} \quad (4.10)$$

Viimeiset kolme termiä (alin rivi) supistuu pois, mikäli alkutilavektorin jakauma ei ole selvillä ja $\mathbf{P}_0^{-1} = \mathbf{0}$. EM algoritmi toimii iteratiivisesti evaluoimalla odotusarvoa

$$E \left[\frac{\partial \log L}{\partial \psi} \middle| y_t \right]. \quad (4.11)$$

Yhtälö asetetaan yhtäsuureksi nollavektorin kanssa ja ratkaistaan uudet estimaatit ψ :lle. Tätä toistetaan kunnes prosessi konvergoi. Dempster et al. (1977), Wu (1983) ja Boyles (1983) ovat osoittaneet, että tietyin oletuksin iterointi konvergoi kohti funktion paikallista maksimia.

Sijoittamalla (4.11) funktioon (4.10) saadaan varianssin h estimaattoriksi:

$$\hat{h} = T^{-1} \sum_{t=1}^T \left[e_{qt}^2 + \mathbf{z}' \mathbf{P}_{qt} \mathbf{z} \right], \quad (4.12a)$$

jossa

$$e_{qt} = y_t - \mathbf{z}' \mathbf{a}_{qt}, \quad t = 1, \dots, T \quad (4.12b)$$

ja \mathbf{a}_{qt} on α_t :n **tasoitettu** (*smoothed*) estimaattori ja \mathbf{P}_{qt} sen MSE-matriisi. Vaikka EM-algoritmia joudutaan modifioimaan joidenkin rakenneyhtälömallien kohdalla, on sen käytöllä etuja puolellaan. Harveyn mukaan on osoitettavissa, että EM-algoritmin omien ominaisuuksien perusteella algoritmin tuottamat varianssiestimaatit täyttävät aina negatiivisuuteen liittyvät rajoitukset.

4.2.3 Determistiset komponentit

Mikäli malliin sisältyy deterministisiä komponentteja, ne voidaan poistaa tilavektorista ja siirtää mittausyhtälöön. Tilannetta voidaan pitää oikeastaan selittävien muuttujien mallittamisen erikoistapauksena. Selittäviä muuttujia tarkastellaan lähemmin luvussa 8. Malli voidaan tällöin kirjoittaa muodossa

$$y_t = \mathbf{z}' \alpha_t^* + \mathbf{x}_t' \delta^* + \varepsilon_t \quad (4.13)$$

$$\alpha_t^* = \mathbf{T}^* \alpha_{t-1}^* + \mathbf{R}^* \eta_t^*, \quad (4.14)$$

jossa α_t^* on tilavektori, pituudeltaan $m^* < m$, \mathbf{x}_t^* on $(m - m^*) \times 1$ vektori, joka sisältää aika-funktion. δ^* on $(m - m^*) \times 1$ vektori sisältäen alkuperäisestä tilavektorista pois siirretyt deterministiset parametrit. Parametrien estimointi voidaan tässä yhteydessä suorittaa käyttäen GLS-menetelmää.

4.3. Huomautuksia

Identifioitavuus on tilastollisen mallittamisen peruskysymyksiä. Identifioituvuutta tarkastellaan mallin ja sen rakenteen osalta. Mallissa määritellään kyseessä olevan muuttujan jakauma, kun taas rakenteessa määritellään jakauman parametrit. Jos kahdella eri rakenteella on sama yhteystiheysfunktio, sanotaan rakenteiden olevan ekvivalentit. Rakenne on identifioituva, jos tällaista ekvivalenttia rakennetta ei ole olemassa. Malli puolestaan on identifioituva, jos sen kaikki mahdolliset rakenteet ovat identifioituvat. On osoitettavissa, että kaikki tärkeimmät rakenneyhtälömallit ovat identifioituvia.

Pienten otosten (alle 200 havaintoa) tapauksessa **taajuusalueeseen** (*frequency-domain*) perustuvat mallit ovat laskennallisesti nopeampia kuin **aika-alueeseen** (*time-domain*) sijoittuvat mallit. Joskus nämä näkökulmat tuottavat toisistaan selvästi poikkeavia tuloksia. Erilaisiin tuloksiin on varsin usein kuitenkin syynä väärin spesifioitu malli. Kuitenkin on olemassa tilanteita, joissa taajuusalueen mallit ja aika-alueen mallit tuottavat oikein spesifioituinkin toisistaan selvästi poikkeavia tuloksia ja tällaisissa tapauksissa aika-alueeseen perustuvat mallit näyttävät Harveyn mukaan tuottavan luotettavampia tuloksia.

Tilanteissa, joissa havainnot eivät noudata normaalijakaumaa, estimaattorit lasketaan **kvasi-uskottavuuteen** (*quasi-maximum likelihood*) perustuen. Kvasi-uskottavuuteen perustuvien estimaattorien asymptoottiset ominaisuudet ovat johdettavissa, kunhan mallin asymptoottiset oletukset ovat voimassa (Harvey 1989, 210-212, 220).

5. Mallin spesifiointi, testaaminen ja diagnostiikka

5.1 Mallin spesifiointi

Aineistoon parhaiten sopivan **mallin valinta** on Harveyn (1989, s 11) mukaan aikasarja-analyysin vaikein aspekti. Mallin valinta perustuu diagnostisten testien hyväksikäyttöön. Näiden testien avulla pyritään varmistumaan, ei vain mallin jäännöksiä ominaisuuksien oikeellisuudesta, vaan myös siitä, että estimoidut komponentit ovat järkeviä tutkittavan ilmiön ominaisuuksiin nähden. Mallin diagnostiset tarkastelut voidaan jakaa kahteen osaan, toisaalta ollaan kiinnostuneita mallin vaatimien oletusten paikkaansapitävyydestä ja toisaalta eri mallien välisestä paremmuudesta.

Ekonometrisessa kirjallisuudessa (esimerkiksi Harvey, 1981) hyvälle mallille on esitetty seuraavat kriteerit:

(a) Mallin tulisi olla mahdollisimman **vähäparametrinen** (*parsimonious*), eli vertailtavien mallien ollessa muilta ominaisuuksiltaan yhtä hyvät, vähäparametrisuus on ratkaiseva kriteeri.

(b) Mallin tulisi olla **konsistentti** dataan verrattuna (*data coherence*). Toisin sanoen, mallin tulisi sopia hyvin aineistoon ja jäännösten tulisi olla suhteellisen pieniä ja normaalisti jakautuneita.

(c) Mallin tulee olla **ristiriidaton** tutkittavan ilmiön ominaisuuksien kanssa (*consistency with prior knowledge*).

(d) Mallin **ennustuskyvyn on rajoitettava** ilmiölle määriteltyjen **rajoitteitten määrittämille arvoille** (*data admissibility*).

(e) Mallin **tulisi sopia** hyvin paitsi otosaineistoon, myös tutkittavaan ilmiöön **otoksen ulkopuolella** (*structural stability*).

(f) Mallin tulisi olla käytettävissä olevaan informaatioon nähden **paras mahdollinen**, toisin sanoen mallin sisältämällä informaatiolla ei tulisi olla mahdollista muodostaa parempaa mallia (*encompassing*).

Rakenneyhtälömallien kohdalla mallin valinnan olennaisimmat kysymykset liittyvät mallin komponenttien luonteeseen, toisin sanoen ovatko komponentit stokastisia vai deterministisiä.

Ennen kuin niin pitkälle päästään, on analysoitavan sarjan käyttäytymistä tarkasteltava yleisemmällä tasolla. Keskeisin tekijä tässä on analysoitavan sarjan graafinen tarkastelu. Tarkastelu tehdään analysoitavan ilmiön tunnettuihin ominaisuuksiin perustuen. Usein on hyödyllistä tarkastella analysoitavaa sarjaa myös tehtyjen muuttujamuunnosten jälkeen. Graafisen tarkastelun avulla saadaan käsitys siitä, onko analysoitava sarja aikainvariantti, vai tapahtuuko tutkittavassa ilmiössä tarkasteluajankohtana rakenteellisia muutoksia. Myös mahdolliset virheeliset havainnot paikallistuvat helposti.

Pääsääntöisesti oikean mallin valinnan tulisi tapahtua siten, että lähdetään liikkeelle mallista, joka sisältää kaikki keskeiset komponentit. Näihin kuuluvat taso-, trendi- ja satunnaiskomponentti, sekä joissakin tapauksissa kausivaihtelukomponentti. Kausivaihtelukomponentin mukaanottaminen voi tapahtua esimerkiksi graafisten tarkastelujen tai sisältötietämyksen perusteella. Tämän jälkeen mallia pyritään supistamaan erilaisten diagnostisten testien perusteella, kunnes päästään optimaalinen malliin.

5.2 Testaaminen ja diagnostiikka

5.2.1 Testaamisen perusteet

Kaikki rakenneyhtälömallien testaamisessa käytetyt tilastolliset testit perustuvat klassisen testiteorian tarjoamiin perustestisiin. Klassinen testiteoria saadaan yhdistettyä rakenneyhtälömalleihin Kalmanin suotimen avulla johdetun uskottavuusfunktion kautta. Näin ollen normaali testausproseduuri on **uskottavuussuhteen testi** (*likelihood ratio test*). Uskottavuussuhteen testin (*us-testi*) ohella käytetään kahta muuta testiä, **Waldin testiä** (*Wald test*) ja **Lag-rangen testiä** (*Lagrange multiplier test*). Testien vaatimat oletukset ovat käytännössä samat kuin uskottavuusfunktion asymptoottiset ominaisuudet.

Uskottavuussuhteen testiä käytetään pääasiassa testattaessa $n \times 1$ parametrivektorin ψ rajoitteiden merkitsevyyttä. Uskottavuussuhde voidaan määritellä muodossa

$$\lambda = \frac{L(\tilde{\psi}_0)}{L(\tilde{\psi})}, \quad (5.1)$$

jossa osoittajassa on nollahypoteesia H_0 vastaava uskottavuusfunktio ja nimittäjässä vaihtoeh-toisen hypoteesin mukainen uskottavuusfunktio. Testisuure (US) voidaan kirjoittaa muodossa

$$US = -2 \log \lambda, \quad (5.2)$$

joka noudattaa asymptoottisesti $\chi^2(m)$ -jakaumaa nollahypoteesin ollessa voimassa. Haittapuo-lena uskottavuussuhteen testissä on se, että uskottavuusfunktio joudutaan estimoimaan, sekä nolla-, että vaihtoeh-tohypoteesin tilanteessa.

Waldin testin etuna on, että uskottavuusfunktiota ei tarvitse estimoida kuin vaihtoeh-toisen hypoteesin mukaisessa tilanteessa. Testisuureen yleinen muoto on

$$\mathbf{W} = [\mathbf{R}\tilde{\psi} - \mathbf{r}]' [\mathbf{R}\mathbf{I}^{-1}(\tilde{\psi})\mathbf{R}']^{-1} [\mathbf{R}\tilde{\psi} - \mathbf{r}], \quad (5.3)$$

jossa $\mathbf{I}(\tilde{\psi})$ on vaihtoeh-toisen hypoteesin mukainen informaatiomatriisi. Kuten uskottavuussuh-teen testikin, Waldin testi noudattaa asymptoottisesti $\chi^2(m)$ -jakaumaa nollahypoteesin ollessa voimassa. Tilanteissa, joissa uskottavuusfunktio on helpointa estimoida nollahypoteesin mu-kaisessa tilanteessa, on käytännöllisintä soveltaa Lagrangen testiä, joka noudattaa nollahypo-teesin voimassaollessa niin ikään $\chi^2(m)$ -jakaumaa. Testisuure voidaan esittää muodossa

$$Lg = \left[\frac{\partial \log L}{\partial \psi} \right]' \mathbf{I}^{-1}(\tilde{\psi}_0) \left[\frac{\partial \log L}{\partial \psi} \right], \quad (5.4)$$

jossa oikealla puolella olevat osamäärätermit on ratkaistu nollahypoteesin mukaisessa tilanteessa. Testattaessa aika-alueessa mallia, jossa useampi jäännöstermi, Lagrangen testi voi osoittautua työlääksi. Taajuusalueessa vastaavaa ongelmaa ei esiinny.

Rakenneyhtälömalleihin liittyvissä testaustilanteissa törmätään joskus tilanteisiin, joissa US-testin vaatimat oletukset eivät pidä. Tällöin voidaan soveltaa **voimakkaimman invariantin testiä** (*most-powerful-invariant test*). Testin konstruointia ovat tarkastelleet esimerkiksi Franzini ja Harvey (1983).

Tilanteissa, joissa joku tai jotkut parametreista ovat määritellyn parametriavaruuden reunalla, ollaan yksisuuntaisessa testaustilanteessa, jolloin uskottavuussuhteen- ja Waldin testit eivät noudata $\chi^2(m)$ -jakaumaa. Lagrangen testi noudattaa kuitenkin alkuperäistä jakaumaansa. Uskottavuussuhteen- ja Waldin testin jakaumat noudattavat χ^2 -jakaumien sekoitusta, joka on sama kuin Kuhn-Tucker-testisuureen noudattama. Näin ollen jakauma on johdettavissa ja US- ja Waldin testit ovat käytettävissä myös parametriavaruuden reunaparametreille. Reunaparametrien testauksessa on kuitenkin huomioitava se, että testaustilanne on yksisuuntainen. Lagrangen testi ei kykene huomioimaan tätä yksisuuntaisuutta ja näin testin voimakkuus jää US- ja Waldin testejä heikommaksi, ellei testiä modifioida. Harveyn mukaan simuloimalla on voitu osoittaa, että modifioitu yksisuuntainen Lagrangen testi on voimakkuudeltaan samaa luokkaa kuin uskottavuussuhteen testi (Harvey ja Hotta 1982).

5.2 Mallin diagnostiset testit

Mallin spesifiointi jakautuu toisaalta mallin vaatimien oletusten tarkasteluun ja toisaalta mallin riittävyuden arviointiin. Johtopäätökset oletusten voimassaolosta ja mallien hyvyydestä tehdään tilastollisten testien antaman informaation perusteella.

Keskeinen kysymys rakenneyhtälömallien spesifiointissa on määrittää ovatko mallin komponentit deterministisiä vai stokastisia, eli toisin sanoen ovatko komponenttien satunnaistekijöiden varianssit nolli. Perusrakennemalli sisältää neljä estimoitavaa parametria, σ_ε^2 , σ_η^2 , σ_ζ^2 ja σ_ω^2 riippumatta siitä, onko mallin kausivaihtelukomponentti esitetty aika- tai taajuusalueen avulla. Kaikki neljä oletetaan positiivisiksi, koska ne ovat variansseja. Mikäli termeistä joku on nolla, se sijoittuu parametriavaruuden reunalle, jolloin ollaan luvun 5.2.1 lopussa kuvatussa tilanteessa. US-testi saadaan tällöinkin modifioitua, mikäli satunnaismuuttuja noudattaa asymp-toottisesti normaalijakaumaa. US-, Waldin- tai Lagrangen testillä voidaan testata, onko σ_ε^2 tai σ_η^2 nolla. US- ja Waldin testi lasketaan tavanomaiseen tapaan ja jakauma johdetaan Kuhn-Tucker-testisuureen noudattaman jakauman perusteella. Yksipuolinen Lagrangen testi joudutaan modifioimaan erikseen.

Mallin vaatimista oletuksista ensimmäinen on oletus aineiston normaalijakautuneisuudesta. Oletus toteutuu, kun mallin kaikki satunnaistekijät ovat normaalisti jakautuneita. Aineiston normalisuuden ohella keskeinen oletus on jäännösten normaalijakautuneisuus siten, että

jäännösten odotusarvo on nolla ja varianssi vakio. Rakenneyhtälömallien osalta jäännöksiksi voidaan tulkita mallin ennustevirhe

$$\hat{v}_t = y_t - \hat{y}_{t|t-1}, \quad t = d + 1, \dots, T \quad (5.5)$$

Jäännösten normaalijakautuneisuutta voidaan tutkia mittaamalla aineiston huipukkuutta ja vinoutta. Olkoon μ y :n keskiarvo ja σ^2 sen varianssi ja merkitään

$$\mu_i = E[y - \mu]^i,$$

eli $\sigma^2 = \mu_2$. Nyt vinous (5.6a) ja huipukkuus (5.6b) voidaan määritellä muodossa

$$\sqrt{\beta_1} = \frac{\mu_3}{\mu_2^{3/2}}, \quad (5.6a)$$

ja

$$\beta_2 = \frac{\mu_4}{\mu_2^2}. \quad (5.6b)$$

Jäännöksille voidaan määritellä vastaavat termit:

$$\bar{v} = \frac{1}{T} \sum_{t=1}^T \tilde{v}_t, \quad (5.7a)$$

$$m_i = \frac{1}{T} \sum_{t=1}^T (\tilde{v}_t - \bar{v})^i, \quad (5.7b)$$

$$\sqrt{b_1} = \frac{m_3}{m_2^{3/2}} \quad (5.7c)$$

ja

$$b_2 = \frac{m_4}{m_2^2}. \quad (5.7d)$$

Normaalijakautuneelle muuttujalle pätee $\sqrt{\beta_1} = 0$ ja $\beta_2 = 3$. Normaalisuutta mittava testi on Bowmanin ja Shentonin (1975) modifioima normaalisuustesti, jonka testisuure on muotoa

$$N_{BS} = s + k, \quad (5.8a)$$

jossa

$$s = \frac{T(\sqrt{b_1})^2}{6} \quad (5.8b)$$

ja

$$k = \frac{T(b_2 - 3)^2}{24}. \quad (5.8c)$$

Testisuure N_{BS} noudattaa χ^2 -jakaumaa vapausastein 2.

Bowmanin ja Shentonin testi soveltuu vain suurille otoksille. Pienille otoksille ($T < 250$) paremmin soveltuvan testin (N_{DH}) ovat Harveyn mukaan esittäneet Doornik ja Hansen (1994). Mikäli normalisuusoletus ei päde, voidaan aineistoa yrittää normalisoida erilaisin muunnoksilla. Mikäli muunnoksinkaan ei päästä normalisuuteen, joudutaan mallin estimoinnissa jatkossa hyödyntämään kvasi-uskottavuutta (luku 4.3).

Autokorrelaatiofunktio ja siihen liittyvät graafiset esitykset, lähinnä korrelogrammi, voivat myös olla käyttökelpoisia informaatiolähteitä, vaikka korrelogrammin ja autokorrelaatiofunktion merkitys ei olekaan yhtä merkittävässä roolissa kuin ARIMA-mallien kohdalla.

Jäännösvarianssin vakioisuutta testataan varianssitestillä, joka on muotoa

$$H(h) = \frac{\sum_{t=T-h+1}^T \tilde{v}_t^2}{\sum_{t=d+1}^T \tilde{v}_t^2}. \quad (5.9)$$

Testataan siis h :n ensimmäisen ja h :n viimeisen jäännöksen suhdetta. Testiin mukaan tulevat jäännökset määräävä kokonaisluku h määritellään siten, että $h \approx T^*/3$. Yleisessä muodossa testisuure noudattaa $F(h,h)$ jakaumaa. Vaihtoehtona on testisuure $hH(h)$, joka noudattaa asymptoottista $\chi^2(h)$ -jakaumaa nollahypoteesin voimassaollessa.

Normaalisuuden ja homoskedastisuuden lisäksi jäännökset oletetaan myös keskenään korreloimattomiksi. Jäännösten autokorreloituneisuuden testaamiseen on useita keinoja, riippuen siitä, tarkastellaanko yksittäistä autokorrelaatiofunktioita vai autokorrelaatiofunktioita kokonaisuutena.

