

<http://www.jyu.fi/library/tutkielmat/354/>

Säännöllistämismenetelmät regressioanalyysissä

Jaana Jokiniemi

Tilastotieteen pro gradu -tutkielma

29. huhtikuuta 1997

Jyväskylän yliopisto
tilastotieteen laitos

Tiivistelmä

Jaana Jokiniemi: *Säännöllistämismenetelmät regressioanalyysissä.*

Tilastotieteen pro gradu -tutkielma, Jyväskylän yliopisto, 29. huhtikuuta 1997.

Sivuja 88, liitteitä 7.

Säännöllistämismenetelmiä käytetään regressioanalyysissä, kun selittävien muuttujien x_1, \dots, x_p välillä on lähes lineaarisia tai lineaarisia riippuvuuksia ts. selittävät muuttujat x_1, \dots, x_p ovat multikollineaarisia, jolloin rakennematriisin $X = (x_1, \dots, x_p)$ neliömatriisi XX ei ole oikein kääntyvä. Vaikka regressiomatriisin B pienimmän neliösumman estimaatilla on minimivarianssi harhattomien estimaattoreiden joukossa, niin selittävien muuttujien ollessa multikollineaarisia sen varianssi ja täten myös sen keskineliövirhe on saattanut kasvaa kohtuuttoman suureksi, jolloin sen pienimmän neliösumman estimaattorin $B_{PNS} = (XX)^{-1}XY$ harhattomuudella ei ole merkitystä. Tällöin säännöllistämismenetelmillä voidaan approksimoida käänteismatriisia $(XX)^{-1}$ matriisilla G , joka ei ole yhtä herkkä selittävien muuttujien multikollineaarisuuden vaikutukselle kuin matriisi $(XX)^{-1}$. Matriisi G pyritään valitsemaan siten, että säännöllistämismenetelmillä saadaan harhainen regressiomatriisin B estimaattori $B_S = GXY$, joka on keskimäärin tarkempi kuin estimaattori B_{PNS} , tai matriisin G käänteismatriisi sisältää suurimman osan rakennematriisin X vaihtelusta, jolla voidaan parhaiten ennustaa havaintomatriisia Y .

Tässä tutkielmassa esitetyt säännöllistämismenetelmät regressioanalyysissä ovat harjaregressio, pääkomponenttiregressio ja osittainen pienimmän neliösumman regressio eli PLS-regressio. Kustakin menetelmästä on kuvaus, tarvittavien parametrien estimointi, selittävien muuttujien tai komponenttien valinta malliin ja yhteys toisiin menetelmiin. Menetelmiä sovelletaan teräksen karkenevuuden mallittamiseen.

Avainsanat: lineaarinen regressio; säännöllistämismenetelmät; harjaregressio; pääkomponentti-regressio; osittainen pienimmän neliösumman regressio

Sisällysluettelo

1 Johdanto.....	3
1.1 Tutkielman sisältö.....	5
2 Regressioanalyysi.....	7
2.1 Kanoninen regressiomalli.....	8
2.2 Pienimmän neliösumman menetelmä, kun selitettäviä muuttujia on yksi.....	9
2.3 Pienimmän neliösumman menetelmä, kun selitettäviä muuttujia on useita.....	13
2.4 Multikollinearisuus.....	15
2.5 Selittävien muuttujien multikollinearisuuden vaikutus pienimmän neliösummanmenetelmässä.....	18
3 Harjaregressio.....	21
3.1 Yksimuuttujaiset harjaestimaattorit $b_H(k)$ ja $a_H(k)$	21
3.2 Yksimuuttujaisen harjaestimaattoreiden $b_H(k)$ ja $a_H(k)$ keskineliövirhe.....	23
3.3 Harjaregression geometrinen tulkinta.....	26
3.4 Ennusteet harjaregressiossa.....	27
3.5 Harjaparametrin k valintamenetelmiä yksimuuttujaisessa tapauksessa.....	28
3.6 Monimuuttujainen harjaregressio.....	30
3.7 Harjaparametrimatriisin K valinta monimuuttujaisessa tapauksessa.....	32
3.7.1 Harjaparametrimatriisin K valinta Brownin ja Paynen (1975) mukaan.....	32
3.7.2. Harjaparametrimatriisin K valinta Brownin ja Zidekin (1982) mukaan.....	33
3.7.3 Harjaparametrimatriisin K valinta Frankin ja Friedmanin mukaan (1993)....	34
3.7.4 Harjaparametrimatriisin K valinta <i>cross-validation</i> -menetelmällä.....	35
3.7.5 Harjaparametrimatriisin K valinta Fülen (1995) mukaan.....	35
4 Pääkomponenttiregressio.....	37
4.1 Yksimuuttujainen pääkomponenttiestimaattori a_{APK}	39
4.2 Yksimuuttujaisen pääkomponenttiestimaattorin a_{APK} keskineliövirhe.....	40
4.3 Ennusteet pääkomponenttiregressiossa	41