Perättäisten jäännösten autokorrelaatio viiveellä k on muotoa

$$r_v(\tau) = \frac{\sum_{t=d+1+\nu}^T (\tilde{v}_t - \bar{\tilde{v}})(\tilde{v}_{t-\nu} - \bar{\tilde{v}})}{\sum_{t=d+1}^T (\tilde{v}_t - \bar{\tilde{v}})^2}, \quad (5.10)$$

joka normalisuusoletuksen, korreloimattomuuden ja homoskedastisuuden ollessa voimassa noudattaa asymptoottisesti jakaumaa $N(0, 1/(T^* - \nu))$.

Tarkasteltaessa jäännösten yleistä autokorreloituneisuutta, voidaan käyttää Box ja Ljungin testiä (Ljung ja Box, 1978), joka on P :lle ensimmäiselle autokorrelaatiolle muotoa

$$Q^* = T^* (T^* + 2) \sum_{\tau=1}^P (T^* - \tau)^{-1} r_v^2(\tau), \quad (5.11)$$

jossa $T^* = T - d$ ja $r_v(\tau)$ kuten (5.10). Q^* noudattaa asympotoottisesti $\chi^2(P - n^*)$ -jakaumaa, jossa n^* on estimoitavien hyperparametrien lukumäärä. Boxin ja Ljungin mukaan testi on sangen robusti jäännösten jakauman suhteen.

Ensimmäisen asteen autokorreloituneisuutta voidaan testata Durbinin ja Watsonin konstruoidulla testillä. Testisuure on muotoa

$$DW = \frac{\sum_{t=d+2}^T (\tilde{v}_t - \tilde{v}_{t-1})^2}{\sum_{t=d+1}^T \tilde{v}_t^2}, \quad DW \sim N(2, 4/T^*). \quad (5.12)$$

Hyödyllistä tietoa jäännöksistä saadaan myös käyttäen hyväksi jäännösten kumulatiivisen summan (*CUSUM*) tai kumulatiivisen neliöiden summan (*CUSUMSQ*) testiä ja kumulatiivisten summamuuttujien graafisia esityksiä. Kumulatiivisen summan testisuure on

$$CUSUM(t) = \hat{\sigma}_*^{-1} \sum_{j=d+1}^t \tilde{v}_j, \quad t = d + 1, \dots, T, \quad (5.13a)$$

jossa

$$\hat{\sigma}_*^2 = (T - d - 1)^{-1} \sum_{t=d+1}^T (\tilde{v}_t - \bar{\tilde{v}})^2 \quad (5.13b)$$

on jäännösvarianssi. Suurilla otoksilla $\hat{\sigma}_*^2$ voidaan korvata estimaattorilla $\tilde{\sigma}_*^2$, joka on muotoa.

$$\tilde{\sigma}_*^2 = (T - d)^{-1} \sum_{t=d+1}^T \tilde{v}_t^2. \quad (5.13c)$$

Käytännössä testin avulla muodostetaan graafi, jota verrataan kahteen ennalta määrättyyn merkitsevyyskäyrään, jotka saadaan yhtälöstä

$$CUSUM = \pm \left[a\sqrt{T-d} + \frac{2a(t-d)}{\sqrt{T-d}} \right], \quad (5.14)$$

jossa $a = 0.948$ 5%:n ja 0.850 10%:n merkitsevyystasolle. *CUSUM*-proseduuri on erittäin käyttökelpoinen aikasarjan rakenteellisten muutosten havainnoinnissa.

Kumulatiivisen neliöiden summan testisuure on muotoa

$$\text{CUSUMSQ}(t) = \frac{\sum_{j=d+1}^t \tilde{v}_j^2}{\sum_{t=d+1}^T \tilde{v}_t^2}. \quad (5.15)$$

Kumulatiivisen neliöiden summatestin avulla voidaan myös jäljittää aikasarjassa tapahtuvia rakenteellisia muutoksia, mutta myös testata jäännösten heteroskedastisuutta. Testisuureen arvoista ajan suhteen piirretyn graafin tulisi nollahypoteesin (homoskedastisuus, ei rakenteellisia muutoksia) voimassa ollessa muodostaa 45° kulmassa olevan suoran.

Voidaan ajatella, että mallin valintaprosessin alkaessa on kaikki mahdolliset rakenneyhtälömallin vaihtoehdot valittavissa, riippumatta siitä ovatko ne realistisia tai oletusten mukaisia. Kun ollaan varmistuttu mallin vaatimien oletusten voimassaolosta valittavana olevien mallivaihtoehtoista supistuvat kaikki perusominaisuuksiltaan epäkelvot mallit pois ja jäljelle jää teoreettisten oletusten osalta kelvolliset rakenneyhtälömallit. Mallivaihtoehtojen lukumäärää saadaan supistettua myös selvittämällä komponenttien luonne stokastisuuden ja deterministisuuden osalta. Jäljelle jääneiden mallivaihtoehtojen lukumäärää pyritään supistamaan edelleen riittävyystarkastelujen ja sisältötietämyksen avulla, kunnes jäljelle jää vain yksi malli, joka sitten valitaan.

Lopullinen mallivalinta teoreettisten oletusten suhteen kelvollisten mallien kohdalla tehdään yleisiä mallin valinnan periaatteita (luku 5.1) noudattaen. Mallin sopivuutta aineistoon voidaan mitata **ennustevirheen varianssilla** (*prediction error variance, p.e.v.*), jonka estimaattori on aika-alueessa operoitaessa muotoa

$$\tilde{\sigma}_v^2 = \frac{1}{T-d} \sum_{t=d+1}^T \tilde{v}_t^2. \quad (5.16)$$

Determinaatiokerrointa (*coefficient of determination, R²*) eli selitystasetta käytetään kuvaamaan kuinka suuren osan kokonaisvaihtelusta malli kykenee selittämään. Determinaatiokerroin määritellään eri tavoin riippuen siitä, millainen malli on kyseessä. Paikallisen tasomallin (2.26) determinaatiokerroin saadaan kaavalla

$$R_L^2 = 1 - \frac{(T-d)\hat{\sigma}_v^2}{\sum_{t=1}^T (y_t - \bar{y})^2}. \quad (5.17)$$

Trendimallin (2.28) kerroin pystytään puolestaan ratkaisemaan kaavalla

$$R_D^2 = 1 - \frac{(T-d)\hat{\sigma}_v^2}{\sum_{t=d+1}^T (w_t - \bar{w})^2}, \quad (5.18)$$

jossa w_t on y_t :n differenssi. Kausivaihtelukomponentin sisältävälle mallille kerroin on muotoa

$$R_S^2 = 1 - \frac{(T-d)\hat{\sigma}_v^2}{\sum_{t=d+1}^T (w_t - \bar{w}_{ts})^2}, \quad (5.19)$$

jossa \bar{w}_{ts} on kausivaihtelukomponentin s sisältävien differenssien w_t keskiarvo.

Samalle havaintoaineistolle tehtyjä malleja voidaan vertailla, paitsi ennustevirheen varianssin, myös informaatiokriteerien avulla. **Akaiken informaatiokriteeri** (*Akaike information criterium*, *AIC*) on muotoa

$$AIC = \tilde{\sigma}_v^2 \exp\left[\frac{2(n+d)}{T}\right], \quad (5.20)$$

jossa d on mallin epästationaaristen komponenttien ja n hyperparametrien lukumäärä. **Bayesin informaatiokriteeri** (*Bayesian information criterium*, *BIC*) on AIC:stä muodoltaan vain hi-
venen poikkeava:

$$BIC = \tilde{\sigma}_v^2 \exp\left[\frac{(n+d) \log T}{T}\right]. \quad (5.21)$$

Koska parametrien määrää lisäämällä ennustevirheen varianssia on mahdollista pienentää, on usein tuloksena liian runsasparametrisia malleja. Informaatiokriteerit rakentuvat vähäparametrisuuden periaatteelle (luku 5.1.). Kaavoissa (5.20) ja (5.21) nähtävä summatekijä $(n+d)$ toimii parametrien määrää kontrolloivana tekijänä. Mikäli uuden parametrin lisäyksestä saata-
va ennustevirheen varianssin pieneneminen ei ole tarpeeksi suurta, informaatiokriteerin arvo kasvaa ja päädytään informaatiokriteerin mukaan huonompaan malliin kuin vähäparametrisemmassa tilanteessa.

Oikean mallin valinnassa voidaan soveltaa myös ennustamismenetelmiä, joita tarkastellaan lähemmin luvussa 8. Ennustamista voidaan hyödyntää mallin valinnassa kahdella tavalla.

(1) Jaetaan havaintojakso kahteen osaan: y_1, \dots, y_{T-l} ja y_{T-l+1}, \dots, y_T . Estimoidaan malli havaintoihin y_1, \dots, y_{T-l} perustuen ja tuotetaan ennusteet jälkimmäiselle havaintojaksolle y_{T-l+1}, \dots, y_T . Saadaan siis aitoja ennusteita, joita voidaan sitten verrata todellisiin havaintoihin ja arvioida näin mallin ennustuskykyä.

(2) Laaditaan havaintoja y_{T-l+1}, \dots, y_T vastaavat ennusteet koko havaintoaineisto y_1, \dots, y_T hyödyntäen eli ennustetaan hetkeltä T taaksepäin. Tulokseksi saatavat ennusteet eivät ole aitoja, sillä ne sisältävät todellisten havaintojen sisältämän informaation myös ennustamisperiodin ajalta. Saaduilla epäaidoilla ennusteilla voidaan kuitenkin arvioida mallin ennustuskykyä.

Mallin riittävyttä ennustekyvyn valossa voidaan testata **Chow'n testillä**, joka on muotoa

$$\text{Chow} = \frac{T - l - d^*}{l} \frac{\sum_{t=T-l+1}^T v_t^2}{\sum_{t=d^*+1}^{T-1} v_t^2}, \quad (5.22)$$

jossa v_t :llä merkitään ennustevirhettä, kun $t = d^* + 1, \dots, T$ ja $d^* = d + k$, jossa k on mahdollisten selittävien muuttujien lukumäärä. Chow'n testisuure noudattaa $F(l, T - l - d^*)$ -jakaumaa. Merkitään standartoituja jäännöksiä \tilde{v} :llä, eli

$$\tilde{v}_{T+j} = \frac{v_{T+j}}{f_{T+j}^{1/2}}. \quad (5.23)$$

Mallin ennustekykyä voidaan arvioida **ennustamisen epäonnistumistestin** (*post-sample predictive failure test statistic*) avulla. Testisuure on muotoa

$$\text{pft} = \sum_{j=1}^L \tilde{v}_{T+j}^2, \quad (5.24)$$

joka noudattaa approksimatiivisesti χ^2 -jakaumaa vapausastein L . Standardoitujen jäännösten osalta voidaan soveltaa myös **Cusum t-testiä**, jonka testisuure on likimäärin t -jakautunut vapausastein $T - L - d^*$ ja muotoa

$$\text{cusum } t = L^{-1/2} \sum_{j=1}^L \tilde{v}_{T+j}. \quad (5.25)$$

Standardoituja jäännöksiä (5.23) hyödyntäen voidaan myös laskea havaintoperiodin jälkeisen ennustetestin testisuure

$$\xi(l) = \sum_{j=1}^l \tilde{v}_{T+j}^2 / l \Big/ \sum_{t=d+1}^T \tilde{v}_t^2 / (T - d), \quad (5.26)$$

joka noudattaa $F(l, T-d)$ -jakaumaa. Väärin spesifioidun mallin tapauksessa testisuure saa tilastollisesti merkitsevän arvon, joka johtaa mallin hylkäämiseen. Tilastollisesti ei-merkitsevä arvo sen sijaan ei automaattisesti kerro mallin olevan hyvä.

Kun riittävän hyvä malli on saatu konstruoitua, voidaan mallia käyttää esimerkiksi ennustamiseen, jota tarkastellaan seuraavassa.

6. Ennustaminen

Tilastollinen malli perustuu tiettyihin tilastollisiin ominaisuuksiin. Näiden ominaisuuksien myötä aukeaa useita mahdollisuuksia, esimerkiksi mahdollisuus konstruoida laadituille ennusteille ennustevälit (luottamusvälit) tilastollisen mallin pohjalta. Seuraavaksi tarkastellaankin ennustefunktiota, ennustamista ja ennustevälien määrittämistä.

6.1 Perusrakennemallilla ennustaminen

Kun rakenneyhtälömalli muokataan tila-aika-muotoon saadaan Kalmanin suotimen avulla estimoitua tilavektorin arvo hetkellä T . Tämä estimaatti on samalla sarjan komponenttien ennustefunktioiden alkuarvo. **Yhden askeleen ennuste** (*one-step-ahead prediction*) aikainvariantille mallille (3.19) voidaan kirjoittaa tila-aika-mallia soveltaen muodossa

$$\mathbf{a}_{T+1|T} = \mathbf{T}\mathbf{a}_T + \mathbf{c}_{T+1}, \quad (6.1)$$

$$\tilde{\mathbf{y}}_{T+1|T} = \mathbf{z}'\mathbf{a}_{T+1|T} + d_{T+1}. \quad (6.2)$$

Yhtälöt (6.1) ja (6.2) yhdistämällä saadaan yksi ennustefunktio muodossa

$$\tilde{\mathbf{y}}_{T+1|T} = \mathbf{z}'\mathbf{T}\mathbf{a}_T + \mathbf{z}'\mathbf{c}_{T+1} + d_{T+1}. \quad (6.3)$$

Seuraavalle ajanhetkelle saadaan ennuste siirtymäyhtälön

$$\alpha_{T+2} = \mathbf{T}^2\alpha_{T+1} + \mathbf{c}_{T+2} + \mathbf{R}\eta \quad (6.4)$$

avulla ottamalla ehdollinen odotusarvo ehdolla T . Saadaan tilavektorin ennuste hetkelle $T + 2$:

$$\mathbf{a}_{T+2|T} = \mathbb{E}_T(\alpha_{T+2|T}) = \mathbf{T}^2\mathbf{a}_{T+1|T} + \mathbf{c}_{T+2}. \quad (6.5)$$

Yhtälö (6.5) voidaan yleistää tilavektorin **usean askeleen ennustefunktioksi** (*multi-step-ahead prediction*), jolloin yhtälö kirjoitetaan muodossa

$$\mathbf{a}_{T+l|T} = \mathbb{E}_T(\alpha_{T+l}) = \mathbf{T}^l\mathbf{a}_{T+l-l|T} + \mathbf{c}_{T+l}. \quad (6.6)$$

Tällöin ennuste y_{T+l} :lle on muotoa

$$\tilde{\mathbf{y}}_{T+l|T} = \mathbf{z}'\mathbf{T}^l\mathbf{a}_{T+l-l} + \mathbf{z}'\mathbf{c}_{T+l} + d_{T+l}. \quad (6.7)$$

Ennuste on y_{T+l} :n **MMSLE** (*minimum mean square linear estimator*).

Tilavektorin α_{T+l} ehdollinen jakauma on normaali ja estimointivirheen ($\alpha_{T+l} - \mathbf{a}_{T+l}$) kovarianssimatriisi $\mathbf{P}_{T+l|T}$ voidaan esittää muodossa

$$\mathbf{P}_{T+l|T} = \mathbf{T}'\mathbf{P}_T\mathbf{T}' + \sum_{j=0}^{l-1} \mathbf{T}^j\mathbf{R}\mathbf{Q}\mathbf{R}'\mathbf{T}^{j'}, \quad l = 1, 2, \dots \quad (6.8)$$

MSE-matriisi aliestimoi todellista keskineliöpoikkeamaa, mikäli mallissa on estimoitavia tuntemattomia parametreja (ψ), sillä MSE ei ota huomioon estimoitavien parametrien sisältämää vaihtelua. Ennusteen MSE voidaan esittää ehdollistettuna estimoitavien hyperparametrien suhteen, tällöin MSE on muotoa

$$\text{MSE}(\tilde{y}_{T+l|T}) = \left[\mathbf{z}'\mathbf{T}'\mathbf{P}_T\mathbf{T}'\mathbf{z} + \mathbf{z}'\left(\sum_{j=0}^{l-1} \mathbf{T}^j\mathbf{R}\mathbf{Q}\mathbf{R}'\mathbf{T}^{j'}\right)\mathbf{z} + 1 \right] \sigma_*^2. \quad (6.9)$$

Ennusteen **ennustevälit** (*prediction intervals*) saadaan yhtälön (6.9) neliöjuuren **RMSE:n** (*root mean square error*) avulla. Kun ψ_* tunnetaan, $100(1 - \alpha)\%$:n ennustevälit (EV) saadaan kaavalla

$$\text{EV}(\tilde{y}_{T+l|T}) = \tilde{y}_{T+l|T} \pm t_{T-d}^{(\alpha/2)} \cdot \text{RMSE}(\tilde{y}_{T+l|T}), \quad (6.10)$$

jossa $t_{T-d}^{(\alpha/2)}$ on t -jakauman arvo merkitsevyytasolla $\alpha/2$. RMSE on yhtälön (6.9) oikean puolen neliöjuuri, jossa σ_*^2 on korvattu harhattomalla estimaattorillaan.

6.2 Huomautuksia

Usein ollaan tilanteessa, jossa aikahomogeenisuus- ja normaalisuusoletusten toteutumiseksi analysoitavaa sarjaa on jouduttu muuntamaan, esimerkiksi logaritimuunnoksella. Yleensä kuitenkin halutaan ennusteet alkuperäiselle sarjalle. Saadut ennusteet voidaan muuntaa suoraan takaisin alkuperäiseen muotoon, mutta takaisinmuunnokseen sisältyy harhaa, eli ennusteet eivät enää ole harhattomia. Harhan osuus on Harveyn mukaan niin pientä, ettei se merkittävästi muuta ennusteita. Harhan korjaamiseen on olemassa muunnoskaava, jota mm. STAMP-ohjelmisto käyttää. Muunnoskaava logaritmisoidusta ennusteesta logaritmisoimattomaan on

$$\hat{y}_{T+j}^* = \exp(\hat{y}_{T+j} + 0.5\hat{r}_{T+j}) - \exp(\hat{y}_{T+j}), \quad j = 1, \dots, L, \quad (6.11)$$

jossa

$$\hat{r}_{T+j} = \hat{\sigma}^2 \hat{f}_{T+j|T}. \quad (6.12)$$

Jälkimmäisessä f noudattaa aiemmin määriteltyä muotoa (3.12c) ja $\hat{\sigma}^2$ on varianssin estimaattori.

Jatkossa ennustamista käsitellään vielä selittävien muuttujien yhteydessä, jolloin tarkastellaan, kuinka selittävien muuttujien lisääminen malliin vaikuttaa ennustamiseen.

7. Perusrakennemallin käytöstä

Edellä on käsitelty rakenneyhtälömallin rakentamista ja ominaisuuksia, lähinnä perusrakennemallin tapauksessa. Seuraavassa tarkastellaan, mitä käyttömahdollisuuksia perusrakennemallilla on.

7.1. Ennustaminen

Edellisessä luvussa tarkasteltiin jo kuinka rakenneyhtälömallia voidaan soveltaa ennustamistarkoituksiin. Usein aikasarjamallin rakentamisen ja koko aikasarja-analyysin tavoitteena on juuri ennusteiden laatiminen. Monissa tilanteissa tehdään kauaskantoisiakin päätöksiä pitkälti juuri ennusteisiin perustuen. Tämän vuoksi ennusteiden perustana olevan tilastollisen mallin rakentajalla ja ennusteiden laatijalla voi olla välillisesti hyvinkin suuri vastuu tehtävistä päätöksistä. Mallin ja ennusteiden laatijan tehtävänä on siis löytää paras mahdollinen malli, johon laadittavat ennusteet pohjautuvat. Toisin sanoen laadittujen ennusteiden pohjana olevaan tilastolliseen malliin on saatava sisällytettyä kaikki ilmiöstä käytettävissä oleva informaatio.

Oleellinen tekijä on myös mallin perustana olevan aineiston luotettavuudesta varmistuminen. Mikäli tutkimusaineiston validiteetti on huono, toisin sanoen aineisto kuvaa tutkittavaa ilmiötä huonosti, on mahdotonta tuottaa tutkittavaa ilmiötä luotettavasti kuvaavia ennusteita, vaikka aineistosta muodostettu malli olisi kuinka hyvä tahansa.

Vaikka tilastollisen mallin pohjana oleva aineisto olisi validi ja muodostettu malli paras mahdollinen, ei pelkkä laadittu ennuste tarjoa riittävästi informaatiota, sillä ennusteen käyttäjä tarvitsee myös tietoa ennusteen tarkkuudesta. Aikasarjan rakenneyhtälömallien etuna tässä suhteessa on se, että rakenneyhtälömallin teoria on yhdistetty klassiseen testiteoriaan. Näin laadituille ennusteille voidaan tuottaa myös ennuste- eli luottamusvälit (luku 6.1), jotka kertovat laaditun ennusteen tarkkuudesta.

7.2. Aikasarjatasoitus

Aikasarjatasoituksella (*smoothing*) pyritään pääsääntöisesti selvittämään tilavektorin α_t olemassaolevalla informaatiolla ehdollistettu odotustusarvo, eli $E(\alpha_t | y_t)$. Tämä voi antaa arvokasta informaatiota, kun halutaan esimerkiksi tietää miten komponentti kuten trendi on kehittynyt aiemmin. Tilavektorin ehdollistettua odotusarvoa voidaan kutsua myös **tasoite- tuksi estimaatiksi** (*smoothed estimate*) ja sen estimaattoria **tasoittimeksi** (*smoother*). Tasoitettu estimaatti voidaan kirjoittaa muodossa

$$\mathbf{a}_{t|T} = E(\alpha_t) = E(\alpha_t | y_t). \quad (7.1)$$

Lineaarille malleille on olemassa kolme eri tasoitusalgoritmia: **pistetasoitus** (*fixed-point smoothing*), **viivetasoitus** (*fixed-lag smoothing*) ja **välitasoitus** (*fixed-interval smoothing*). Ensin mainittua käytetään, kun halutaan tilavektorin tasoitettuja estimaatteja jollakin kiinteällä ajanhetkellä, eli saadaan $\mathbf{a}_{t|T}$ tietyille τ :n arvoille kaikilla hetkillä ajanjaksolla $t > \tau$. Viivetasoi-

tuksella saadaan tasoitettu estimaatti $\mathbf{a}_{t-j|t}$, kun $j = 1, \dots, M$ ja M on jokin maksimaalinen viive. Välitasoituksella saadaan puolestaan tasoitetut estimaatit jollekin kiinteälle jänneväliille.

Koopman (1993) on esittänyt menetelmän, jolla jäännösvektorin tasoittamisen kautta päästään myös tehokkaaseen tilavektorin tasoittimeen. Menetelmä hyödyntää **apujäännöksiä** (*auxiliary residuals*) ja EM-algoritmia.

7.3 Trendipuhdistus

Trendipuhdistuksella (*detrending*) tarkoitetaan trendin vaikutuksen poistamista sarjasta. Ennen kuin voidaan tarkemmin tarkastella asiaa, on syytä määritellä trendi. Harveyn mukaan **trendi on se osa aikasarjaa, joka ekstrapoloituna antaa selvimmän kuvan aikasarjan pitkän aikavälin kehityksestä.**

Trendipuhdistettu (T_p) sarja voidaan kirjoittaa muodossa

$$\tilde{y}_t^{Tp} = y_t - \tilde{\mu}_t. \quad (7.2)$$

Etenkin taloudellisen tutkimuksen piirissä trendipuhdistus on yleinen käytäntö ennen aikasarjan tilastollisten analyysien suorittamista. Harveyn mukaan trendin puhdistamista ei tulisi käyttää kuin korkeintaan erittäin harvinaisissa tilanteissa, joissa kyseessä on deterministinen lineaarinen trendi stationarisessa prosessissa. Stokastisen trendikomponentin tapauksessa trendipuhdistus ei ole kuitenkaan viisasta, sillä se johtaa helposti harhaanjohtaviin tuloksiin.

7.4 Kausipuhdistus

Kausipuhdistuksella (*deseasonalising*) tarkoitetaan kausivaihtelun poistamista analysoitavasta sarjasta. Estimoitu kausivaihtelu voidaan kuukausiaineistolla määritellä **siksi osaksi sarjaa, joka ekstrapoloitaessa toistaa itsensä minkä tahansa vuoden aikana ja summautuu nolaksi vuosijaksoittain.**

Rakenneyhtälömallin kausivaihtelun ennustefunktio täyttää ehdon

$$\tilde{Y}_{T+l|T} = -\sum_{j=1}^{s-1} \tilde{Y}_{T+l-j|T}, \quad l = 1, 2, \dots \quad (7.3)$$

Niin dummy-muuttujiin perustuva kuin trigonometrisenkin kausivaihtelu täyttävät ehdon (7.3). Trigonometrisen kausivaihtelun osalta ehto voidaan kirjoittaa muodossa

$$\tilde{Y}_{T+l|T} = \sum_{j=1}^{[s/2]} (\gamma_{jT} \cos \lambda_{jt} + \gamma_{jT}^* \sin \lambda_{jt}), \quad l = 0, 1, 2, \dots \quad (7.4)$$

Kausivaihtelukomponentti ei sisällä mitään informaatiota aikasarjan yleisestä suunnasta tarkasteltiimpa sitten lyhyellä tai pitkällä aikavälillä. Niinpä joskus on järkevää suorittaa kausipuhdistus ja keskittyä tämän puhdistetun sarjan analysointiin. Kausipuhdistettu (K_p) ennuste saadaan kaavalla

$$\tilde{y}_{T+l|T}^{kp} = \tilde{y}_{T+l|T} - \tilde{\gamma}_{T+l|T}, \quad l = 1, 2, \dots \quad (7.5)$$

Kausivaihtelupuhdistus voi joissakin tapauksissa helpottaa rakenteellisten muutosten paikallistamista samoin kuin vaimenevan trendin havaitsemista. Tämän lisäksi kausipuhdistettu sarja on joissakin tapauksissa informatiivisempi kuin alkuperäinen sarja, koska kausivaihtelupuhdistus usein edesauttaa sarjan tulkinnallisuutta. Kysymys siitä, tulisiko esimerkiksi virallisten tilastojen osalta julkaista alkuperäinen vai kausipuhdistettu sarja, on monisyinen eikä yksikäsitteistä oikeaa ratkaisua liene olemassa.

7.5 Komponenttikohtainen tarkastelu

Komponenttikohtainen tarkastelu on luvuissa 7.2 - 7.4 kuvattuihin tarkastelutapoihin nähden erilainen lähestymistapa. Edellä ollaan oltu kiinnostuneita ilmiön käyttäytymisestä ilman jonkin yksittäisen komponentin vaikutusta ja tämän komponentin vaikutus on pyritty poistamaan mallista. Komponenttikohtaisessa analysoinnissa ei pyritä poistamaan minkään komponentin vaikutusta vaan ollaan kiinnostuneita siitä, kuinka tutkittavan ilmiön yksittäiset komponentit käyttäytyvät. Rakenneyhtälömallien tapauksessa yksittäisten komponenttien analysointi on sangen helppoa, koska kullekin komponentille saadaan estimaatti jokaiselle ajanhetkelle. Käytännössä siis analysoitava sarja saadaan purettua useisiin sarjoihin, jotka kukin kuvaavat alkuperäisen sarja yhtä komponenttia.

7.6 Huomautuksia

Kausivaihtelun käsittelemiseen on toki muitakin menetelmiä, kuten X11- (Shiskin et. al., 1967) ja X11-ARIMA-menetelmät (Dagum, 1975). Ensin mainittu perustuu erilaisten suotimien käyttöön ja on virallisen tilastotoimen (Suomessa Tilastokeskus) piirissä paljon käytetty menetelmä. X11-ARIMA perustuu X11-menetelmään. Periaatteena X11-ARIMAssa on se, että muodostetaan tasoitettavaa sarjaa kuvaava ARIMA-malli, laaditaan ennuste ja tämän jälkeen koko sarja (alkuperäinen ja ennuste) käsitellään X11-menetelmällä. Harveyn mukaan X11-menetelmän ongelmana on joustamattomuus, menetelmä käsittelee kaikkia sarjoja samalla tavalla ja joskus X11 tuottaa tasoitettuun sarjaan epätoivottuja ominaisuuksia. Edellisessä kappaleessa yhtälöissä (7.3) - (7.5) lyhyesti kuvattu **mallipohjainen kausitasoitus** (*model-based seasonal adjustment*) on erittäin joustava menetelmä, sillä se on räätälöity tietylle sarjalle, eikä mallin koostumus tarvitse olla samanlainen kaikille sarjoille. Verrattaessa esimerkiksi X11-menetelmää ja mallipohjaista tasoitusta on osoitettavissa, että X11 tuottaa vaimempia trendikomponentteja kuin mallipohjainen menetelmä. Lisäksi X11-menetelmällä tuotetun tasoitetun sarjan satunnaistekijä on varianssiltaan suurempi kuin mallipohjaisen menetelmän tuottama.

8. Selittävät muuttujat

Seuraavassa tarkastellaan selittävien muuttujien lisäämisen vaikutusta rakenneyhtälömalliin ja mallin käyttöön. Selittäviä muuttujia sisältävä rakenneyhtälömalli supistuu perinteiseksi regressiomalliksi, kun mallin stokastiset komponentit jäännöstermiä lukuunottamatta poistetaan.

8.1 Johdanto

8.1.1 Eksogeisuus

Tähän saakka analysoitavaa sarjaa kuvaava malli on rakennettu ilmiön aiempiin havaintoihin ja vallitsevaan ajanhetkeen perustuen. Nyt malliin lisätään k kappaletta muuttujia, jotka kykenevät jollakin tapaa selittämään analysoitavan sarjan liikkeitä. Oletetaan, että lisättävät muuttujat ovat ainakin **heikosti eksogeenisiä** (*weakly exogenous*). Heikon eksogeisuuden muodollinen määritelmä on se, että voidaan muodostaa yhteistiheysfunktio malliin lisättävien X :ien ja y :n välille, toisin sanoen on olemassa

$$F(y, X; \lambda) = \prod_{t=1}^T p(y_t, x_t | y_{t-1}, X_{t-1}; \lambda), \quad (8.1)$$

jossa λ :lla merkitään parametreja, joista y ja X riippuvat.

Heikon eksogeisuuden lisäksi oletetaan, että suhde selitettävän sarjan y_t ja selittäjien X_t välillä on lineaarinen. Tällöin perusrakennemalli voidaan laajentaa muotoon

$$y_t = \mu_t + \gamma_t + x_t' \delta + \varepsilon_t, \quad (8.2)$$

jossa x_t on k kappaletta selittäviä muuttujia sisältävä $k \times 1$ vektori ja δ on selittäviin muuttujiin liittyviä parametreja sisältävä $k \times 1$ vektori. Malliin voidaan lisätä jo aiemmin käsiteltyyn tapaan myös muita komponentteja, kuten sykli tai päivittäisvaikutukset. Oletus heikosta eksogeisuudesta merkitsee sitä, että selittävät muuttujat voidaan ehdollistaa ilman, että menetetään estimoinnin kannalta olennaista informaatiota. Selitettävän muuttujan y_t ennusteet tehdään ehdollistettuna x_t :llä eli laaditut ennusteet voivat perustua selittäivistä muuttujista laadittuihin ennusteisiin tai erilaisiin tulevaisuuden skenaarioihin.

Vahva eksogeisuus (*strong exogeneity*) on kyseessä silloin, kun selitettävän ja selittävien muuttujien välillä ei ole **moleminsuuntaista riippuvuutta** (*feedback*). Vahvan eksogeisuuden oletuksen ollessa voimassa voidaan mallia käyttää myös ennustustarkoituksiin. Formaalisti sama voidaan esittää ehdon

$$p(x_t | y_{t-1}, X_{t-1}; \theta_x) = p(x_t | X_{t-1}; \theta_x), \quad t = 1, \dots, T \quad (8.3)$$

muodossa. Harveyn mukaan käytännössä on erittäin vaikea testata suoranaisesti muuttujien eksogeisuutta, mutta epäsuorasti se on mahdollista testaamalla mallin stabiiliutta. Jos

muuttujat eivät ole heikosti eksogeenisiä ei malli myöskään yleensä ole stabiili. Tarkemmin eksogeenisuutta ja stabiiliutta ovat käsitelleet Engle et al. (1983) ja Engle (1984).

8.1.2 Erilaisia selitysmalleja

Perusrakennemalli (8.1) voidaan kirjoittaa tila-aika-mallin muodossa, jolloin mittausyhtälö on muotoa

$$y_t = \mathbf{z}'_t \alpha_t + \mathbf{x}'_t \delta + \varepsilon_t, \quad t = 1, \dots, T. \quad (8.4)$$

Luvussa 3 esitettyyn muotoon verrattuna laajennuksena on se, että tekijä d_t korvataan tekijällä $\mathbf{x}'_t \delta$, jolloin yhtälön avulla voidaan esittää aiempaa monimutkaisempia malleja. Koska δ on yleensä tuntematon, on kannattavaa ottaa se mukaan tilavektoriin α_t , jolloin saadaan **laajennettu tilavektori** (*augmented state vector*) $\alpha_t^\diamond = [\alpha'_t \delta'_t]'$. Tällöin tila-aika-malli voidaan kirjoittaa muodossa

$$y_t = [\mathbf{z}'_t \ \mathbf{x}'_t] \alpha_t^\diamond + \varepsilon_t, \quad t = 1, \dots, T \quad (8.5)$$

ja

$$\alpha_t^\diamond = \begin{bmatrix} \alpha_t \\ \delta_t \end{bmatrix} = \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \alpha_{t-1} \\ \delta_{t-1} \end{bmatrix} + \begin{bmatrix} \eta_t \\ \mathbf{0} \end{bmatrix}. \quad (8.6)$$

Siirtymäyhtälön (8.6) alempi osa merkitsee käytännössä sitä, että $\delta_t = \delta$, eli δ on aikainvariantti. Mallia voitaisiin laajentaa vielä stokastisoimalla termi δ , jolloin

$$\delta_t = \delta_{t-1} + \mathbf{v}_t,$$

jossa \mathbf{v}_t oletetaan normaalijakautuneeksi odotusarvovektorilla $\mathbf{0}$ ja kovarianssimatriisilla σ_δ^2 . Jatkossa kuitenkin oletetaan termi aikainvariantiksi ja deterministiseksi.

Joskus voidaan käyttää myös selitettävän muuttujan viiveitä selittäjinä, jolloin kyseessä on **dynaaminen malli**. Tällöin malli voidaan kirjoittaa muodossa

$$y_t = \phi_1 y_{t-1} + \dots + \phi_r y_{t-r} + \mathbf{z}'_t \alpha_t + \mathbf{x}'_t \delta + \varepsilon_t, \quad t = r + 1, \dots, T. \quad (8.7)$$

Joissakin tapauksissa mallissa käytetään myös selittävien muuttujien viiveitä, jolloin kyseessä on **jakautuneitten viiveiden malli** (*distributed lags model*). Viivemallien rakentamisessa käytetään apuna sisältöteoreettista tietämystä, graafisia tarkasteluja ja ristikorrelaatioanalyysiä. Ongelmana viivemallien estimoinnissa saattaa olla eri viiveiden välinen multikollineaarisuus.

Jakautuneiden viiveiden mallin rakenne riippuu luonnollisesti siitä, kuinka viiveet ovat jakautuneet. Viiveet voivat olla esimerkiksi polynomialisesti tai suhteellisesti jakautuneet. Ensin mainitussa tapauksessa viiveiden välillä on polynomirajoitteet ja malli voidaan kirjoittaa yksinkertaisimmillaan yhden selittävän muuttujan tapauksessa muodossa

$$y_t = \mathbf{z}'_t \boldsymbol{\alpha}_t + \sum_{\tau=0}^h \delta_\tau x_{t-\tau} + \varepsilon_t. \quad (8.8)$$

Tällöin kyseessä on **polynomialisesti jakautuneiden viiveiden malli** (*polynomially distributed lags*), jonka on esittänyt Almon (1965). Sen sijaan, että ryhdyttäisiin estimoimaan yhtälön (8.8) keskimmäisen termin kaikkia $(h + 1)$:tä kerrointa, oletetaan kertoimien muodostavan g :nnen ($< h$) asteen polynomien, jolloin viivemallissa estimoitavaksi jää vain $g + 1$ parametria.

Jorgenssonin (1966) esittämän **suhteellisesti jakautuneitten viiveiden mallin** (*rationaly distributed lags*) viiverakenne voidaan esittää viiveoperaattorin kahden polynomien suhteen avulla. Näin muotoiltu rakenneyhtälömalli voidaan kirjoittaa muodossa

$$y_t = \mathbf{z}'_t \boldsymbol{\alpha}_t + \sum_{i=0}^k \frac{\omega_i(L)}{\lambda_i(L)} x_{t-\tau} + \varepsilon_t, \quad (8.9)$$

jossa $\omega_i(L)$ ja $\lambda_i(L)$ ovat äärellistä astetta olevia polynomiviiveitä. Mallin supistettu muoto

$$y_t = \sum_{i=0}^k \frac{\omega_i(L)}{\lambda_i(L)} x_{t-\tau} + u_t, \quad (8.10)$$

jossa u_t on ARIMA-prosessi, tunnetaan myös Boxin ja Jenkinsin (1976) esittämänä **siirtofunktiomallina** (*transfer function model*). Useiden selittävien muuttujien tapauksessa operaatio aika-alueessa tulee Harveyn mukaan sekavaksi ja epäsystemaattiseksi, jonka vuoksi siirtofunktioviivemallin (8.10) käyttö ei ole järkevää.

Selittävän muuttujan viivemalli (*lagged dependent variable*) on muotoa

$$y_t = \varphi_1 y_{t-1} + \dots + \varphi_r y_{t-r} + \mathbf{z}'_t \boldsymbol{\alpha}_t + \sum_{i=1}^k \omega_i(L) x_{it} + \varepsilon_t, \quad (8.11)$$

jossa $\omega_i(L)$:t ovat joukko polynomiviiveitä ja $\varphi_1, \dots, \varphi_r$ ovat polynomien $\varphi(L)$ kertoimia. Mikäli viiveiden luonteeseen liittyviä rajoituksia ei ole, on kyseessä **rajoittamattomien viiveiden malli** (*unconstrained lags*). Tällöin mallin viiverakenne voidaan muotoilla muotoon

$$y_t = \mathbf{z}'_t \boldsymbol{\alpha}_t + \mathbf{x}'_t \boldsymbol{\delta} + \sum_{\tau=0}^{h-1} \Delta \mathbf{x}'_{t-\tau} \boldsymbol{\delta}^\diamond + \varepsilon_t, \quad (8.12a)$$

jossa

$$\boldsymbol{\delta}^\diamond_\tau = - \sum_{j=\tau+1}^h \boldsymbol{\delta}_j, \quad \tau = 0, 1, \dots, h-1 \quad (8.12b)$$

ja

$$\delta = \sum_{j=0}^h \delta_j. \quad (8.12c)$$

Kaavan (8.12) muoto johtuu siitä, että muuntamalla x -muuttujat differensseiksi

$$x_t, \Delta x_t, \dots, \Delta x_{t-h+1}$$

niiden välillä esiintyy alkuperäistä vähemmän multikollineaarisuutta.