4.4 Yksimuuttujaisessa tapauksessa valittujen pääkomponenttien indeksijoukon S_A valintamenetelmiä.....	42
4.5 Monimuuttujainen pääkomponenttiregressio.....	44
4.6 Monimuuttujaisessa tapauksessa valittujen pääkomponenttien indeksijoukon S_A valintamenetelmiä.....	46
4.6.1 <i>Cross-validation</i>	46
5 Osittainen pienimmän neliösumman regressio eli PLS-regressio.....	47
5.1 PLS-algoritmi ja sen geometrinen tulkinta.....	48
5.2 Muita osittaisen pienimmän neliösumman tulkintoja.....	54
5.3 PLS-algoritmin matemaattinen perusta.....	59
5.4 Ennusteet PLS-regressiossa.....	62
5.5 PLS-regressiomallin ulottuvuuden valinta.....	63
5.5.1 PLS-regressiomallin jäännösten graafinen tarkastelu.....	63
5.5.2 <i>Cross-validation</i>	63
6 Teräksen karkenevuuden mallittamisesta.....	65
6.1 Sovitteen ja ennusteen selityskertoimet.....	68
6.2 Selittävien muuttujien multikollinearisuus.....	70
6.3 Pienimmän neliösumman regressio.....	71
6.4 Harjaregressio.....	73
6.5 Pääkomponenttiregressio.....	77
6.6 Osittainen pienimmän neliösumman regressio.....	80
7 Yhteenveto.....	85
Lähteet.....	87
Liite A.....	89

1 Johdanto

Teräksen karkenevuus on eräs sen laatuominaisuus, joka voidaan määritellä sen kykyinä muuttua kiderakenteeltaan kokonaan tai osittain austeniitista martensiitiksi, kun se karkaisuvehkutuksen jälkeen sammutetaan tietyllä tavalla. Näin määriteltyyn teräksen karkenevuuteen vaikuttavat teräksen seostus eli seosaineiden (C, Si,...) pitoisuudet, kappaleen koko ja muoto sekä hehkutuslämpötila ja jäähdytysnopeus. Teräksen karkenevuutta voidaan mitata Jominy-kokeella, jossa standardin SFS 2375 (tai iso 642-1979) mukaan terässauva kuumennetaan austenitointilämpötilaan, siirretään se nopeasti erityiseen jäähdytyslaitteeseen ja sammutetaan suuntaamalla vesisuihkun alapäähän. Näin terässauva jäähtyy sitä hitaammin mitä ylemmäs sen alapäästä nousee ja erityisesti koesauvan alapää jäähtyy äärimmäisen nopeasti muuttuen kiderakenteeltaan täysin marteensiittiseksi. Terässauvan jäähdytyä täysin hiotaan sen sivuille mittaustasot, joista sen kovuus mitataan joko Rockwellin C-yksikkönä tai Vickers-yksikkönä. Kovuusmittaukset tehdään esimerkiksi 1.5, 3, 5, 7, 9, 11, 13, 15, 20, 25,...,50 mm etäisyyksiltä terässauvan alapäästä. Mittaukset suoritetaan erikseen kummastakin tasosta ja mikäli tulokset poikkeavat toisistaan alle 2 HRC (tai alle 20 HV), vastaavat lukemat keskiarvostetaan. Saadut etäisyys-kovuus -parit muodostavat Jominy-käyrän.

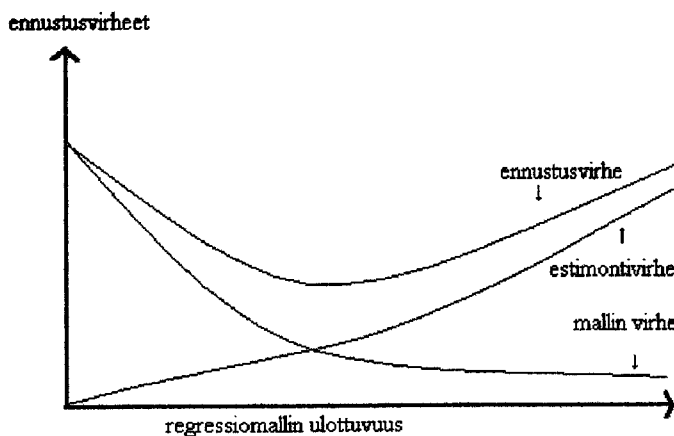
Tehtävänä on muodostaa lineaarinen regressiomalli teräksen karkenevuudelle. Koska Jominy-koetta varten on vakioitu kaikki muut teräksen karkaisu tulokseen vaikuttavat tekijät paitsi seosaineiden pitoisuudet, niin muodostettavassa teräksen karkenevuuden regressiomallissa selitettävänä muuttujina ovat eri etäisyyksiltä terässauvan alapäästä mitatut teräksen kovuudet ja selittävinä muuttujina ovat teräksen seosaineiden pitoisuudet. (Huhtala, 1991.)

Kun X on selittävien muuttujien rakennematriisi, Y on selitettävien muuttujien havaintomatriisi ja muodostettava teräksen karkenevuuden regressiomalli ratkaistaan pienimmän neliösumman menetelmällä, niin tällöin joudutaan kääntämään neliömatriisi $X'X$. Koska teräksen seosaineiden pitoisuudet ovat multikollineaarisia, niin neliömatriisi $X'X$ ei ole oikein kääntyvä. Teräksen seosaineiden pitoisuuksien multikollineaarisuuden ongelmaa ei voida ratkaista poistamalla multikollineaarisia teräksen seosaineiden pitoisuuksia, sillä muodostettavan regressiomallin käytön kannalta kaikkien teräksen seosaineiden mukanaolo on tärkeitä. Tässä tutkielmassa teräksen karkenevuuden regressiomalli estimoidaan harjaregressiolla, pääkomponenttiregressiolla ja osittaisella pienimmän neliösumman regressiolla eli PLS-regressiolla. Nämä menetelmät ovat säännöllistämismenetelmiä, joilla saadaan käänteismatriisille $(X'X)^{-1}$ approksimaatio G , joka ei ole

yhtä herkkä selittävien muuttujien multikollineaarisuuden vaikutukselle kuin matriisi $(X'X)^{-1}$.

Vaikka regressiomatriisin B pienimmän neliösumman estimaatilla on minimivarianssi harhattomien estimaattoreiden joukossa, niin selittävien muuttujien ollessa multikollineaarisia sen varianssi on saattanut kasvaa kohtuuttoman suureksi, koska pienimmän neliösumman estimaattorissa $B_{PNS} = (X'X)^{-1}X'Y$ joudutaan kääntämään singulaarinen tai lähes singulaarinen neliömatriisi $X'X$. Jos selittävät muuttujat ovat multikollineaarisia, niin säännöllistämismenetelmillä voidaan saada regressiomatriisin B estimaattori $B_s = GXY$, joka on hiukan harhainen, mutta jonka varianssi on huomattavasti pienempi kuin harhattoman estimaattorin B_{PNS} varianssi. Estimaattorin hyvyttä voidaan mitata estimaattorin keskineliövirheellä, joka on estimaattorin harhan neliön ja varianssin summa. Säännöllistämismenetelmillä pyritään valitsemaan matriisi G siten, että estimaattorin B_s keskineliövirhe on pienempi kuin estimaattorin B_{PNS} keskineliövirhe.

Säännöllistämismenetelmiä on käytetty usein ennustamisessa, koska säännöllistämismenetelmillä voidaan vähentää regressiomallin ulottuvuutta. Jotta saadaan ennustevirheeltään pienin mahdollinen ennuste, niin regressiomallin ulottuvuuden tulee olla sopiva (ks. kuva 1).



Kuva 1: Ennustevirhe regressiomallin ulottuvuuden funktiona (Naes & Irgens, 1986)

Kemiassa ja monissa teollisissa sovelluksissa on havaittu todellisten selittävien piilomuuttujien ja havaittujen muuttujien välinen ero, sillä kemiallisen prosessin ulottuvuus saattaa muuttua, jos sitä tutkitaan uudesta näkökulmasta. Sosiologiassa ja psykologiassa vastaava ongelma on jo kauan

tiedostettu, minkä seurauksena faktorimalleista on tullut yksi tilastotieteen keskeinen osa-alue. Selittävien muuttujien valinnan ongelma on usein helpommin ratkaistavissa kemiallisissa sovelluksissa kuin psykologisissa ja sosiologisissa sovelluksissa, koska kemiallisissa sovelluksissa kuten Jominy-kokeessa selittävät muuttujat ovat kokeellisesti tulkittavissa ja kiistattomasti mitattavissa.