8.2. Selittävät muuttujat ja parametrien estimointi

Estimointimenetelmät noudattelevat jo aiemmin estimointia käsitelleessä kappaleessa läpikäytyjä periaatteita, mutta selittävien muuttujien mukaanotto aiheuttaa estimointiin joitakin laajennuksia, joita tarkastellaan seuraavassa.

Aika-alueessa parametrien estimoimiseksi sovelletaan ennustevirrehajotelmaa, jonka avulla uskottavuusfunktio saadaan konstruoitua (luku 3.4). Aika-alueessa estimointaessa on olemassa kaksi mahdollista su-estimaattoria. Estimaattori riippuu siitä, miten δ -parametrivektoria käsitellään. Vektori voidaan lisätä tilavektoriin, tai se voidaan estimoida **GLS-menetelmällä** (*GLS transformation method*). Mikäli jälkimmäistä käytetään, voidaan malli (8.4) kirjoittaa muotoon

$$y_t = x_t' \delta + u_t, \quad t = 1, \dots, T \quad (8.13a)$$

$$u_t = z_t' \alpha + \varepsilon_t, \quad (8.13b)$$

olettaen, että u_t :n odotusarvo on nolla. Termi on Harveyn mukaan yleensä heteroskedastinen ja korreloitunut. GLS-estimaattori δ :lle on

$$\tilde{\delta} = (x' v^{-1} x)^{-1} x' v^{-1} y, \quad (8.14a)$$

jossa

$$v^{-1} = \frac{1^2}{f}. \quad (8.14b)$$

Määritellään $y^* = ly$, $x^* = lx$ ja $u^* = lu$, saadaan regressioyhtälö

$$y_t^* = x_t^* \delta + u_t^*. \quad (8.15)$$

GLS-estimaattori δ :lle on tällöin muotoa

$$\tilde{\delta} = \left[\sum_{t=1}^T \frac{x_t^* x_t^{*'}}{f_t} \right]^{-1} \sum_{t=1}^T \frac{x_t^* y_t^*}{f_t}. \quad (8.16)$$

Jos oletetaan, että δ tunnetaan, voidaan mittausyhtälö kirjoittaa muotoon

$$y_t - \mathbf{x}_t' \delta = \mathbf{z}_t' \alpha_t + \varepsilon_t, \quad t = 1, \dots, T. \quad (8.17)$$

Soveltamalla Kalmanin suodinta yhtälöön (8.17), saadaan yhtälön 'havainnoille' ennustevirheet v_t , jotka voidaan määritellä ehdolla ψ muodossa

$$v_t = y_t^* - \mathbf{x}_t^{*'} \delta, \quad t = 1, \dots, T. \quad (8.18)$$

Koska

$$l(y - \mathbf{x}\delta) = ly - l\mathbf{x}\delta = y^* - \mathbf{x}^{*'} \delta,$$

samaa Kalmanin suodinta voidaan käyttää erikseen havainnoille y_t ja selittäville muuttujille \mathbf{x}_t -vektorissa. Näin siis GLS-estimaattori $\tilde{\delta}(\psi)$ saadaan regressoimalla y_t :n ennustevirheet y_t^* \mathbf{x}_t :n ennustevirheillä \mathbf{x}_t^{*} .

Toinen tapa δ :n estimaatin laskemiseksi on δ :n sisällyttäminen tilavektoriin. Tila-aika-malli on tällöin muotoa

$$y_t = [\mathbf{z}_t' \quad \mathbf{x}_t'] \alpha_t^\diamond + \varepsilon_t, \quad t = 1, \dots, T \quad (8.19a)$$

$$\alpha_t^\diamond = \begin{bmatrix} \alpha_t \\ \delta_t \end{bmatrix} = \begin{bmatrix} \mathbf{T}_t & 0 \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \alpha_{t-1} \\ \delta_{t-1} \end{bmatrix} + \begin{bmatrix} \eta_t \\ \mathbf{0} \end{bmatrix}. \quad (8.19b)$$

Tila-aika-mallista (8.19) on johdettavissa uskottavuusfunktio, joka on muotoa

$$\log L(\psi; \mathbf{y}) = -\frac{T}{2} \log 2\pi - \frac{1}{2} \sum_{t=1}^T \log f_t - \frac{1}{2} \sum_{t=1}^T \frac{(y_t^* - \mathbf{x}_t^{*'} \delta)^2}{f_t}, \quad (8.20)$$

joka voidaan saattaa tiivistettyyn (*concentrated*) muotoon

$$\log L(\psi) = -\frac{T}{2} \log 2\pi - \frac{1}{2} \sum_{t=1}^T \log f_t - \frac{1}{2} \sum_{t=1}^T \frac{v_t^2}{f_t}, \quad (8.21)$$

jossa v_t on määritelty kuten yhtälössä (8.18).

Tilavektorin laajennuksen kautta tapahtuma estimointi tuottaa Harveyn mukaan teoreettisesti korrektein tuloksen, kun δ :aa käsitellään satunnaismuuttujana. GLS-menetelmä puolestaan tuottaa korrektein tuloksen, kun δ :aa käsitellään kiinteänä vektorina. Käytännössä menetelmien tuottamat tulokset eivät Harveyn mukaan juurikaan poikkea toisistaan.

Estimoitavana on kuitenkin vielä σ_*^2 . Annetuilla suhteellisten hyperparametrien arvoilla (ψ_*) molemmat menetelmät tuottavat saman jäännösneliösumman $S(\psi_*)$. σ_*^2 :n su-estimaattorit poikkeavat kuitenkin toisistaan, käsiteltäessä δ :aa satunnaismuuttujana (laajennetun tilavektorin menetelmä), estimaattori on muotoa

$$s_*^2 = (T - d - k)^{-1} \sum_{t=d+1}^T \tilde{v}_t^{02}, \quad (8.22)$$

kun taas GLS-menetelmän mukainen estimaattori on

$$\tilde{\sigma}_*^2 = (T - d)^{-1} \sum_{t=d+1}^T \tilde{v}_t^2. \quad (8.23)$$

Kuten edellä jo todettiin molemmat menetelmät tuottavat saman jäännösneliösumman eli

$$S(\psi_*) = \sum_{t=d+1}^T \tilde{v}_t^{02} = \sum_{t=d+1}^T \tilde{v}_t^2, \quad (8.24)$$

joten ero estimaattorien välille muodostuu jakotermissä. Jäännöksiin perustuen voidaan mallin riittävyttä arvioida useilla eri testeillä. Osaa testeistä on käsitelty jo luvussa 5 ja osaa käsitellään seuraavassa.

8.3 Selittävien muuttujien valinta ja mallin diagnostiikka

Seuraavaksi tarkastellaan mihin kriteereihin selittävien muuttujien malliin ottaminen perustuu ja miten se vaikuttaa luvussa 5 tarkasteltuihin mallin diagnostisiin tarkasteluihin.

Ennen kuin selittäviä muuttujia sisältävää mallia ryhdytään rakentamaan on syytä rakentaa malli, jossa selittäviä muuttujia ei vielä ole. Tällainen **yksimuuttujainen** (*univariate*) malli antaa usein hyödyllistä informaatiota myös selittäviä muuttujia sisältävän mallin rakentamiseen. Joskus malliin tuotavat selittävät muuttujat korvaavat yksimuuttujaisen mallin komponentteja, kuten esimerkiksi kausivaihtelun. Myös selitettävän ja selittävien sarjojen graafinen tarkastelu voi paljastaa, kuinka selittävät muuttujat mahdollisesti pystyvät selitettävää sarjaa selittämään. Mikäli selitettävä sarja on jouduttu esimerkiksi differensoimaan, on selittävät sarjat differensoitava vastaavasti.

Selittäviä muuttujia sisältävän mallin valinta perustuu paljolti samoihin kriteereihin kuin yksimuuttujainenkin mallin tapauksessa. Näiden jo aiemmin luvussa 5 käsiteltyjen kriteerien lisäksi mallilta edellytetään sitä, että selittäville muuttujille estimoidut parametrit eivät ole ristiriidassa selittävien sarjojen määrittelyjen kanssa. Arvioitaessa mallin riittävyttä kuvaavilla mittareilla selittäville muuttujilla varustetun mallin tulisi olla myös vähintään yhtä hyvä kuin 'kilpailevat' mallit.

Jäännöstesteillä voidaan myös selittävien mallien yhteydessä testata onko aikasarjassa rakenteellisia muutoksia ja onko malli virheellisesti spesifioitu. Tähän tarkoitukseen voidaan käyt-

tää CUSUM(T) -testistä (5.9) modifioitua **rekursiivista t-testiä** (*recursive t-test*). Sen testisuure on muotoa

$$\psi = \frac{\sum_{t=d+k+1}^T \tilde{v}_t^2}{\hat{s}_* \sqrt{T-d-k}} = \frac{\text{CUSUM}(T)}{\sqrt{T-d-k}} \quad (8.25a)$$

ja noudattaa t-jakaumaa vapausastein $T - d - k - 1$, mikäli mallissa ei ole viivästettyjä selitettäviä muuttujia ja ψ_* on tunnettu. Varianssiestimaatti testisuureessa on muotoa

$$\hat{s}_*^2 = (T - d - k - 1)^{-1} \sum_{t=d+k+1}^T (\tilde{v}_t^\circ - \bar{\tilde{v}}^\circ). \quad (8.25b)$$

Ennustevirheen varianssi voidaan estimoida muodossa

$$s^2 = s_*^2 \bar{f}, \quad (8.26)$$

tai

$$\tilde{\sigma}^2 = \tilde{\sigma}_*^2 \bar{f}. \quad (8.27)$$

Termit s_*^2 ja $\tilde{\sigma}_*^2$ ovat kuten yhtälöissä (8.22) ja (8.23) ja \bar{f} on f :n keskiarvo.

Ennustekyvyn testi (*post-sample predictive test*) määriteltiin yhtälössä (5.25). Selittävien muuttujien mallille testisuure on muotoa

$$\xi(l) = \frac{\sum_{j=1}^n v_{T+j}^{\circ 2}}{ls_*^2}. \quad (8.28)$$

Kun ψ_* on tunnettu, eikä selitettävän muuttujan viiveitä ole mallissa, standardoidut jäännökset ovat riippumattomia ja normaalijakautuneita odotusarvolla 0 ja varianssilla σ_*^2 . Tällöin testisuure noudattaa $F(l, T - n - k)$ -jakaumaa mallin ollessa oikein spesifioitu.

Selittävät muuttujat voivat olla keskenään voimakkaasti korreloituneita, kuten Harvey ja Phillips (1979) varoittavat. Tällöin pyörästysvirheillä on taipumusta kasvaa Kalmanin suotimessa. Tämän välttämiseksi voidaan Harveyn ja Phillipsin mukaan soveltaa esimerkiksi Kaminskin, Brysonin ja Schmidtin (1971) esittämää **neliöjuurialgoritmia** (*square root algorithm*).

8.4 Interventiot

Interventioanalyysissä on tehdään johtopäätöksiä tunnettujen tapahtumien vaikutuksesta tutkittavaan ilmiöön. Näitä tunnettuja tapahtumia mittaavat interventiot ovat itse asiassa dummy-muuttujia. Mallia (8.4) voidaan laajentaa interventiokomponentilla, jolloin malli on kirjoitettavissa muotoon

$$y_t = \mathbf{z}'_t \alpha_t + \mathbf{x}'_t \delta + \lambda w_t + \varepsilon_t, \quad (8.29)$$

jossa w_t on interventiomuuttuja ja λ sen kerroin. Muuttujan w_t määritelmä riippuu interventiön luonteesta. **Hetkellinen interventio** (*transitory effect*) vaikuttaa ainoastaan jollakin tietyllä hetkellä $t = \tau$, jolloin w_t on muotoa

$$w_t = \begin{cases} 0, & t \neq \tau \\ 1, & t = \tau \end{cases} \quad (8.30)$$

Hetkellisten interventioiden avulla voidaan myös käsitellä **poikkeavat havainnot** (*outliers*). **Tasomuutos** (*level change*) on puolestaan sarjan tasossa tapahtuva muutos, w_t on muotoa

$$w_t = \begin{cases} 0, & t < \tau \\ 1, & t \geq \tau \end{cases} \quad (8.31)$$

Tasomuutos voidaan esittää myös hetkellisenä muutoksena, joka kohdistuu trendikomponentin tasoparametriin. **Kulmakertoimen muutos** (*slope change*) aiheutuu kulmakerroinparametriin kohdistuvasta interventiosta, trendikomponentin kulmakerroin on tällöin muotoa

$$\mu_t = \mu_{t-1} + \beta_{t-1} + \lambda w_t + \eta_t. \quad (8.32)$$

Joskus saattaa tapahtua myös **muutos kausivaihtelussa** (*change in seasonal pattern*). Tällöin tarvitaan $s - 1$ dummy-muuttujaa kausivaihtelukomponentteja varten muutoshetkestä τ eteenpäin. Lisäksi oletetaan, että dummy-muuttujat summautuvat nolliin.

Myös muunlaisia interventiota voi tapahtua, w_t on muotoa

$$w_t = \begin{cases} 0, & t < \tau \\ \varphi^{t-\tau}, & t \geq \tau \end{cases} \quad (8.33)$$

kun hetkellisen intervention vaikutus vaimenee vähitellen.

Joskus, etenkin kuukausiaineiston tapauksessa, on syytä huomioida **juhlapyhien ajoittuminen eri kuukausille eri vuosina** (*trading day effect*). Tämä voidaan huomioida lisäämällä malliin muuttuja, joka sisältää informaation juhlapyhien sijoittumisesta eri kuukausille. Tilanne vastaa pitkälti normaalia selittävän muuttujan lisäämistä malliin, mutta joissakin yksinkertaisissa tapauksissa ongelma ratkaistavissa myös interventioiden avulla. Tätä ongelmaa, tosin aikasarjatasoituksen yhteydessä, ovat tarkastelleet Kitagawa ja Gersch (1984).

Interventioita voidaan siis käyttää tunnettujen sarjaan vaikuttavien tapahtumien, poikkeavien havaintojen ja rakenteellisten muutosten mallittamiseen. Poikkeavien havaintojen ja rakenteellisten muutosten havaitseminen tapahtuu yleensä graafisten tarkastelujen perusteella. On kuitenkin tilanteita, joissa graafisten tarkastelujen avulla ei voida nähdä vaikkapa rakenteellista muutosta. Tällöin voi olla hyödyllistä tarkastella **apujäännöksiä** (*auxiliary residuals*).

Apujäännökset ovat mallin **satunnaiskomponentin** (*irregular*) ja **tasoparametrin satunnaistekijän** (*level disturbances*) tasoitetut estimaatit (Harvey et al. 1992).

Intervention estimointi on mahdollista tehdä niin taajuus- kuin aika-alueessa operoitaessa. Estimointi voidaan suorittaa liittämällä interventiomuuttuja tilavektoriin muiden selittävien muuttujien tapaan.

Intervention luonteen määrittelyminen voi olla vaikeaa, mikäli intervention dynamiikkaa ei tunneta. Hetkellinen interventio sykäysluonteisena ei aiheuta pysyviä muutoksia tutkittavaan ilmiöön kuten taso- tai kulmakerroininterventiot. Intervention määrittelyssä Harveyn mukaan ainoa järkevä strategia on hyödyntää kaikki mahdollinen käytettävissä oleva informaatio, jonka perusteella muodostetaan interventiokomponentti, joka oikeellisuus pyritään varmistamaan diagnostisten testien avulla. Tarkoitukseen sopivia intervention väärin spesifioinnin testejä ovat esimerkiksi rekursiivinen t-testi tai CUSUM_τ(h)-testi. Ensin mainitun testisuure voidaan johtaa jälkimmäisen avulla ja on tässä tapauksessa muotoa

$$\psi_{\tau}(l) = \frac{\sum_{t=\tau+1}^{\tau+l} \tilde{v}_t^{\circ}}{l^{1/2} s_{*}(\tau-1)} = \frac{\text{CUSUM}_{\tau}(l)}{\sqrt{l}}, \quad (8.34)$$

jossa s_{*} on määritelty kuten (8.24), jossa summaus T:n sijasta $(\tau - 1)$:een saakka. Nollahypoteesin ollessa voimassa testisuure noudattaa t-jakaumaa vapausastein $(\tau - 1 - d - k)$.

Luvuissa 1 - 8 on nyt käsitelty rakenneyhtälömallin peruskehikko teoriatasolla ja tämän lisäksi on tarkasteltu joitakin rakenneyhtälömallin käytön sovellusmahdollisuuksia. Seuraavaksi sovelletaan käsiteltyä teoriakehikkoa empiiriseen aineistoon, joka koostuu kaukolämpöön ja sen kulutukseen liittyvistä aikasarjoista. Empiiristä tutkimusongelmaa ja tutkimuksen tavoitteita tarkastellaan seuraavassa.

9. Empiirinen ongelma ja tutkimusaineisto

Tutkimuksen teoriaosassa esiteltiin aikasarjojen rakenneyhtälömallin teoriaa, lähinnä perusrakennemallien osalta. Jatkossa käytetään tätä teoriakehikkoa empiirisen ongelman ratkaisuun. Empiirisessä sovelluksessa hyödynnetään rakenneyhtälömallien teoriaa niin ennustamisessa kuin mallin eri komponenttien analysoinnissa. Aluksi tarkastellaan sovellusesimerkkinä olevaa empiiristä ongelmaa. Esitellään tutkimuksen tavoitteet ja tarkastellaan tutkittavan ilmiön luonnetta. Tämän jälkeen tarkastellaan analyysien kohteena olevaa havaintoaineistoa ja sen ominaisuuksia.

9.1 Empiirinen ongelma

9.1.1 Tutkimuksen tavoitteet

Empiirisenä sovelluksena laaditaan erilaisia Jyväskylän kaupungin kaukolämmön (*district heating*) kulutusta kuvaavia malleja. Mallien perusteella laaditaan sitten ennusteita seuraavan vuoden kulutukselle. Jotta mallien ennustekykyä pystytään arvioimaan, laaditaan malli vuosien 1989 - 1995 aineistolla, jolla sitten tuotetaan ennuste vuoden 1996 kulutukselle. Koska vuoden 1996 osalta toteutuneet kulutustiedot on käytettävissä, voidaan laadittuja ennusteita verrata toteutuneeseen kulutukseen.

Lopuksi tuotetaan vielä aito ennuste vuoden 1997 kaukolämmön kulutukselle. Jotta saataisiin kohtuullinen kokonaiskuva empiirisen ongelman luonteesta on ensin kuitenkin syytä tarkastella mitä kaukolämmöllä oikein tarkoitetaan, miten se toimii ja missä sitä käytetään.

9.1.2 Kaukolämpö

Kaukolämpö on Suomen yleisin lämmitysmuoto. Sen markkinaosuus vuonna 1995 oli 46% rakennuskannasta ja myynnin arvo 4.4 miljardia markkaa. Kaukolämmön kulutus jakautui vuonna 1995 siten, että asumisen osuus oli 56%, teollisuuden 10% ja muun kulutuksen 34%. Suomen ohella Tanska on voimakas kaukolämmityksen hyödyntäjä. Kaukolämpöä tuotetaan useissa tapauksissa yhdessä sähkön kanssa. Tuotannossa eniten käytetyt polttoaineet ovat kiivihiili (n. 40%), maakaasu (n. 20%) ja turve (n. 20%). Kehitys on suuntautunut kohti ympäristöystävällisempiä raaka-aineita, sillä vielä vuonna 1981 lähes puolet tuotetusta energiasta tuli öljystä, kun 1995 sen osuus oli alle 10 prosenttia (Suomen Kaukolämpö ry, 1995).