Tässä tutkielmassa muodostettavan teräksen karkenevuuden regressiomallin ulottuvuutta approksimoidaan säännöllistämismenetelmillä. PLS-regressiossa ja pääkomponenttiregressiossa regressoidaan havaintomatriisi Y rakennematriisia X vasten, mutta lisäksi niissä mallitetaan rakennematriisia X ja PLS-regressiossa myös havaintomatriisia Y . Pääkomponenttiregressiossa approksimoitu regressiomallin ulottuvuus on rakennematriisin X valittujen pääkomponenttien lukumäärä, joka on pienempi tai yhtä pieni kuin selittävien muuttujien lukumäärä. PLS-regressiossa muodostetaan PLS-komponentit siten, että ne selittävät mahdollisimman vähällä komponenttien lukumäärällä mahdollisimman paljon rakennematriisin X kokonaisvaihtelusta, jolla voidaan parhaiten ennustaa havaintomatriisia Y . Harjaregressiossa, kuten myös pienimmän neliösumman regressiossa approksimoitu regressiomallin ulottuvuus on regressiomallissa olevien selittävien muuttujien lukumäärä ts. harjaregressiossa ja pienimmän neliösumman regressiossa ei regressiomallin ulottuvuutta vähennetä.

1.2 Tutkielman sisältö

Tässä tutkielmassa perehdytään harjaregressioon, pääkomponenttiregressioon ja osittaisen pienimmän neliösumman regressioon. Näitä menetelmiä käytetään siis silloin, kun selittävät muuttujat ovat multikollineaarisia.

Tämän tutkielman luvussa 2 on regressioanalyysin perusteita. Kappaleessa 2.1 esitetään kanoninen regressiomalli, joka on yleinen regressiomallin muoto, kun säännöllistämismenetelmiä käytetään. Kappaleessa 2.2 esitetään pienimmän neliösumman menetelmä, jota verrataan säännöllistämismenetelmiin. Koska selittävien muuttujien multikollineaarisuutta ei voida nähdä suoraan selittävien muuttujien korrelaatiomatriisista, niin kappaleessa 2.4 on menetelmiä selittävien muuttujien multikollineaarisuuden havaitsemiseen. Vaikka regressiomatriisin B pienimmän neliösumman estimaatilla on minimivarianssi harhattomien estimaattorien joukossa, niin kappaleessa 2.5 esitetään, miksi pienimmän neliösumman menetelmän hyvyys voidaan kyseenalaistaa, kun selittävät muuttujat ovat multikollineaarisia.

Luvussa 3 on esitetty harjaregressio, luvussa 4 pääkomponenttiregressio ja luvussa 5

osittainen pienimmän neliösumman regressio. Kustakin menetelmästä on kuvaus, tarvittavien parametrien estimointi ja selittävien muuttujien tai komponenttien valinta malliin. Lisäksi näissä luvuissa osoitetaan menetelmien yhteys toisiin säännöllistämismenetelmiin.

Säännöllistämismenetelmiä sovelletaan teräksen karkenevuusaineistoon luvussa 6. Muodostettavassa regressiomallissa on selitettävinä muuttujina eri etäisyyksiltä terässauvan alapäästä mitatut teräksen kovuudet, joten muodostettava regressiomalli on monimuuttujainen. Koska ennalta jo tiedetään, että teräksen seosaineiden pitoisuudet korreloivat voimakkaasti keskenään, niin muodostettavaa teräksen karkenevuuden regressiomallia estimoidaan säännöllistämismenetelmillä. Eri menetelmillä estimoituja teräksen karkenevuuden regressiomallien riittävyttä arvioidaan selityskertoimien avulla. Muodostettujen mallien selityskertoimet lasketaan mallin sovitteelle eli havainnoille, joista regressiokertoimien estimaattit on laskettu, ja ennusteelle eli havainnoille, jotka eivät olleet mukana estimoitaessa regressiomallin regressiokertoimien estimaatteja.

2 Regressioanalyysi

Regressioanalyysissä tutkitaan satunnaismuuttujien $Y^{(1)}, \dots, Y^{(q)}$ riippuvuutta eli regressiota annetuista muuttujista $X = (x_1, \dots, x_p)$. Regressioanalyysissä satunnaismuuttujia $Y^{(1)}, \dots, Y^{(q)}$ kutsutaan selitettäviksi muuttujiksi (engl. explained variables), riippuviksi muuttujiksi (engl. dependent variables) tai vastemuuttujiksi (engl. response variables) ja muuttujia x_1, \dots, x_p selittäviksi muuttujiksi (engl. explanatory variables), ennustaviksi muuttujiksi (engl. predictor variables) tai riippumattomiksi muuttujiksi (engl. independent variables). Kun selittäviä muuttujia on enemmän kuin yksi eli $p > 1$, niin regressioanalyysiä kutsutaan usean selittäjän regressioanalyysiksi tai moniulotteiseksi regressioanalyysiksi (engl. multiple regression). Kun selitettäviä muuttujia on enemmän kuin yksi eli $q > 1$, niin regressioanalyysiä kutsutaan tässä monimuuttujaiseksi regressioanalyysiksi (engl. multivariate regression). (Wang & Chow, 1994.)

Oletetaan, että satunnaismuuttujien $Y^{(1)}, \dots, Y^{(q)}$ ehdollinen odotusarvo ehdolla x_1, \dots, x_p on

$$E(Y|X) = I\beta_0 + XB,$$

missä $Y = (Y^{(1)}, \dots, Y^{(q)})$, $X = (x_1, \dots, x_p)$, vektori β_0 on ns. vakiovektori, jonka alkiot ovat vakiota (engl. constant term), I on ykkösvektori ja B on regressiomatriisi (engl. regression coefficients matrix). Satunnaismuuttujat $Y^{(1)}, \dots, Y^{(q)}$ ja annetut muuttujat x_1, \dots, x_p ovat mitattavissa olevia suureita tai joitain niiden funktiota. Vakiovektori β_0 ja regressiomatriisi B ovat tuntemattomia parametreja, jotka on kiinnitetty, ja ne tulee estimoida havaintoaineistosta. Jos jokaisessa havainnossa $i = 1, \dots, n$ on havaittu sekä $y_i = (y_i^{(1)}, \dots, y_i^{(q)})$ että $x_i = (x_{i1}, \dots, x_{ip})$ ja ehdollinen odotusarvo $E(Y|X)$ on lineaarinen annetulla X , niin

$$Y = E(Y|X) + E = I\beta_0 + XB + E, \quad (2.1)$$

missä siis Y on $(n \times q)$ -havaintomatriisi, vakiovektorin β_0 ulottuvuus on $1 \times q$, I on $(n \times 1)$ -ykkösvektori, X on $(n \times p)$ -rakennematriisi, regressiomatriisin B ulottuvuus on $p \times q$ ja E on $(n \times q)$ -jäännösmatriisi. Kun merkitään $\beta_0 = (\beta_0^{(1)}, \dots, \beta_0^{(q)})$, $B = (\beta^{(1)}, \dots, \beta^{(q)})$ ja $E = (\epsilon^{(1)}, \dots, \epsilon^{(q)})$, niin regressiomalli (2.1) voidaan esittää myös muodossa

$$y_i^{(j)} = I\beta_0^{(j)} + \beta_1^{(j)}x_{i1} + \dots + \beta_p^{(j)}x_{ip} + \epsilon_i^{(j)}, \quad i = 1, \dots, n \text{ ja } j = 1, \dots, q. \quad (2.2)$$

(Brown, 1993.)

Regressiomallissa (2.1) tuntemattomat parametrit β_0 ja B ovat lineaarisia satunnaismuuttujan Y ehdolliselle odotusarvolle $E(Y|X)$, minkä seurauksena regressiomallia (2.1) voidaan pitää lineaarisena. Odotusarvon $E(Y|X)$ lineaarisuus selittävien muuttujien x_1, \dots, x_p suhteen ei siis ole oleellinen vaatimus lineaarisessa regressiomallissa. (Brown, 1993.)

Regressiomallissa (2.1) oletetaan, että jäännösmatriisin E odotusarvo on nollamatriisi. Merkitään jäännösmatriisin E kovarianssimatriisia merkinällä Σ_E . Kovarianssimatriisin Σ_E ulottuvuus on $nq \times nq$ ja se on positiivisesti definiitti matriisi, mitä merkitään $\Sigma_E > \theta$. Jos muuttujien X ja Y välinen lineaarinen riippuvuus on täydellistä, niin $E = \theta$. Tuskin koskaan E on nollamatriisi, sillä yleensä kokeessa on aina satunnaisvaihtelua, joka poikkeuttaa arvoja muuttujien X ja Y välisestä täydellistä lineaarisesta riippuvuudesta. Täten havaintomatriisin Y ja rakenne-matriisin X riippuvuussuhteen kuvaaminen regressiomallin avulla on käytännössä likimääräistä.