Tässä työssä tarkastellaan Jyväskylän Energia Oy:n tuottaman kaukolämmön kulutusta. Jyväskylän Energia Oy myy kaukolämpöä myös Jyväskylän Maalaiskunnalle, mutta tässä työssä rajoitutaan Jyväskylän kaukolämmön kulutuksen tarkasteluun. Jyväskylässä kaukolämmön piiriin kuuluu peräti yli 70% rakennuskannasta.

9.1.3 Kaukolämpöterminologiaa

Ennen kuin kaukolämpöjärjestelmän toimintaa voidaan tarkastella, on syytä määritellä eräitä kaukolämpötekniikan keskeisiä käsitteitä.

Jäähtymällä tarkoitetaan kiinteistöön tulevan ja sieltä lähtevän kaukolämpöveden lämpötilojen erotusta. Jäähtymä voidaan laskea kaavalla

$$\text{jäähtymä (}^\circ\text{C)} = (\text{lämmön kulutus (Mwh)} \times 860) / \text{kiertovesimäärä}, \quad (9.1)$$

jossa 860 on veden ominaiskerroin. Mitä suurempi jäähtymä on, sitä paremmin kaukolämpöjärjestelmä toimii.

Ominaiskulutus on lämmönkulutus suhteutettuna lämmitettävään tilavuuteen. Ominaiskulutuksen avulla saadaan eri kokoiset kiinteistöt keskenään vertailukelpoisiksi. Kiinteistön kaukolämpölaitteiston kautta virranneen veden määrä on **kiertovesimäärä**, jonka yksikkö on kuutiometri (m^3).

Astepäiväluvulla kuvataan säätilan, lähinnä lämpötilan vaihteluja. Astepäiväluku voidaan laskea kaavalla

$$\text{APL} = \sum_{i=1}^n (17 - \bar{t}_i), \quad (9.2)$$

jossa n on lämmitysvuorokausien lukumäärä ja \bar{t}_i lämmitysvuorokauden i keskilämpötila. Lämmitysvuorokausi on vuorokausi, jolloin

$$\bar{t}_i < 10^\circ\text{C} \quad (\text{tammi} - \text{kesä kuu})$$

$$\bar{t}_i < 12^\circ\text{C} \quad (\text{heinä} - \text{joulukuu}).$$

Loppuvuodelle on lämmitysvuorokaudelle asetettu alkuvuotta lämpimämpi yläraja. Tätä perustellaan sillä, että keväällä auringon lämmittävä vaikutus on syksyä suurempi (Helminen, 1989).

Kaavassa (9.2) sisälämpötilaksi on asetettu 17. Normaalisti sisälämpötilaksi on määritelty 20°C . Poikkeama määrittelyssä johtuu siitä, että huonelämpötilan oletetaan kohoavan muista kuin lämmityksestä johtuvista syistä kolme astetta.

9.1.4 Kaukolämmön toimintaperiaate

Kaukolämpöä saadaan lämpöä ja sähköä tuottavista lämpövoimalaitoksista ja lämpökeskuksista. Lämpö siirretään asiakkaille kaukolämpöverkossa kiertävän kuumen veden avulla. Kaukolämpöveden lämpötila vaihtelee säätilan mukaan $65^\circ - 115^\circ\text{C}$ välillä. Veden kiehumispisteen ylittävä veden lämpötila mahdollistuu kaukolämpöverkostossa vallitseva korkeamman paineen avulla. Alhaisimmillaan lämpötila on kesällä jolloin lämmitystarve rajoittuu vain lämpi-

män käyttöveden lämmittämiseen. Asiakkailta tuotantolaitoksiin palaavan veden lämpötila vaihtelee 25° - 55°C välillä.

Kiinteistön käyttämä lämpömäärä mitataan lämpöenergiamittarilla. Mittarin osat ovat virtausanturi, lämpömäärälaskin ja lämpötila-anturit. Virtausanturi mittaa kiertävän kaukolämpöveden määrän. Lämpötila-anturit mittaavat kiinteistöön tulevan ja kiinteistöstä lähtevän kaukolämpöveden lämpötilan. Lämpömäärälaskin laskee kiinteistön ja lämpimän käyttöveden lämmitykseen käytetyn lämpöenergian virtausanturilta ja lämpötila-antureilta saatujen tietojen perusteella. Käytetyn lämpöenergian määrä lasketaan kaavalla

$$\text{lämmönkulutus (kWh)} = \text{jäähtymä (° C)} \times \text{kiertovesimäärä (m}^3\text{)} \times 1.163, \quad (9.3)$$

jossa 1.163 on muunnosvakio (1 kcal / h = 1.163 W). Verrattaessa kaavaa (9.3) kaavaan (9.1) huomataan, että ratkaisemalla kaava (9.1) lämmön kulutuksen suhteen päädytään kaavaan (9.3). Kulutuksen suuresta volyyymistä johtuen käytännöllisin mittayksikkö lämpöenergian määrälle on tässä yhteydessä gigawattitunti (GWh), joka vastaa tuhatta megawattituntia. Kaukolämmön kokonaiskulutukseksi saadaan näin ollen kaikkien kaukolämmön kuluttajien eli kulutuspaikkojen yhteenlaskettu kaukolämmön kulutus. Kokonaiskulutus voidaan siten esittää muodossa

$$\text{Kokonaiskulutus(GWh)} = \frac{\sum_{i=1}^N (j_i \times k_i \times 1.163)}{1000}, \quad (9.4)$$

jossa j_i = kulutuspaikan i jäähtymä,
 k_i = kulutuspaikan i kiertovesimäärä,
 N = kulutuspaikkojen lukumäärä,
1.163 = vakio.

9.2 Tutkimusaineisto

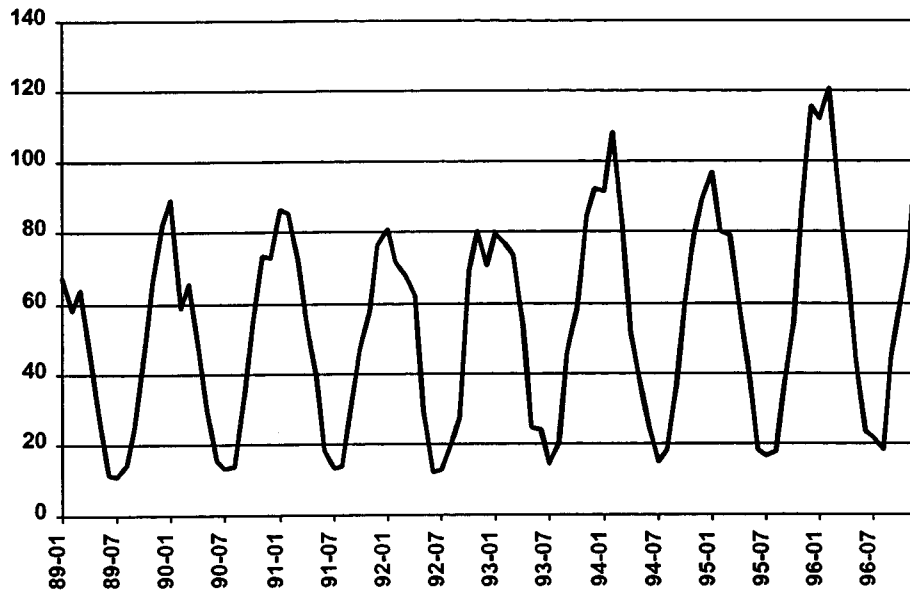
9.2.1 Aineistosta yleisesti

Tässä työssä tarkasteltava aineisto käsittää kaukolämmön kokonaiskulutuksen kuukausiaikasarjan vuoden 1989 alusta vuoden 1996 joulukuuhun ($T = 96$). Lisäksi hyödynnetään Jyväskylän astepäivälukujen aikasarjaa samalta ajanjaksolta. Astepäiväluku mitataan Ilmatieteen laitoksen toimesta Luonetjärven lentokentällä.

9.2.2 Analysoitavat muuttujat

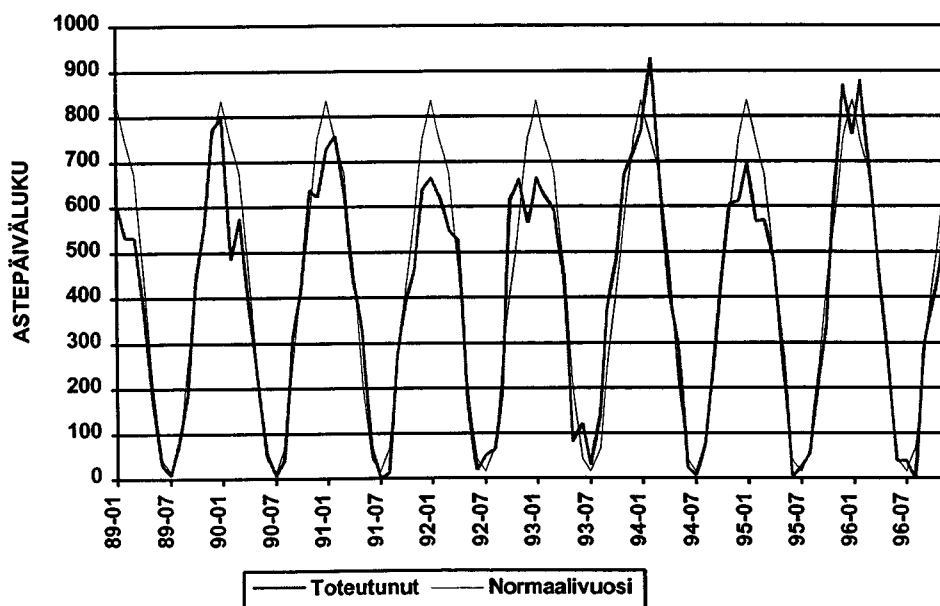
Selitettävä muuttuja on siis Jyväskylän Energia Oy:n Jyväskylän alueelle myymän **kaukolämmön kokonaiskulutus**. Selittävänä muuttujana käytetään puolestaan **astepäivälukua**, joka kuvaa säätilan aiheuttamaa lämmityksen tarvetta. Havaintoaineiston perusteella laaditaan erilaisia kulutusta kuvaavia malleja, joilla tuotetaan sitten kokonaiskulutusennusteita vuodeksi eteenpäin. Kuten edellä jo todettiin aluksi laaditaan malleja, jotka perustuvat vuosien 1989 - 1995 havaintoihin. Näin voidaan vertailla erilaisten mallien tuottamia ennusteita toteutuneeseen.

seen kulutukseen. Vuoden 1996 todellisesta kulutuksesta on tässä yhteydessä poistettu vuoden 1996 alusta verkkoon liittyneen Säynätsalon osuus, koska mallien perustana olevassa aineistossakaan ei Säynätsalo ole mukana (kuvio 9.1).



Kuvio 9.1. Kaukolämmön kokonaiskulutus vuosina 1989 - 1996 (Säynätsalon osuus poistettu).

Lopuksi laaditaan vielä ennuste vuodelle 1997. Mallissa hyödynnetään selittävänä muuttujana astepäivälukua, jonka osalta laaditaan kolme erilaista säätilaskenaariota vuotta 1997 varten. Laaditaan ennuste vuosien 1961 - 1990 astepäivälukujen keskiarvoon perustuvan **normaalivuoden** sääolosuhteille sekä tätä kymmenen prosenttia lämpimämmälle ja viileämmälle säätyypille. Normaalivuoden ja havaintojakson astepäiväluvut nähdään kuviossa 9.2.



Kuvio 9.2. Toteutuneet astepäiväluvut 1989 - 1996 sekä normaalivuoden astepäiväluku.

9.2.3 Aineistoon tehdyistä korjauksista

Kuten nähdään aineisto on melko lyhyt ($T = 96$), mutta sen analysointia edesauttaa sarjojen sangen systemaattinen käyttäytyminen, verrattaessa kaukolämmön kulutussarjaa (kuvio 9.1) astepäivälukusarjaan (kuvio 9.2), nähdään sarjojen käyttäytyvän hyvin samankaltaisesti. On siis oletettavissa, että astepäiväluku on erittäin hyvä kaukolämmön kulutuksen selittäjä.

Alkuperäiseen aineistoon jouduttiin tekemään joitakin korjauksia. Vuosina 1989 - 1991 laskutuskäytännöstä johtuen kesä- ja heinäkuu laskutettiin yhteislaskulla, jolloin myös näiden kuukausien kulutus kirjautui heinäkuulle. Näiltä osin aineistoa korjattiin siten, että kulutus jaettiin kesä- ja heinäkuulle samassa suhteessa kuin kyseisten kuukausien lämmönhankinnan kanssa. Esimerkiksi vuonna 1991 Jyväskylän Energian lämmönhankinta oli kesäkuussa 21.453 gigawattituntia ja heinäkuussa 15.123 GWh. Näin ollen kesä-heinäkuun kokonaiskulutuksesta $58.7\% (21.453 / (15.123 + 21.453) * 100\% = 58.7\%)$ sijoitettiin kesäkuulle ja loput heinäkuulle. Vastaava korjaus tehtiin touko- ja kesäkuun kohdalla vuonna 1992, jolloin nämä kuukaudet olivat yhteislaskutuksessa.

Edellä mainittujen lisäksi korjattiin vuoden 1996 marras- ja joulukuun kulutuksia. Tämä johtui osaksi siitä, että Jyväskylän kaupungin omistamien kiinteistöjen osalta marraskuun laskutus tehtiin hintatariffin vaihtumisen seurauksena joulukuun laskun yhteydessä, jolloin myös kulutustilastoon muodostui vastaava virhe. Tämän lisäksi joulukuun kulutuslukemaan sisältyi selittämätöntä virhettä sillä kuukauden kulutustilasto näytti kulutetun huomattavan paljon enemmän kuin kuukauden lämmönhankinta oli. Tämän vuoksi marras- ja joulukuu korjattiin lämmönhankintatilaston ja arvioidun lämpöhävikin perusteella. Marraskuun lämpöhäviö on arvioitu yleisesti noin seitsemäksi prosentiksi kokonaishankinnasta ja joulukuussa hävikki on prosenttiyksikköä pienempää. Näin marras- ja joulukuun korjatut kulutuslukemat saatiin vähentämällä kuukauden lämmönhankinnasta arvioidun hävikin osuus, eli saatiin:

marraskuu: $0.93 \times 80.641 \text{ GWh} = 74.996 \text{ GWh}$

joulukuu: $0.94 \times 119.370 \text{ GWh} = 112.208 \text{ GWh}$.

Sama korjaustoimenpide suoritettiin joulukuun havainnolle vuosien 1993 ja 1995 osalta. Kummassakin tapauksessa havaintoarvo alkuperäisessä aineistossa poikkesi huomattavasti lämmönhankintatilastosta, joka kerätään lämmön tuotantolaitoksista kalenterikuukausittain. Lämmönhankintatilasto on etenkin kalenterikuukausittain tarkasteltuna huomattavasti luotettavampi kuin asiakkaiden mittarilukemailoituksiin perustuva kulutustilasto. Siksi lämmönhankintatilastosta olennaisesti poikkeavat havainnot on perusteltua korjata edellä kuvatulla tavalla. Liitteessä 1 on esitetty alkuperäinen aineisto, korjattu aineisto sekä aineisto, josta on poistettu Säynätsalon osuus. Lisäksi liitteessä esitetään lämmönhankintatilasto sekä astepäivälukutilasto niin normaalivuoden kuin havaintujen astepäivälukujenkin osalta.

10. Analyysitulokset

Tässä luvussa tarkastellaan erilaisten ennustemallien toimivuutta laatimalla ennusteet vuoden 1996 kulutukselle ja pohditaan eri mallien ominaisuuksia. Lopuksi toimitaan realistisessa ennustamistilanteessa ja laaditaan vuoden 1997 kulutukselle ennuste soveltaen erilaisia skenaarioita vuoden 1997 astepäiväluvulle.

10.1 Kaukolämmön kokonaiskulutusennuste vuodelle 1996

Seuraavassa muodostetaan erilaisia kaukolämmön kulutusta kuvaavia malleja ja arvioidaan mallien toimivuutta ennustamistarkoituksessa. Kaikilla esitetyillä malleilla laaditaan kuukausitason ennuste vuodelle 1996, jota sitten voidaan verrata niin toteutuneeseen kulutukseen kuin muiden mallien tuottamiin ennusteisiin.

10.1.1 Yksinkertainen perusrakennemalli

Aineistoon sovitetaan seuraavassa perusrakennemallia, jolloin saadaan kokonaiskulutusta kuvaavaksi malliksi **malli 1a**:

$$\text{kulutus (GWh)} = \text{taso}^d + \text{kulmakerroin}^d + \text{kausiv.}^d + \text{AR(1)-jäännös}, \quad (10.1)$$

jossa d kertoo termin olevan luonteeltaan deterministinen. AR(1)-tekijä on muotoa

$$v_t = 0.289 \times v_{t-1} + \hat{\xi}_t, \quad \text{var}(\hat{\xi}_t) = 55.14. \quad (10.2)$$

Kausivaihtelukomponentti on rakennettu dummy-muuttujatekniikalla. Mallin diagnostisten testien tuloksia on esitetty taulukossa 10.1.

Taulukko 10.1 Mallin 1a diagnostisia tuloksia (p-arvot sulkeissa)

	malli 1a	
R_s^2	0.31	
p.e.v.	44.27	
N_{DH}	20.70	(.000)
$Q(8,7)$	3.75	(.808)
$H(27)$	2.03	(.036)
AIC	4.12	
D-W	1.81	
Cusum $t(70)$	0.11	(.455)
Chow $F(12,70)$	1.67	(.093)

Selitysaste 0.31 on melko heikko ja ennustevirheen varianssi (p.e.v.) on varsin suuri. Jäännösten normaalisuutta testataan Bowmanin ja Shentonin testin sijasta Durbinin ja Hansenin testillä (N_{DH}), joka soveltuu paremmin pienille aineistolle. Jäännökset eivät testin mukaan noudata normaalijakaumaa. Syynä tähän ovat eräät sääoloiltaan poikkeavat kuukaudet, mm. helmikuu

1994 ja joulukuu 1995, jotka kumpikin olivat poikkeuksellisen kylmiä (ks. kuvio 9.2). Näiden poikkeavien havaintojen mukanaolo osaltaan heikentää myös selitystasetta.

Jäännökset ovat asianmukaisesti korreloimattomia, sillä Boxin ja Ljungin Q-testi jäännösten autokorreloituneisuudelle antaa p-arvon 0.808. Jäännösten 1. asteen autokorreloituneisuutta testaava Durbinin ja Watsonin testisuure poikkeaa vain hivenen optimaalisesta arvostaan (= 2.00). Tältä osin malli on kunnossa. Akaiken informaatiokriteerin arvoksi saadaan 4.12. Ennustekykyä mittaavien testien (Cusum-t ja Chow'n F) mukaan malli on riittävän ennustekykyinen. Jäännösvarianssi ei H(27)-testisuureen mukaan ole vakio, vaan varianssi kasvaa tarkasteluajanjakson loppupuolella. Malli ei siis ole tältä osin riittävän hyvä, koska jäännökset eivät noudata normaalijakaumaa ja ovat heteroskedastisia.