2.1 Kanoninen regressiomalli

Siirrytään käyttämään standardoituja muuttujia, joille on tehty muunnokset

$$x_{ij} = \frac{x_{ij}^* - \bar{x}_j^*}{\sqrt{\text{var}(x_j^*)(n-1)}}, \quad j = 1, \dots, p \quad (2.3)$$

missä x_j^* on alkuperäisen selittävän muuttujan x_j^* keskiarvo, $\text{var}(x_j^*)$ on alkuperäisen selittävän muuttujan x_j^* varianssi, n on havaintoaineiston otoskoko ja $i = 1, \dots, n$, ja

$$y_i^{(j)} = \frac{y_i^{*(j)} - \bar{y}^{*(j)}}{\sqrt{\text{var}(y^{*(j)})(n-1)}}, \quad j = 1, \dots, q, \quad (2.4)$$

missä $\bar{y}^{*(j)}$ on alkuperäisen selitettävän muuttujan $y^{*(j)}$ keskiarvo, $\text{var}(y^{*(j)})$ on alkuperäisen selitettävän muuttujan $y^{*(j)}$ varianssi, n on havaintoaineiston otoskoko, ja $i = 1, \dots, n$. Nyt regressiomallissa (2.1) vakiovektori β_0 on nollavektori. Täten regressiomalli (2.1) voidaan kirjoittaa muotoon

$$Y = E(Y|X) + E = I\beta_0 + XB + E = XB + E . \quad (2.5)$$

Määritellään seuraavaksi regressiomallin (2.5) kanoninen muoto. Olkoon matriisit T , jonka ulottuvuus on $n \times n$, ja V , jonka ulottuvuus on $p \times p$, ortogonalia siten, että

$$U = T'Y = T'XVV'B + T'E , \quad (2.6)$$

missä matriisin T sarakkeina ovat neliömatriisin XX' ominaisvektorit ja matriisin V sarakkeina neliömatriisin XX ominaisvektorit. Oletetaan tästä lähtien, että neliömatriisin XX ominaisarvot on järjestetty seuraavaan suuruusjärjestykseen $\lambda_1 \geq \dots \geq \lambda_p$, minkä mukaan on järjestetty matriiseissa V ja T ominaisvektorit $v_j, j = 1, \dots, p$, ja $t_j, j = 1, \dots, n$. Tulo $T'XV$ on

$$T'XV = \begin{pmatrix} \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_p}) \\ \mathbf{0} \end{pmatrix} , \quad (2.7)$$

missä nollamatriisin $\mathbf{0}$ ulottuvuus on $n - p \times p$ ja neliömatriisin XX ominaisarvojen $\lambda_1, \dots, \lambda_p$ neliöjuuret ovat matriisin X singulaariarvoja. Matriisi X voidaan kirjoittaa muotoon

$$X = T \begin{pmatrix} \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_p}) \\ \mathbf{0} \end{pmatrix} V' , \quad (2.8)$$

jota kutsutaan matriisin X singulaariarvohajotelmaksi. Matriisin X asteelle r on voimassa epäyhtälö $r \leq \min(n, p)$ ja se on neliömatriisin $X'X$ nollassa eroavien ominaisarvojen lukumäärä. Täten matriisi X voidaan esittää muodossa

$$X = T_r \begin{pmatrix} \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_r}) \\ \mathbf{0} \end{pmatrix} V_r' \quad (2.9)$$

ja tulo (2.7) on muotoa

$$T_r' X V_r = \begin{pmatrix} \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_r}) \\ \mathbf{0} \end{pmatrix} , \quad (2.10)$$

missä nollamatriisin $\mathbf{0}$ ulottuvuus on $n - r \times r$.

Kun merkitään

$$V'B = A , \quad (2.11)$$

niin regressiomalli (2.6) voidaan kirjoittaa muotoon

$$u_i^{(j)} = \sqrt{\lambda_i} \alpha_i^{(j)} + \epsilon_i^{(j)}, \quad i = 1, \dots, r, \quad \text{ja} \quad u_i^{(j)} = \epsilon_i^{(j)}, \quad i = r + 1, \dots, n, \quad (2.12)$$

missä $\epsilon^{(j)} = T\epsilon^{(j)}$, r on matriisin X aste ja $j = 1, \dots, q$. Tätä regressiomallia kutsutaan kanoniseksi regressiomalliksi. (Brown, 1993.)

2.2 Pienimmän neliösumman menetelmä, kun selitettäviä muuttujia on yksi

Oletetaan, että estimoitavana on regressiomalli (2.5), missä selitettäviä muuttujia on yksi eli $q = 1$. Pienimmän neliösumman menetelmässä määrätään parametrin β estimaatti siten, että estimoidun mallin antamien sovitteiden \hat{y} ja havaintojen y välinen jäännösneliösumma

$$S(\beta) = \sum_{i=1}^n \epsilon_i^2 = (y - \hat{y})(y - \hat{y}) = \sum_{i=1}^n (y_i - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 \quad (2.13)$$

on mahdollisimman pieni. Asettamalla neliösumman osittaisderivaatat

$$\frac{\partial S(\beta)}{\partial \beta} = -2X'y + 2X'X\beta$$

nolliksi saadaan normaaliyhtälöt

$$(X'X)b_{PNS} = X'y, \quad (2.14)$$

missä b_{PNS} on parametrin β pienimmän neliösumman ratkaisu. Kun $r = p$, niin neliömatriisi XX on epäsingulaarinen, jolloin pienimmän neliösumman estimaattori b_{PNS} voidaan kirjoittaa muotoon

$$b_{PNS} = (X'X)^{-1}X'y. \quad (2.15)$$

(Liski & Puntanen, 1976.)