Jäännösten ongelmaa voidaan yrittää ratkaista hyödyntämällä hetkellisten interventioiden tekniikkaa (luku 8.4). Tällöin malliin lisätään sääoloiltaan poikkeuksellisten kuukausien yhteyteen dummy-muuttujat. Mallitetaan kaksi poikkeavilta vaikuttavaa havaintoa dummy-muuttujien avulla. Havainnot ovat vuoden 1992 joulukuu ja vuoden 1994 helmikuu. Näistä ensin mainittu oli poikkeuksellisen lauha, astepäiväluku vain 563, kun pitkän aikavälin keskiarvo on 751. 1994 helmikuu oli puolestaan poikkeuksellisen kylmä astepäiväluvun ollessa 928, kun normaalivuoden lukema on 747. Mallittamalla nämä kaksi havaintoa kumpikin omalla dummy-muuttujalla saadaan **malli 1b**,

$$\text{kulutus (GWh)} = \text{taso}^d + \text{kulmakerroin}^d + \text{kausiv.}^d + \text{interventiot} + \text{AR(1)-jäännös}, \quad (10.3)$$

jossa AR(1)-jäännöstermi on muotoa

$$v_t = 0.261 \times v_{t-1} + \hat{\xi}_t, \quad \text{var}(\hat{\xi}_t) = 39.27. \quad (10.4)$$

Lisäyksen merkitystä mallin diagnostisille ominaisuuksille voidaan arvioida seuraavassa olevan taulukon 10.2 avulla.

Taulukko 10.2 Mallien 1a ja 1b diagnostisia tuloksia (p-arvot sulkeissa)

	malli 1a		malli 1b	
R^2_s	0.31		0.50	
p.e.v.	44.27		32.53	
N_{DH}	20.70	(.000)	13.26	(.001)
$Q(8,7)^1$ $Q(8,6)^2$	3.75 ¹	(.808)	4.72 ²	(.580)
H(27)	2.03	(.036)	1.12	(.385)
AIC	4.12		3.89	
D-W	1.81		1.82	
Cusum t(70)	0.11	(.455)	0.43	(.335)
Chow F(12,70)	1.67	(.093)	2.02	(.035)

Kuten nähdään mallin 1b selitystaseta on huomattavasti parempi kuin mallissa 1a. Myös ennustevirheen varianssi on pienentynyt huomattavasti. Jäännösten normalisuus on parantunut myös melko paljon, testisuureen arvo on pienentynyt noin kolmasosalla. Tämä ei kuitenkaan riitä normalisoimaan jäännöksiä. Autokorreloituneisuuden osalta tilanne säilyy ongelmattomana, niin Q-testin kuin Durbinin ja Watsonin testisuureenkin avulla mitattuna.

Jäännösvarianssi on oletusten mukaisesti homogeeninen mallissa 1b, kun mallissa 1a saatiin heteroskedastisuudesta kertova tulos. Parannus johtuu siitä, että dummy-muuttujilla mallitetut havainnot pienentävät jäännösvarianssia havaintojoukon loppupuolelta. Akaiken informaatiokriteeri paranee niin ikään. Sen sijaan ennustekykyä mittaavien testisuureiden arvot muuttuvat huonompaan suuntaan, eikä malli ole Chow'n testisuureen mukaan enää riittävän ennustuskyykyinen. Kaiken kaikkiaan interventioiden lisäämisellä on tässä tapauksessa merkittävä vaikutus mallin diagnostisiin ominaisuuksiin. Taulukossa 10.3 nähdään viimeisen estimoinnissa käytetyn havainnon (*final state*) mukaiset estimaatit.

Taulukko 10.3 Mallien 1a ja 1b estimointitulokset.

malli 1a

termi	estimaatti	r.m.s.e.	t-arvo	p-arvo
taso	60.974	2.179	27.981	.000
kulmak.	0.223	0.045	4.976	.000
AR(1)	22.164*	3.301	6.714	.000
kausi-1	35.094	2.627	13.361	.000
kausi-2	26.768	2.628	10.185	.000
kausi-3	21.311	2.624	8.124	.000
kausi-4	2.742	2.621	1.046	.298
kausi-5	-18.373	2.619	-7.016	.000
kausi-6	-33.852	2.618	-12.931	.000
kausi-7	-38.183	2.618	-14.585	.000
kausi-8	-35.097	2.619	-13.402	.000
kausi-9	-18.255	2.621	-6.967	.000
kausi-10	2.836	2.624	1.081	.283
kausi-11	22.592	2.628	8.596	.000
kausi-12	32.417			

malli 1b

termi	estimaatti	r.m.s.e.	t-arvo	p-arvo
taso	60.148	1.841	32.667	.000
kulmak.	0.209	0.038	5.485	.000
AR(1)	20.353*	2.976	6.839	.000
kausi-1	35.007	2.281	15.346	.000
kausi-2	22.751	2.441	9.320	.000
kausi-3	21.386	2.277	9.394	.000
kausi-4	2.849	2.274	1.253	.214
kausi-5	-18.244	2.273	-8.027	.000
kausi-6	-33.705	2.272	-14.833	.000
kausi-7	-38.019	2.273	-16.730	.000
kausi-8	-34.917	2.273	-15.360	.000
kausi-9	-18.056	2.275	-7.937	.000
kausi-10	3.056	2.277	1.342	.183
kausi-11	22.838	2.281	10.012	.000
kausi-12	35.054			
12/1992	-16.582	6.551	-2.531	.013
2/1994	28.246	6.568	4.301	.000

* AR-komponentin arvo hetkellä 12/1995.

Kuten nähdään dummy-muuttujat ovat selkeästi tilastollisesti merkitseviä. Joulukuun 1992 lauhan sään seurauksena kulutus laski vajaat 17 gigawattituntia joulukuun normaalitasosta. Vastaavasti helmikuun 1994 poikkeuksellisen kylmät olosuhteet lisäsivät kulutusta lähes 30 gigawattituntia. Dummy-muuttujien lisääminen malliin vaikuttaa ennen kaikkea helmi- ja joulukuun kausivaihtelukomponentteihin. Mallissa 1b joulukuun komponentin arvo on nousut ja vastaavasti helmikuun arvo on laskenut. Tämä johtuu luonnollisesti siitä, että dummy-muuttujilla mallitettujen havaintojen painoarvo mallin muissa komponenteissa vähenee.

Mallien 1a ja 1b avulla saadaan vuodelle 1996 taulukon 10.4 mukaiset ennusteet ja STAMP-ohjelmiston oletusarvoina tuottamat 68%:n ennustevälit, sekä ennustevirhe, joka kertoo prosentteina ennusteen poikkeaman toteutuneesta kulutuksesta.

Taulukko 10.4 Perusrakennemallien 1a ja 1b tuottamat ennusteet vuodelle 1996.

malli 1a

	toteutunut	ennuste	r.m.s.e.	68% - ennusteväli		virhe-%
tammikuu	112.191	102.690	7.374	95.313	110.060	-8.47
helmikuu	120.835	90.037	7.775	82.262	97.811	-25.49
maaliskuu	95.684	83.493	7.844	75.65	91.337	-12.74
huhtikuu	69.221	64.772	7.867	56.905	72.640	-6.43
toukokuu	44.651	43.773	7.883	35.891	51.656	-1.97
kesäkuu	23.380	28.489	7.895	20.593	36.384	21.85
heinäkuu	21.822	24.374	7.908	16.466	32.282	11.69
elokuu	17.972	27.682	7.921	19.762	35.603	54.03
syyskuu	44.627	44.750	7.933	36.817	52.683	0.28
lokakuu	61.408	66.066	7.946	58.120	74.012	7.59
marraskuu	74.130	86.048	7.959	78.088	94.007	16.08
joulukuu	111.491	96.098	7.973	88.125	104.070	-13.81
yht.	797.412	758.272		663.992	852.545	-4.91

malli 1b

	toteutunut	ennuste	r.m.s.e.	68% - ennusteväli		virhe-%
tammikuu	112.191	100.680	6.412	94.272	107.100	-10.26
helmikuu	120.835	84.708	6.708	78.000	91.415	-29.90
maaliskuu	95.684	82.526	6.755	75.771	89.280	-13.75
huhtikuu	69.221	63.930	6.772	57.159	70.702	-7.64
toukokuu	44.651	42.976	6.783	36.193	49.759	-3.75
kesäkuu	23.380	27.706	6.794	20.913	34.500	18.50
heinäkuu	21.822	23.597	6.804	16.794	30.401	8.13
elokuu	17.972	26.908	6.814	20.094	33.722	49.72
syyskuu	44.627	43.978	6.824	37.153	50.802	-1.45
lokakuu	61.408	65.299	6.835	58.464	72.134	6.34
marraskuu	74.130	85.290	6.846	78.444	92.136	15.05
joulukuu	111.491	97.715	6.857	90.858	104.570	-12.36
yht.	797.412	745.313		664.115	826.521	-6.53

Kuten taulukosta nähdään, diagnostiikaltaan selvästi parempi malli 1b tuottaa parempia ennusteita alkuvuotta lukuunottamatta. Myös suppeampi ennustevirheen varianssi ja parempi selitysaste kielivät siitä, että mallin 1b tuottamat ennusteet osoittautuvat pitkällä tähtäimellä luotettavammiksi. Vuositasolle aggregoituna malli 1a tuottaa kuitenkin tässä tapauksessa paremman ennusteen. Tämä johtuu kuitenkin puhtaasti sattumasta. Molempien mallien osalta helmikuun ennuste on sangen virheellinen. Tämä johtuu siitä, että helmikuu 1996 oli poikkeuksellisen kylmä, kiinteistöjen lämmitystarvetta kuvaava astepäiväluku sai tällöin arvon 880, kun vuosien 1961 - 1990 keskiarvoon perustuvan normaalivuoden astepäiväluku helmikuulle on vain 747.

Mallien tuottamat ennusteet näyttävät jäävän talvikuukausina toteutuneen tason alapuolelle ja kesäkuukausina malli yliestimoi kulutusta. Ennusteet ovat kuukausitasolla arvioituna valtaosin heikohkoja, sillä vain viidessä tapauksessa kahdestatoista piste-ennusteen virhe jää alle kymmenen prosenttiyksikön. Vuositasolla sen sijaan päästään melko hyvään tulokseen virheen ollessa alle viisi prosenttia. Ennustevalit ovat sangen leveät ja toteutunut lämmönkulutus sijoittuu ennustevalille kahdessa tapauksessa kolmesta. Jatkossa ennustevaliä saadaan supistettua hyödyntämällä lisäinformaatiota (*astepäiväluku*). Myös piste-ennusteiden voidaan olettaa tarkentuvan.

10.1.2 Regressiomalli

Laajennetaan tarkastelua nyt siten, että hyödynnetään lisäinformaatiota selittävän muuttujan (*astepäiväluku*, *apl*) muodossa. Koska seuraavassa on tarkoituksena vertailla erilaisia malleja, hyödynnetään tietoa vallinneista astepäiväluvuista myös ennustusperiodin aikana, vaikka realistisessa ennustetilanteessa ennustettavan periodin täsmällisiä astepäivälukuja ei tietenkään olisi saatavilla.

Mallin rakentaminen voidaan aloittaa muodostamalla yksinkertainen klassinen regressiomalli, joka sisältää siis kiinteän tasokomponentin, selittävän muuttujan ja jäännöstekijän. Näin ollen **malli 2** on siis muotoa

$$\text{kulutus (GWh)} = \text{taso}^d + \beta \times \text{astepäiväluku} + \text{jäännös}, \quad (10.5)$$

Taulukossa 10.5 on esitetty mallin 2 (regressiomalli) diagnostisten testien tuloksia.

Taulukko 10.5 Mallin 2 diagnostisten testien tuloksia (p-arvot sulkeissa).

	malli 2	
R^2_L	0.97	
p.e.v.	26.58	
N_{DH}	9.52	(.009)
$Q(8,8)$	34.47	(.000)
$H(27)$	6.49	(.000)
AIC	3.33	
D-W	0.77	
Cusum t(71)	5.09	(.000)
Chow F(12,71)	4.50	(.000)

Kuten nähdään mallin selityssaste on erittäin hyvä, joskaan selityssastetta ei tulisi verrata mallin 1 selityssasteeseen, sillä mallissa 1 selityssaste laskettiin eri tavalla kuin tässä tapauksessa. Syytä tähän on se, että kuten yhtälöistä 5.17 - 5.19 nähdään, selityssasteella on erilainen laskenta-kaava riippuen siitä, millaisia komponentteja malli sisältää. Myös ennustevirheen varianssi on pienentynyt verrattuna malleihin 1a ja 1b. Tämä osaltaan kertoo siitä, että ennusteiden ennustevalit ovat sängen pienet.

Jäännökset ovat normaalisti jakautuneita, mutta aikasarjaregressiomallille tyypillisesti selvästi autokorreloituneita ($D-W = 0.77$ ja $Q(8,8) = 34.47$). Durbinin ja Watsonin testisuureen arvo viittaisi selvästi AR(1)-rakenteen olemassaoloon. Koska tässä yhteydessä on tarkoitus kuitenkin tuottaa vain klassisen regressiomallin mukaiset ennusteet, ei mallin täydentämistä AR-komponentilla suoriteta.

Mallin jäännökset ovat myös heteroskedastisia ($H(27) = 6.49$). Myöskään Chow'n testisuure ja kumulatiivisen summan t-testi antavat tilastollisesti merkitsevän tuloksen, joten malli ei tässäkään mielessä ole riittävän hyvä. Regressiomallille saadaan yhtälön 10.6 mukaiset estimointitulokset:

$$\text{kulutus} = 10.457 + 0.105 \times \text{apl} + \hat{\varepsilon}, \quad \text{var}(\hat{\varepsilon}) = 27.17, \quad R^2 = 0.97. \quad (10.6)$$

(10.167) (47.949)

Estimoitujen komponenttien t-arvot on esitetty suluissa kyseisen komponentin alla. Kumpikin komponentti on tilastollisesti erittäin merkitsevä. Nähdään, että astepäiväluvussa yhden asteen lisäys lisää kaukolämmön kulutusta noin 105 MWh ja kiinteä tasokomponentti on noin 10.5 GWh. Koska kesällä lämmitystarvetta ei esiinny, tasokomponentti voidaan tulkita käyttöveden lämmityksestä aiheutuvan kulutuksen komponentiksi. Diagnostisten tarkastelujen nojalla mallia ei voida missään tapauksessa pitää riittävän hyvänä, mutta vertailun vuoksi konstruoidaan mallilla kuitenkin ennusteet vuodelle 1996.

Taulukko 10.7 Mallin 2 (regressiomalli) tuottamat ennusteet vuodelle 1996.

	toteutunut	ennuste	r.m.s.e.	68% - ennustevali		virhe-%
tammikuu	112.191	97.964	5.243	92.722	103.210	-12.68
helmikuu	120.835	88.649	5.243	83.406	93.891	-26.64
maaliskuu	95.684	80.798	5.243	75.555	86.041	-15.56
huhtikuu	69.221	59.549	5.243	54.306	64.792	-13.97
toukokuu	44.651	33.067	5.243	27.824	38.310	-25.94
kesäkuu	23.38	14.749	5.243	9.506	19.992	-36.92
heinäkuu	21.822	11.923	5.243	6.680	17.166	-45.36
elokuu	17.972	17.261	5.243	12.018	22.504	-3.96
syyskuu	44.627	36.416	5.243	31.174	41.659	-18.40
lokakuu	61.408	54.316	5.243	49.073	59.559	-11.55
marraskuu	74.13	70.645	5.243	65.402	75.888	-4.70
joulukuu	111.491	89.067	5.243	83.824	94.310	-20.11
yht.	797.412	654.404		591.490	717.322	-17.93

Ennusteet ovat keskimäärin huonompia kuin malleilla 1a ja 1b. Kuukausitasolla tarkasteltuna ennusteiden virhe vaihtelee neljästä prosenttiyksiköstä neljänkymmeneen viiteen. Mallin tuottamat ennusteet aliestimoivat kulutusta systemaattisesti, eli asteapäiväluku ei yksin kykene selittämään kulutuksen kuukausittaisia vaihteluja. Kokeillaan siis mallin laajentamista.

11.1.3 Laajennettu perusrakennemalli

Seuraavaksi yhdistetään perusrakennemallin antama informaatio perinteisen regressiomallin informaatioon. Käytännössä siis yhdistetään edellä laaditut mallit. Voidaan tulkita, että selittävä muuttuja selittää kaukolämmön kulutuksesta säätilan vaihtelun aiheuttamat lämmitystarpeen vaihtelut ja perusrakennemalli vastaa muista lämmönkulutukseen vaikuttavista tekijöistä. Näin saadaan tuloksena **malli 3**, joka on yhtälön 10.7 mukaista muotoa.

$$\text{kulutus (GWh)} = \text{taso}^s + \text{kulmak.}^d + \text{kausiv.}^s + \beta \times \text{apl} + \text{AR(1)-jäännös}, \quad (10.7)$$

(0.08) (0.45) (7.96)

Stokastisten komponenttien varianssiestimaatit on esitetty vastaavien komponenttien alla. Muilta osin estimointitulokset on esitetty taulukossa 10.9. Yhtälössä s merkitsee stokastista ja d determinististä termiä. Kausivaihtelu on rakennettu dummy-muuttujien avulla. AR(1)-jäännöstermi on muotoa

$$v_t = 0.358 \times v_{t-1} + \hat{\xi}_t, \quad \text{var}(\hat{\xi}_t) = 7.96. \quad (10.8)$$

Mallin diagnostisten testien tulokset esitetään seuraavassa taulukossa.

Taulukko 10.8 Mallin 3 diagnostisten testien tulokset (p-arvot sulkeissa).

	malli 3	
R^2_s	0.84	
p.e.v.	10.03	
N_{DH}	3.81	(.149)
Q(9,6)	16.84	(.001)
H(23)	2.57	(.014)
AIC	2.71	
D-W	1.94	
Cusum t(59)	1.68	(.095)
Chow F(12,59)	0.97	(.168)

Selitysaste 0.84 on hyvä ja ennustevirheen varianssi (p.e.v.) on sangen pieni verrattuna aikaisempiin malleihin. Ensimmäisen asteen autokorreloituneisuutta mittaava Durbinin ja Watsonin testisuure on myös kiitettävän lähellä odotusarvoaan, joka on kaksi. Ongelmana mallissa on jäännösten autokorreloituneisuus, joka liittyy selittävän muuttujan mukaanottoon, sillä mallissa 1 ongelmaa ei ollut, mutta jo regressiomallissa ongelma esiintyi. Jäännösten autokorreloituneisuutta mittaavan Boxin ja Ljungin Q-testin mukaan jäännökset ovat selkeästi autokorreloituneita (p = .001). Jäännösten autokorrelaatioiden yksityiskohtaisemmassa tarkastelussa havaitaan, että suurimmat autokorrelaatiot ovat viiveillä 5, 7, 8, 12, 13 ja 17. Korrelaa-

tiot ovat itseisarvoltaan suurempia kuin 0.15. Myös jäännösten varianssin vakioisuutta mittaava H(h)-testisuure antaa tilastollisesti merkitsevän tuloksen ($p = .014$), eli jäännösvarianssiin ei vaikuta vakiolta. Mallille saadaan taulukossa 10.9 nähtävät estimointitulokset, joissa AR(1)-termi vastaa AR-komponentin arvoa joulukuussa 1995. Komponentin arvo pienenee aikaa myöten yhtälössä 10.8 nähtävän kertoimen mukaisesti.