Pienimmän neliösumman menetelmässä sovite \hat{y}_{PNS} on

$$\hat{y}_{PNS} = X b_{PNS} = X(X'X)^{-1}X'y = Hy, \quad (2.16)$$

missä $(n \times n)$ -projektiomatriisi $H = X(X'X)^{-1}X'$ on ns. hattumatriisi. Koska pienimmän neliösumman menetelmässä jäännös e voidaan kirjoittaa muotoon

$$\mathbf{e} = (\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{y} = (\mathbf{I}_n - \mathbf{H})\mathbf{y} = \mathbf{P}\mathbf{y}, \quad (2.17)$$

missä matriisi $\mathbf{P} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ on $(n \times n)$ -projektiomatriisi, niin jäännöksen \mathbf{e} riippuvuus selitettävästä muuttujasta \mathbf{y} on lineaarista. Sovitetut jäännökset \mathbf{e} ja rakennematriisin \mathbf{X} sarakevektorit ovat keskenään ortogonaalisia, sillä

$$\mathbf{X}'\mathbf{e} = \mathbf{X}'(\mathbf{y} - \mathbf{X}\mathbf{b}_{PNS}) = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{y} = \mathbf{0} \quad (2.18)$$

ja jäännöksen \mathbf{e} kovarianssimatriisi on

$$\text{cov}(\mathbf{e}, \mathbf{e}) = \mathbf{E}\mathbf{e}\mathbf{e}' = \mathbf{E}(\mathbf{P}\mathbf{y}\mathbf{y}'\mathbf{P}') = \mathbf{P}\mathbf{E}(\mathbf{y}\mathbf{y}')\mathbf{P}' = \mathbf{P}\text{cov}(\mathbf{y})\mathbf{P}'. \quad (2.19)$$

Koska $\mathbf{E}(\mathbf{e}'\mathbf{e}) = \mathbf{E}(\mathbf{y}'\mathbf{P}'\mathbf{P}\mathbf{y}) = \sigma_y^2(n - p)$, niin virhevarianssin σ_y^2 pienimmän neliösumman estimaattori on

$$\hat{\sigma}_y^2 = \frac{\mathbf{e}'\mathbf{e}}{n - p}. \quad (2.20)$$

Estimaattorin \mathbf{b}_{PNS} kovarianssimatriisi on

$$\text{cov}(\mathbf{b}_{PNS}) = \sigma_y^2(\mathbf{X}'\mathbf{X})^{-1}. \quad (2.21)$$

Koska estimaattorin \mathbf{b} ja sen todellisen arvon $\boldsymbol{\beta}$ poikkeuman neliön odotusarvo on estimaattorin keskineliövirhe

$$\text{MSE}(\mathbf{b}) = \mathbf{E}(\mathbf{b} - \boldsymbol{\beta})(\mathbf{b} - \boldsymbol{\beta})' = \mathbf{E}(\mathbf{b} - \mathbf{E}\mathbf{b})(\mathbf{b} - \mathbf{E}\mathbf{b})' + (\mathbf{E}\mathbf{b} - \boldsymbol{\beta})(\mathbf{E}\mathbf{b} - \boldsymbol{\beta})' = \text{var}(\mathbf{b}) + \mathbf{B}^2(\mathbf{b}),$$

missä $\text{var}(\mathbf{b})$ on estimaattorin varianssi ja $\mathbf{B}^2(\mathbf{b})$ on estimaattorin harhan neliö, ja estimaattori \mathbf{b}_{PNS} on harhaton, niin estimaattorin \mathbf{b}_{PNS} kokonaisvarienssi ja samalla sen keskineliövirhe on

$$\text{var}(\mathbf{b}_{PNS}) = \sigma_y^2 \text{tr}(\mathbf{X}'\mathbf{X})^{-1} = \sigma_y^2 \sum_{j=1}^p \lambda_j^{-1} = \text{MSE}(\mathbf{b}_{PNS}). \quad (2.22)$$

Kun lineaarisessa regressiomallissa (2.5) $q = 1$, $r = p$, $\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}$ ja $\text{cov}(\boldsymbol{\epsilon}) = \sigma_y^2 \mathbf{I}_n$, niin Gauss-Markovin lauseen mukaan pienimmän neliösumman estimaattori \mathbf{b}_{PNS} on tarkentuva ja harhaton

estimaattori, jolla on pienin mahdollinen varianssi, eli estimaattori b_{PNS} on BLUE (engl. Best Linear Unbiased Estimator) (Wang & Chow, 1994).