Taulukko 10.9 Estimointitulokset mallille 3.

termi	estimaatti	r.m.s.e.	t-arvo	p-arvo
taso	25.672	2.323	11.053	.000
kulmak.	0.173	0.039	4.450	.000
$\hat{\beta}$	0.085	0.005	18.240	.000
AR(1)	5.999*	1.924	3.118	.002
kausi-1	6.736			
kausi-2	5.353	1.900	2.818	.006
kausi-3	4.874	1.682	2.899	.005
kausi-4	-1.874	1.433	-1.308	.195
kausi-5	-4.193	1.602	-2.618	.011
kausi-6	-2.840	2.217	-1.281	.204
kausi-7	-6.485	2.241	-2.895	.005
kausi-8	-7.904	2.077	-3.806	.000
kausi-9	-5.905	1.551	-3.808	.000
kausi-10	-1.387	1.424	-0.974	.333
kausi-11	3.855	1.739	2.216	.030
kausi-12	9.771	2.027	4.820	.000

* AR-komponentin arvo hetkellä 12/1995.

Kuten nähdään mallin termit ovat kolmea kausivaihtelukomponenttia lukuunottamatta tilastollisesti merkitseviä. Kausivaihtelukomponenteista ensimmäinen on saatu yhtälön 2.14 mukaisen summaussäännön perusteella, jolloin komponenttien summa on nolla. Komponenttien avulla voidaan konstruoida ennusteet tuleville ajanjaksoille. Taulukon 10.9 estimaatit vastaavat siis komponenttien arvoja viimeisen havainnon hetkellä eli joulukuussa 1995. Näin ollen ensimmäiseksi saadaan ennuste tammikuulle 1996. Tämä tapahtuu sijoittamalla edellä esitetyt komponentit yhtälöihin 10.7 ja 10.8. Ennusteeksi saadaan

$$99.462 = 25.672 + 0.172634 + 6.736 + (0.0852865 \times 759) + (0.3583 \times 5.9985), \quad (10.9)$$

$$\text{ennuste } 1/96 = \text{taso} + \text{kulmak} + \text{kausiv.} + \hat{\beta} \times \text{apl} + \text{AR(1)-jäännös.}$$

Yhtälössä (10.9) on käytetty STAMP -ohjelmiston tuottamia estimaattien tarkkoja arvoja taulukossa (10.9) esitettyjen pyöristettyjen arvojen sijaan. Vastaavalla tavalla saadaan ennusteet seuraaville kuukausille sijoittamalla mekaanisesti oikeat komponentit yhtälöön. Huomiotava on vain se, että kulmakerroin komponentti lisätään jokaisella ajanhetkellä aina uudestaan, esimerkiksi hetken T+ k kulmakerroin komponentti olisi tässä tapauksessa $k \times 0.172634$. Vastaavasti AR(1)-komponentti samalla hetkellä olisi $0.3583^k \times 5.9985$. Astepäiväluvun osalta tammikuulle havaitun astepäiväluvun (759) tilalle sijoitettaisiin hetkeä T + k koskeva astepäiväluku.

Astepäiväluvun kasvu yhdellä asteella merkitsee mallin mukaan 85 MWh:n kasvua kulutuksessa. Verrattuna malliin 1 kausivaihtelukijän estimaatit ovat selvästi pienentyneet, koska nyt kausivaihtelukomponenttiin sisältyy vain osa kausivaihtelusta. Tämä osa voidaan tulkita käyttöveden lämmityksen kausivaihteluksi, kun astepäiväluku vastaa varsinaisesta kiinteistöjen lämmittämisen kausivaihtelusta. Mallin avulla saadaan taulukon 10.10 mukaiset ennusteet.

Taulukko 10.10 Mallin 3 (laajennettu perusrakennemalli) tuottamat ennusteet vuodelle 1996.

	toteutunut	ennuste	r.m.s.e.	68% - ennusteväli		virhe-%
tammikuu	112.191	99.462	3.529	95.933	102.990	-11.35
helmikuu	120.835	107.190	3.685	103.510	110.880	-11.29
maaliskuu	95.684	89.078	3.739	85.340	92.817	-6.90
huhtikuu	69.221	64.160	3.763	60.398	67.923	-7.31
toukokuu	44.651	43.273	3.779	39.494	47.051	-3.09
kesäkuu	23.380	27.206	3.792	23.414	30.998	16.36
heinäkuu	21.822	23.725	3.805	19.921	27.530	8.72
elokuu	17.972	19.150	3.817	15.333	22.968	6.55
syyskuu	44.627	45.969	3.830	42.139	49.798	3.01
lokakuu	61.408	58.079	3.843	54.235	61.922	-5.42
marraskuu	74.130	72.363	3.858	68.505	76.221	-2.38
joulukuu	111.491	104.980	3.858	101.120	108.830	-5.84
yht.	797.412	754.635		709.342	799.928	-5.36

Mallin ennustekyky näyttää melko hyvältä, etenkin loppuvuoden ennustevirheet ovat kiitettävän pieniä. Alkuvuoden osalta silmiinpistävää on tammi- ja helmikuun suurehkot ennustevirheet, joista jälkimmäiseen on mitä ilmeisimmin syynä jo aiemmin todettu poikkeuksellisen kylmä sää. Estimoidulle 68%:n ennustevälille toteutuneista havainnoista osuu seitsemän. Tässäkin suhteessa ongelmallisia kuukausia ovat vuoden alkukuukaudet sekä näiden ohella joulukuu. Malli näyttäisi hivenero aliarvioivan kylmempien kuukausien kulutusta ja vastaavasti painottavan liikaa kesän kulutusta. Tämä ongelma saattaa johtua siitä, että astepäiväluvun todellinen beta-kerroin ei ole vuoden kaikille kuukausille sama. Mikäli astepäiväluvun beta-kerroin kasvaa astepäiväluvun kasvaessa päädytään juuri edellä havaittuun tilanteeseen. Ratkaisuksi voitaisiin tulevaisuudessa koettaa sitä, että mallitetaan astepäiväluvun vaikutus käyttäen useampaa muuttujaa. Voitaisiin jakaa astepäiväluvun antama informaatio esimerkiksi kahteen muuttujaan, joista toinen sisältäisi pakkaskauden ja toinen lämpimän ajanjakson informaation, tai voidaan jopa tehdä muuttuja vuoden jokaiselle kuukaudelle erikseen, jolloin saadaan astepäiväluvun beta-kertoimelle oma estimaatti kullekin kuukaudelle.

11.1.4 Mallien vertailua

Edellä on siis muodostettu neljä erilaista mallia, joista ensimmäinen, 1a on yksinkertainen perusrakennemalli. Malli 1b on edellinen laajennettuna kahdella interventiomuuttujalla, joilla mallitetaan kahta sääoloiltaan poikkeuksellista kuukautta (12 / 1992 ja 2 / 1994). Malleista kolmas on perinteinen regressiomalli ja neljäs on perusrakennemalli, jossa hyödynnetään kiinteistöjen lämmitystarvetta kuvaavaa astepäivälukusarjaa selittävänä muuttujana. Taulu-

koista 10.11 ja 10.12 nähdään vuoden 1996 toteutunut kaukolämmön kulutus, edellä tarkasteltujen neljän mallin tuottamat ennusteet vuodelle 1996 sekä ennusteiden virheprosentit.

Taulukko 10.11 Mallien 1 - 3 tuottamat ennusteet vuodelle 1996.

	toteutunut	malli 1a	malli 1b	malli 2	malli 3
tammikuu	112.191	102.690	100.680	97.964	99.318
helmikuu	120.835	90.037	84.708	88.649	107.050
maaliskuu	95.684	83.493	82.526	80.798	88.913
huhtikuu	69.221	64.772	63.930	59.549	63.996
toukokuu	44.651	43.773	42.976	33.067	43.101
kesäkuu	23.380	28.489	27.706	14.749	27.019
heinäkuu	21.822	24.374	23.597	11.923	23.550
elokuu	17.972	27.682	26.908	17.261	18.974
syyskuu	44.627	44.750	43.978	36.416	45.786
lokakuu	61.408	66.066	65.299	54.316	57.891
marraskuu	74.130	86.048	85.290	70.645	72.189
joulukuu	111.491	96.098	97.715	89.067	104.800
yht.	797.412	758.272	745.313	654.404	752.587

Taulukko 10.12 Mallien 1 - 3 tuottamien ennusteiden virheprosentit.

	toteutunut	malli 1a	malli 1b	malli 2	malli 3
		virhe-%	virhe-%	virhe-%	virhe-%
tammikuu	112.191	-8.47	-10.26	-12.68	-11.35
helmikuu	120.835	-25.49	-29.90	-26.64	-11.29
maaliskuu	95.684	-12.74	-13.75	-15.56	-6.90
huhtikuu	69.221	-6.43	-7.64	-13.97	-7.31
toukokuu	44.651	-1.97	-3.75	-25.94	-3.09
kesäkuu	23.380	21.85	18.50	-36.92	16.36
heinäkuu	21.822	11.69	8.13	-45.36	8.72
elokuu	17.972	54.03	49.72	-3.96	6.55
syyskuu	44.627	0.28	-1.45	-18.40	3.01
lokakuu	61.408	7.59	6.34	-11.55	-5.42
marraskuu	74.130	16.08	15.05	-4.70	-2.38
joulukuu	111.491	-13.81	-12.36	-20.11	-5.84
yht.	797.412	-4.91	-6.53	-17.93	-5.36

Kuten taulukosta nähdään pelkän regressiomallin ennustekyky on ennustevirheellä mitattuna muita malleja huonompi. Myös yksinkertaisen rakenneyhtälömallin 1a, sekä mallin 1b ennustevirheet vaihtelevat liian paljon. Tässä tarkasteltavien kahdentoista kuukauden ennusteiden ennustevirheiden otosvarianssi on perusrakennemallilla 1a 438.00 ja perusrakennemallilla 1b 412.19. Regressiomallilla ennustevirheen otosvarianssi on 152.94 ja mallilla 3 varianssi on 72.89. Malli 3 on siten ennustekyvyltään malleista selvästi paras, vaikka ennusteet vuositason la summattuna yksinkertainen perusrakennemalli 1a tuottaakin niukasti paremman tuloksen. Regressiomalli osoittautuu tässä vertailussa yllättäen malleja 1a ja 1b paremmaksi.

10.2 Kulutusennusteet vuodelle 1997

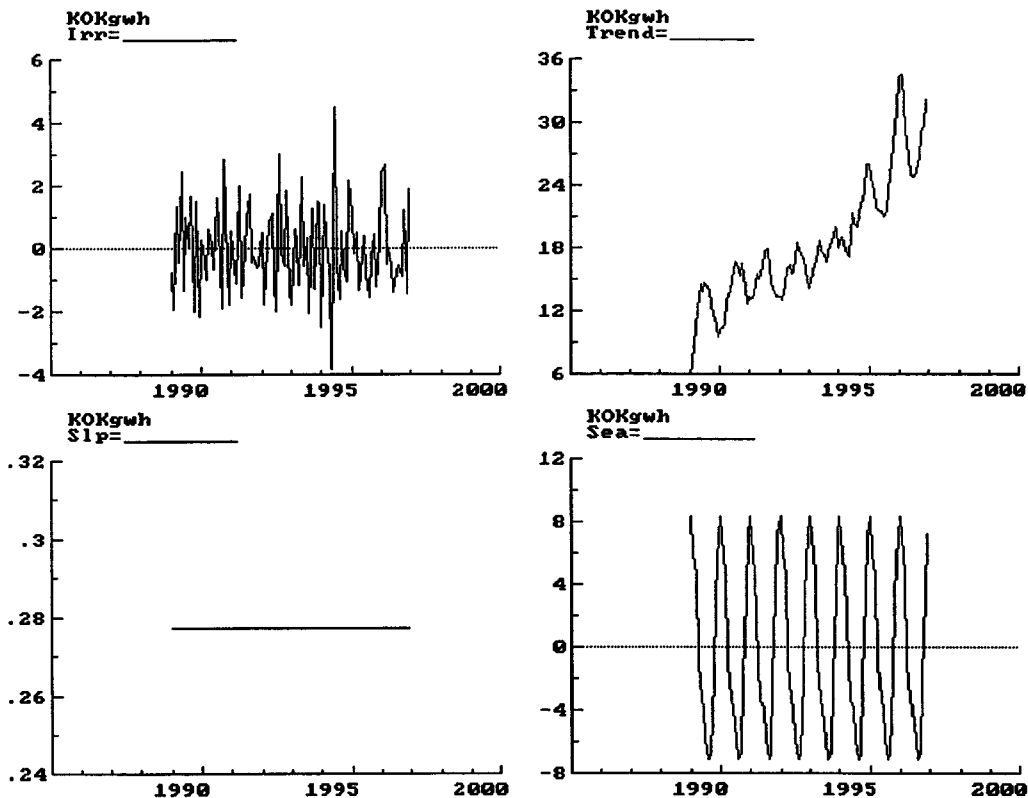
Seuraavaksi laaditaan ennusteet kaukolämmön kuukausikulutukselle vuonna 1997. Ennusteita laaditaan kolmelle erilaiselle säätilaskenaariolle: normaalivuosi, lämmin vuosi ja kylmä vuosi. Lämmitettävän rakennustilavuuden oletetaan kasvavan kuten aiemminkin, joten tästä johtuva kulutuksen kasvu saadaan huomioitua kulmakerroinkomponentin avulla.

Vuoden 1997 kulutukselle muodostettava malli 4 on muotoa

$$\text{kulutus (GWh)} = \text{taso}^s + \text{kulmakerroin}^d + \text{kausiv.}^d + \beta \times \text{apl} + \text{jäännös.} \quad (10.10)$$

(6.26) (3.81)

Kaikkien komponenttien estimointitulokset nähdään taulukossa 10.14 ja komponenttien varianssiestimaatit on esitetty komponenttien alla. Malli sisältää siis stokastisen tasokomponentin, deterministisen kulmakerroinkomponentin ja kausivaihtelutekijän, selittävän muuttujan ja jäännöstermin. Kausivaihtelu on dummy-muuttujaista muotoa (2.15). STAMP ohjelmistolla voidaan tuottaa graafinen esitys eri komponenteille. Kuviossa 10.1 nähdään jäännöskomponentin (*Irr*), trendikomponentin (*Trend*), kulmakerroinkomponentin (*Slp*) ja kausivaihtelukomponentin (*Sea*) graafiset esitykset mallin 4 osalta.



Kuvio 11. 1. Mallin 4 (10.10) komponenttien graafinen esitys

Kuten kuviosta nähdään, kiinteät komponentit (kulmakerroin, kausivaihtelu) pysyvät vakioina toisin kuin esimerkiksi stokastisen tason sisältävä trendikomponentti. Trendikomponentissa havaitaan voimakas muutos vuoden 1996 aikana. Tämä johtuu ainakin osittain poikkeuksellisen kylmästä helmikuusta. Mallin 4 diagnostisia ominaisuuksia mittaavien testien tulokset on esitetty taulukossa 10.13.

Taulukko 10.13 Mallin 4 diagnostisten testien tuloksia (p-arvo sulkeissa).

	malli 4	
R^2_s	0.85	
p.e.v.	10.90	
N_{DH}	2.62	(.269)
Q(8,7)	22.50	(.002)
H(31)	1.43	(.160)
AIC	2.70	
D-W	1.96	
Cusum t(59)	-0.16	(.564)
Chow F(12,71)	1.85	(.053)

Mallin selitysaste on edelleen hyvä ja ennustevirheen varianssi on pieni verrattuna edellä tarkasteltuihin malleihin. Mallin jäännökset ovat normaalisti jakautuneet ja homoskedastisia, mutta Boxin ja Ljungin Q-testin mukaan autokorrelaatioita. Voimakkaimmat autokorrelaatiot havaitaan viiveillä 5, 7, 8, 12, 13, 17, joilla kaikilla korrelaation itseisarvo on yli 0.15.

1. asteen autokorrelaatiota ei kuitenkaan ilmene, sillä Durbinin ja Watsonin testisuure saa arvon 1.96. Myös Akaiken informaatiokriteeri saa sangen pienen arvon. Ennustekykä mittaava kumulatiivisen summan t-testi saa p-arvon 0.564 ja Chow'n testi 0.053. Näin ollen Chow'n testin mukaan mallin ennustekyky olisi juuri ja juuri riittävä ja kumulatiivisen summan t-testin mukaan ennustekyky on hyvä. Saadaan taulukon 10.14 mukaiset estimointitulokset.

Taulukko 10.14 Mallin 4 estimointitulokset

termi	estimaatti	r.m.s.e.	t-arvo	p-arvo
taso	32.319	2.735	11.818	.000
kulmak.	0.277	0.237	1.168	.246
$\hat{\beta}$	0.088	0.005	19.523	.000
kausi-1	8.408	1.829	4.597	.000
kausi-2	5.963	1.689	3.531	.001
kausi-3	4.748	1.439	3.300	.001
kausi-4	-1.071	1.111	-0.964	.337
kausi-5	-3.332	1.332	-2.501	.014
kausi-6	-3.794	1.924	-1.971	.052
kausi-7	-5.829	2.017	-2.890	.005
kausi-8	-7.037	1.882	-3.740	.000
kausi-9	-6.774	1.260	-5.377	.000
kausi-10	-1.997	1.103	-1.810	.074
kausi-11	3.415	1.384	2.469	.015
kausi-12	7.299			

Mallin viimeinen kausivaihtelukomponentti on saatu kausivaihtelukijöiden summausperiaatteen (yhtälö 2.14) mukaisesti. suurin osa tekijöistä on tilastollisesti merkitseviä. Kulmakerroin ja huhtikuun kausivaihtelukijä (kausi-4) eivät selvästikään ole tilastollisesti merkitseviä, mutta kahden muun kausivaihtelukijän (kaudet 6 ja 10) p-arvot jäävät vain niukasti tilastollisesti merkitsevän rajasta (0.05). Kokonaisuudessa kausivaihtelu on kuitenkin selkeästi

merkitsevä komponentti. Kulmakertoimen sisällyttäminen malliin on sisällöllisten seikkojen nojalla perusteltua, vaikka p-arvo ei olekaan tilastollisesti merkitsevä. Lämmitettävä rakennustilavuus kasvaa koko ajan, joten kulmakerroin vastaa tämän piirteen mallittamisesta. Mallin 4 tuottamat kulutusennusteet vuodelle 1997 olettaen, että astepäiväluvut vastaavat normaalivuoden lukemia on esitetty taulukossa 11.14.

Taulukko 10.15 Mallin 4 tuottamat kulutusennusteet vuodelle 1997 olettaen normaalivuoden astepäivälukuja vastaava säätyyppi (myös Säynätsalon kulutus huomioitu).

	ennuste	r.m.s.e.	68% - ennusteväli	
tammikuu	114.670	3.587	111.080	118.250
helmikuu	104.660	4.285	100.370	108.940
maaliskuu	97.111	4.896	92.215	102.010
huhtikuu	73.681	5.448	68.233	79.129
toukokuu	49.404	5.960	43.444	55.364
kesäkuu	33.799	6.439	27.360	40.238
heinäkuu	29.662	6.893	22.768	36.555
elokuu	33.225	7.327	25.898	40.552
syyskuu	49.890	7.744	42.146	57.633
lokakuu	70.011	8.146	61.865	78.157
marraskuu	89.446	8.536	80.910	97.982
joulukuu	109.120	8.915	100.200	118.030
yht.	854.679		776.489	932.840

Normaalivuoden ohella vastaavat ennusteet 10% normaalivuotta viileämmälle ja lämpimämmälle säätyypille nähdään taulukossa 10.16.

Taulukko 10.16 Mallin 4 vuodelle 1997 tuottamat kulutusennusteet eri astepäivälukuskenarioilla (myös Säynätsalon kulutus huomioitu).

	"lämmin vuosi"	normaaliv. ennuste	"viileä vuosi"	r.m.s.e.
tammikuu	107.270	114.670	122.070	3.954
helmikuu	98.049	104.660	111.270	4.730
maaliskuu	91.207	97.111	103.010	5.422
huhtikuu	69.540	73.681	77.822	6.036
toukokuu	47.465	49.404	51.343	6.594
kesäkuu	33.447	33.799	34.152	7.108
heinäkuu	29.573	29.662	29.750	7.587
elokuu	32.696	33.225	33.842	8.038
syyskuu	47.687	49.890	52.093	8.466
lokakuu	65.958	70.011	73.712	8.874
marraskuu	84.336	89.446	94.557	9.266
joulukuu	102.510	109.120	115.720	9.606
yht.	809.738	854.679	899.341	

Kuten taulukosta 10.16 nähdään, kesäajan kulutus pysyy melko vakiona, koska kesällä ei lämmitystarvetta yleensä esiinny ja kaikki kulutus syntyy käyttöveden lämmityksestä. Kylmimpinä talvikuukausina 10% muutos astepäiväluvussa näkyy runsaan kuuden prosentin muutoksena kulutuksessa. Vuositasolle aggregoituna 10% muutos kuukausittaisessa astepäiväluvussa näkyy 5.25% muutoksena kulutuksessa. Tämä tulos perustuu oletukselle, että astepäiväluvun vaikutus kulutukseen on luonteeltaan lineaarista. On kuitenkin mahdollista, että lineaarisuusoletus ei todellisuudessa pidä paikkaansa, vaan astepäiväluvun painoarvo kasvaa lämmitystarpeen kasvaessa. Tämän ongelman ratkaiseminen onkin jatkotutkimuksen keskeisiä kysymyksiä.

Vuositasolla tarkasteltuna normaalivuoden astepäivälukuihin perustuvat ennusteet ennakoivat kulutukselle noin viiden ja puolen prosentin vuotuista kasvua. 1996 kokonaiskulutus Säynät-salo mukaanlukien oli noin 810 GWh ja mallin 4 vuodelle 1997 ennustama kulutus oli noin 855 GWh. Vuoden 1996 astepäivälukujen summa oli 5039, kun normaalivuoden summa on 5053, eli tältä osin vuosi 1996 vastasi melko tarkoin normaalivuotta. Normaalivuotta 10% kylmempi vuosi 1997 johtaa mallin mukaan yhdentoista prosentin kasvuun kulutuksessa ja vastaavasti 10% lämpimämmän vuoden kulutus pysyy edellisvuoden tasolla.

11. Yhteenveto

Edellä on siis esitelty rakenneyhtälömallien teoriakehikkoa perusrakennemallien osalta. Aikaisempaan aikasarja-analyttiseen ajatteluun verrattuna keskeisin muutos on mallin komponenttien stokastisointi. Tämän muutoksen myötä mahdollistuu klassisen testiteorian hyödyntäminen.

Teoriakehikon osalta keskeisimpiä tekijöitä luvussa 2 tapahtuneen mallin rakentamisen ohella on mallin estimoinnin mahdollistava Kalmanin suodin, jonka toimintaa tarkasteltiin luvussa 3. Luvussa 4 osoitettiin, kuinka Kalmanin suotimen avulla klassiseen estimointiteoriaan linkitetyn perusrakennemallin komponentit saadaan estimoitua. Tämän jälkeen luvussa 5 tarkasteltiin perusrakennemallien diagnostiikkaa. Luvuissa 6 ja 7 tarkasteltiin mallien erilaisia käyttö-tarkoituksia, joista keskeisimpänä on mallien käyttö ennustamiseen. Ennustamisen ohella perusrakennemalleja voidaan hyödyntää myös esimerkiksi analysoimalla mallin eri komponentteja erikseen. Luvussa 8 tehtiin perusrakennemalleihin selkeä laajennus ottamalla mukaan lisäinformaatio selittävien muuttujien ja interventiotermien muodossa. Tämän jälkeen sovellettiin perusrakennemalleja empiirisen ongelman ratkaisuun. Empiirisenä ongelmana tässä työssä oli Jyväskylän kaukolämmön kulutuksen ennustaminen.

Perusrakennemallien teoriakehikko on erittäin joustava ja näin ollen sitä voidaan laajentaa monin tavoin. Tässä työssä perusrakennemalleja on laajennettu selittävien muuttujien sekä interventioiden avulla. Myös muita laajennusmahdollisuuksia on, esimerkiksi taloustieteellisissä ongelmissa on hyödynnetty niin kutsuttua päivittäisvaikutusten komponenttia.

Kaukolämmön kulutuksen osalta laadittiin kulutusennusteet vuodelle 1997. Vuositasolla on mallin mukaan odotetavissa, että kulutus kasvaa runsaat 5%. Tämä edellyttää kuitenkin sitä, että säätila noudattelee Ilmatieteen laitoksen pitkän aikavälin keskiarvoja. Mikäli sää on tätä 10% lämpimämpää, pysyy kulutus jotakuinkin vuoden 1996 tasolla ja sään ollessa 10% ”normaalivuotta” viileämpää on odotettu kulutuksen kasvu noin yhdentoista prosenttiyksikön luokkaa. Jos tehdään erittäin rajoittavat oletukset, joiden mukaan lähinnä lämmitystilavuuden kasvusta johtuva kulutuksen kasvu jatkuu ennallaan ja säätila noudattaa ”normaalivuoden” tasoa, kaksinkertaistuu kulutus jo ennen vuotta 2010 mennessä, ja vuonna 2010 kulutus olisi runsaat 1700 GWh.

Kuukausitasolla tarkasteltuna astepäiväluku on keskeisin kulutusta selittävä tekijä, vaikka sekään yksin ei pysty kulutusta luotettavasti selittämään. Asteen muutos astepäiväluvussa merkitsee keskimäärin noin 86 MWh:n kasvua kulutuksessa. Tätä tietoa hyödyntäen on helppo laskea kulutusennusteita mille tahansa astepäivälukuina esitetyille lämpötilaoletuksille. Tulevaisuudessa kaukolämpöverkkoon liitettävän rakennuskannan oletetaan tässä yhteydessä kasvavan edellisten vuosien malliin. Tämän vuoksi mallissa on mukana kulmakerroinkomponentti, vaikka sen mukanaolo ei tilastolliselta kannalta katsottuna olisikaan perusteltua.

Mallien analysointiin ja estimointiin kehitetty STAMP 5.0 -ohjelmisto on sängen helppokäyttöinen, mutta joiltakin osin hivenen liian rajoittunut. Tästä huolimatta ohjelmisto on oivallinen työkalu rakenneyhtälömallien analysointiin.

Kuten luvusta 11 on nähtävissä, eivät empiiristä aineistoa kuvaamaan valitut mallit täytä kaikkia malleilta vaadittuja teoreettisia ominaisuuksia. Tästä huolimatta etenkin lisäinformaatiota selittävien muuttujien muodossa hyväksikäyttävä malli tuottaa sangen tarkkoja ja luotettavia ennusteita. Tämä kertoo osaltaan perusrakennemallien joustavuudesta erilaisten empiiristen ongelmien ratkaisemisessa. Mallin komponentit ovat luonteeltaan sellaisia, että niille on helppo löytää myös empiirinen tulkinta. Myös tämä osaltaan laajentaa perusrakennemalleilla ratkaistavissa olevaa ongelmakenttää. Näin ollen onkin odotettavissa, että perusrakennemallien edustama rakenneyhtälömalliajattelu yleistyy aikasarja-analyysin perustyökaluna ja syrjäyttää myös samalla ennustamisessa yleisesti sovelletun ARIMA-malliajattelun.

Liite 1. Havaintoaineisto

KK	= ajanjakso
ALKUP.	= alkuperäinen kokonaiskulutus (MWh)
MWh	= korjattu kokonaiskulutus (MWh)
MWh-s	= korjattu kokonaiskulutus ilman Säynätsaloa (eroaa edellisestä 1/1996 alkaen).
APL	= astepäiväluku
NORM	= normaalivuoden astepäiväluku
KH	= lämmön kokonaishankinta (MWh)

KK	ALKUP.	MWh	MWh-S	APL	NORM	KH
89-01	67232	67232	67232	608	836	73200
89-02	58100	58100	58100	532	747	6500
89-03	63589	63589	63589	533	672	66900
89-04	45228	45228	45228	390	469	5000
89-05	28799	28799	28799	171	216	29100
89-06	20	11490	11490	30	41	14800
89-07	22105	10636	10636	7	14	13700
89-08	14452	14452	14452	79	65	19000
89-09	25171	25171	25171	185	248	30000
89-10	46449	46449	46449	441	419	55100
89-11	65815	65815	65815	561	575	69200
89-12	82769	82769	82769	772	751	90600
90-01	89091	89091	89091	797		93521
90-02	58684	58684	58684	484		64024
90-03	65483	65483	65483	574		71134
90-04	46921	46921	46921	390		49325
90-05	30236	30236	30236	219		34492
90-06	84	15090	15090	50		17923
90-07	28158	13152	13152	6		15622
90-08	13821	13821	13821	40		18756
90-09	33945	33945	33945	308		41946
90-10	53966	53966	53966	415		56870
90-11	73162	73162	73162	637		79763
90-12	73094	73094	73094	621		78487
91-01	86537	86537	86537	728		89447
91-02	85184	85184	85184	756		90304
91-03	72508	72508	72508	625		77239
91-04	55367	55367	55367	440		56593
91-05	39062	39062	39062	327		43837
91-06	1	18236	18236	64		21453
91-07	31090	12856	12856	0		15123
91-08	13704	13704	13704	12		17193
91-09	32510	32510	32510	268		37362
91-10	46906	46906	46906	390		54016
91-11	57386	57386	57386	463		63201
91-12	69908	76482	76482	639		81364

KK	ALKUP.	MWh	MWh-S	APL	NORM	KH
92-01	90945	80887	80887	664	836	86050
92-02	71870	71870	71870	620	747	78125
92-03	67616	67616	67616	547	672	71282
92-04	62134	62134	62134	526	469	65980
92-05	137	29059	29059	168	216	32585
92-06	41255	12194	12194	17	41	15717
92-07	12333	12472	12472	49	14	16075
92-08	20125	20125	20125	62	65	22302
92-09	27549	27549	27549	189	248	32126
92-10	68674	68674	68674	614	419	77002
92-11	80528	80528	80528	661	575	84480
92-12	70324	70324	70324	563	751	78288
93-01	79444	79444	79444	664		85569
93-02	77110	77110	77110	621		81976
93-03	73423	73423	73423	589		77497
93-04	54241	54241	54241	448		59790
93-05	56	24701	24701	79		27651
93-06	48734	24087	24087	120		26963
93-07	14479	14479	14479	28		17132
93-08	20103	20103	20103	141		24124
93-09	45670	45670	45670	372		52156
93-10	57935	57935	57935	484		66746
93-11	84103	84103	84103	670		93312
93-12	70365	92042	92042	720		97917
94-01	91747	91747	91747	769		104233
94-02	108246	108246	108246	928		119477
94-03	82711	82711	82711	672		90415
94-04	51129	51129	51129	403		58527
94-05	36405	36405	36405	299		42931
94-06	24219	24219	24219	23		24038
94-07	14519	14519	14519	5		16204
94-08	18021	18021	18021	75		20823
94-09	36327	36327	36327	239		40107
94-10	58912	58912	58912	446		67327
94-11	79146	79146	79146	604		89386
94-12	89540	89540	89540	612		90961
95-01	96929	96929	96929	695		102056
95-02	80517	80517	80517	566		81687
95-03	79073	79073	79073	569		84088
95-04	61362	61362	61362	475		66759
95-05	42001	42001	42001	271		48203
95-06	18207	18207	18207	0		18240
95-07	16686	16686	16686	25		20119
95-08	17767	17767	17767	54		21514
95-09	36350	36350	36350	223		41417
95-10	54068	54068	54068	334		56701

KK	ALKUP.	MWh	MWh-S	APL	NORM	KH
95-11	87238	87238	87238	646		98067
95-12	89058	115555	115555	869		122931
96-01	114046	114046	112191	759	836	116427
96-02	122826	122826	120835	880	747	127011
96-03	97230	97230	95684	677	672	98065
96-04	70179	70179	69221	464	469	72925
96-05	42876	45926	44651	245	216	48857
96-06	27716	23797	23380	39	41	25316
96-07	22239	22239	21822	39	14	25068
96-08	18343	18343	17972	0	65	20876
96-09	45355	45355	44627	289	248	48118
96-10	62441	62441	61408	376	419	64951
96-11	70445	74996	74130	480	575	80641
96-12	136353	112208	111491	791	751	11937

Lähteet

- Almon, S. (1965). The distributed lag between capital appropriations and expenditures. *Econometrica*, 33: 178 - 196.
- Anderson, B.D.O. ja Moore, J.B. (1979). *Optimal Filtering*. Englewood Cliffs: Prentice Hall.
- Bowman, K.O. ja Shenton, L.R. (1975). Omnibus test contours for departures from normality based on $\sqrt{\beta_1}$ and β_2 , *Biometrika*, 62: 243 - 250.
- Box, G.E.P. ja Jenkins, G.M. (1970). *Time Series Analysis: Forecasting and Control*. San Francisco, CA: Holden-Day.
- Box, G.E.P. ja Jenkins, G.M. (1976) *Time Series Analysis: Forecasting and Control*, 2. edition. San Francisco, CA: Holden-Day.
- Boyles, R.A. (1983). On the convergence of the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 45: 47 - 50.
- Brown, R.G. (1963). *Smoothing, Forecasting and Prediction*. Englewood Cliffs: Prentice Hall.
- Dagum, E.B. (1975). Seasonal factor forecasts from ARIMA models. *Bulletin of the International Statistical Institute*, 46: 203 - 216.
- de Jong, P. (1988). The likelihood for a state space model. *Biometrika*, 75: 165 - 169.
- de Jong, P. ja Chu-Chun-Lin, S. (1994). Fast likelihood evaluation and prediction for non-stationary state space models. *Biometrika*, 81: 133 - 142.
- Dempster, A.P., Laird, N.M. ja Rubin, D.B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39: 1 - 38.
- Doornik, J.A. ja Hansen, H. (1994). An omnibus test of univariate and multivariate normality, *Unpublished paper*, Nuffield College, Oxford.
- Engle, R.F. (1978). Estimating structural models of seasonality. In A. Zellner (ed.). *Seasonal Analysis of Economic Time Series*, s 281 - 308. Washington D.C: Bureau of the Census.
- Engle, R.F. (1984). Wald, likelihood ratio and Lagrange multiplier model diagnostics. In Z. Griliches and M.D. Intriligator (eds.) *Handbook of Economics*, vol 2, 775 - 828.
- Engle, R.F. ja Hendry, D.F. (1983). Exogeneity. *Econometrica*, 51: 277 - 304.

- Franzini, L. ja Harvey, A.C. (1983). Testing for deterministic trend and seasonal components in time series models. *Biometrika*, 70: 673 - 682.
- Garbade, K. (1977). Two methods for examining the stability of regression coefficients. *Journal of the American Statistical Association*, 72: 54-63.
- Hannan, E.J., Terrel, R.D. ja Tuckwell, N. (1970) Seasonal adjustment of economic time series. *International Economic Review*, 11: 24-52.
- Harrison, P.J. ja Stevens, C.F. (1971). A Bayesian approach to short-term forecasting. *Operational Research Quarterly*, 22: 821-842
- Harrison, P.J. ja Stevens, C.F. (1976) Bayesian forecasting. *The Journal of Royal Statistical Society, Series B*, 38: 205 - 247.
- Harvey, A.C. (1981). *The Econometric Analysis of Time Series*. Deddington, Oxford: Phillip Allan; New York: Wiley.
- Harvey, A.C. (1989). *Forecasting, structural time series model and Kalman filter*. Cambridge: Cambridge University Press.
- Harvey, A.C. ja Hotta, L.K. (1982). Specification tests for dynamic models with unobserved components. *Unpublished paper*, LSE.
- Harvey, A.C. ja Koopman, S.J. (1992). Diagnostic checking of unobserved components time series models. *Journal of Economics and Business Statistics*, 10: 377 - 389.
- Harvey, A.C. ja Phillips, G.D.A. (1979). The maximum likelihood estimation of regression models with autoregressive-moving average disturbances. *Biometrika*, 66: 49-58.
- Helminen, V.A. (1989). Lämmitystä varten laskettu astepäiväluku. *LVI-lehti*, 9/1989: 33 - 35.
- Holt, C.C. (1957). Forecasting seasonals and trends by exponentially weighted moving averages. *ONR Research Memorandum 52*, Carnegie Institute of Technology, Pittsburgh, Pennsylvania.
- Jorgensson, D.W. (1966). Rational distributed lag function. *Econometrica*, 34: 135 - 149.
- Kaminski, P.G., Bryson, A.E. & Schmidt, S.F. (1971). Discrete square root filtering: a survey of current techniques. *I.E.E.E. Trans. Auto. Control* AC-16: 727 - 737.
- Kendall, M.G. (1973). *Time-Series*. London: Charles Griffin & Co. Ltd.
- Kitagawa, G. ja Gersch, W. (1984). A smoothness prior - state space modeling of time series with trend and seasonality. *Journal of American Statistical Association*, 79: 378 - 389.

- Koopman, S.J. (1993). Disturbance smoother for state space models. *Biometrika*, 80: 117 - 126.
- Koopman, S.J., Harvey, A.C., Doornik, J.A. ja Shephard, N. (1995). *STAMP 5.0 Structural Time Series Analyser, modeller and Predictor*. London: Chapman & Hall.
- Ljung, G.M. ja Box, G.E.P. (1978). On a measure of lack of fit in time series models. *Biometrika*, 65: 297 - 305.
- Muth, J.F. (1960). Optimal properties of exponentially weighted forecasts. *Journal of the American Statistical Association*, 55: 299-305.
- Nerlove, M. ja Wage, S. (1964). On the optimality of adaptive forecasting. *Management Science*, 10: 207-209.
- Rosenberg, B. (1973). Random coefficient models: the analysis of a cross-section of time series by stochastically convergent parameter regression. *Annals of Economic and Social Measurement*, 2: 399-428.
- Schweppe, F. (1965). Evaluation of likelihood functions for Gaussian signals. *IEEE Transactions on Information Theory*, 11: 61-70.
- Shephard, N. (1993). Maximum likelihood estimation of regression models with stochastic trend components. *Journal of American Statistical Association*, 88: 590 - 595.
- Shiskin, J., Young, A.H. ja Musgrave, J.C. (1967). The X11 variant of the census method II seasonal adjustment program. *Technical Paper 15*, Bureau of the Census, Washington D.C.
- Suomen Kaukolämpö ry, (1995). *Vuosikertomus 1995*. Espoo: Suomen Kaukolämpö ry.
- Theil, H. ja Wage, S. (1964). Some observations on adaptive forecasting. *Management Science*, 10: 198-206.
- Watson, M.W. ja Engle, R.F. (1983). Alternative algorithms for the estimation of dynamic factor, MIMIC and varying coefficient regression. *Journal of Econometrics*, 23: 385 - 400.
- Winters, P.R. (1960). Forecasting sales by exponentially weighted moving averages. *Management Science*, 6: 324 - 342.
- Wu, C.F.J. (1983). On the convergence of the EM algorithm. *Annals of statistics*, 11: 95 - 103.

Tilastotoimen menetelmien maisteriohjelman pro gradu -tutkielma sarja

1. Salmikuukka, J. (1997) Aikasarjojen perusrakennemalleista ja niiden soveltaminen Jyväskylän kaukolämmön kulutuksen analysointiin ja ennustamiseen. (76 s. , 1 liite) Jyväskylän Energia Oy, Jyväskylä