Oletetaan, että tehtävänä on ennustaa annetulla selittävän muuttujan x_0 arvolla tuntematonta selitettävän muuttujan y_0 arvoa käyttäen mallia

$$y_0 = E(y_0|x_0) + \epsilon = x_0' \beta + \epsilon, \quad E(\epsilon) = 0, \quad \text{var}(\epsilon) = \sigma_y^2. \quad (2.23)$$

Pienimmän neliösumman menetelmässä tuntemattoman selitettävän muuttujan y_0 ennuste on $\hat{y}_{0PNS} = x_0' b_{PNS}$. Sen varianssi ja keskineliövirhe on

$$\text{var}(\hat{y}_{0PNS}) = \sigma_y^2(1 + 1/n) + \sigma_y^2 \text{tr}(x_0'(X'X)^{-1}x_0) = \text{MSE}(\hat{y}_{0PNS}). \quad (2.24)$$

(Brown, 1993.)

Yllä olevassa ennusteen keskineliövirheessä annettu selittävän muuttujan arvo x_0 on siis keskistetty. Tällöin mitä enemmän alkuperäinen selittävän muuttujan arvo x_0^* eroaa arvosta $\bar{x} = (\bar{x}_1^*, \dots, \bar{x}_p^*)$, missä \bar{x}_j^* on alkuperäisen selittävän muuttujan x_j^* keskiarvo, $j = 1, \dots, p$, niin sitä suurempia ovat annetun x_0 alkiot ja sitä suurempia ovat $\text{var}(\hat{y}_{0PNS})$ ja $\text{MSE}(\hat{y}_{0PNS})$. Geometrisesti tulkittuna regressiokertoimien estimaateissa olevat virheet aiheuttavat sen, että estimoitu regressiosuora on ikään kuin kiertynyt suhteessa todelliseen regressiosuoraan. Koska regressiosuora kulkee keskiarvopisteen kautta, niin kiertäminen tapahtuu tämän pisteen ympäri. Tämän vuoksi ennusteen epävarmuus on sitä suurempi mitä kauempana selittävän muuttujan arvo on keskiarvostaan. (Ranta, 1991.)

Jos regressiomalli on esitetty kanonisessa muodossa (2.6), missä selitettäviä muuttujia on yksi, niin pienimmän neliösumman estimaattori a_{PNS} on

$$a_{jPNS} = u_j / \sqrt{\lambda_j}, \quad j = 1, \dots, r, \quad \text{ja} \quad a_{jPNS} = 0, \quad j = r + 1, \dots, n, \quad (2.25)$$

missä r on rakennematriisin X aste, sillä normaaliyhtälöt (2.14) ovat nyt muotoa

$$\Lambda_r a_{PNS} = (\sqrt{\Lambda_r} \ 0) u,$$

joten

$$\mathbf{a}_{PNS} = \Lambda_r^{-1} (\sqrt{\Lambda_r} \mathbf{0}) \mathbf{u} = (\sqrt{\Lambda_r}^{-1} \mathbf{0}) \mathbf{u} ,$$

missä $\Lambda_r = \text{diag}(\lambda_1, \dots, \lambda_r)$ ja nollamatriisin $\mathbf{0}$ ulottuvuus on $n - r \times r$.

Pienimmän neliösumman estimaattorin \mathbf{a}_{PNS} kovarianssimatriisi on

$$\text{cov}(\mathbf{a}_{PNS}) = \sigma_y^2 \Lambda_r^{-1} = \sigma_y^2 \text{diag}(1/\lambda_1, \dots, 1/\lambda_r) . \quad (2.26)$$

Huomattavaa on, että estimaattorit $a_{jPNS}, j = 1, \dots, r$, ovat korreloimattomia ja estimaattorit $b_{jPNS}, j = 1, \dots, r$, korreloivat. Koska estimaattori \mathbf{a}_{PNS} on kanonisen regressiomallin (2.6) parametrin α pienimmän neliösumman estimaattori, niin se on Gauss-Markovin lauseen mukaan BLUE. Koska estimaattori \mathbf{a}_{PNS} on harhaton, sen kokonaisvarianssi ja samalla keskineliövirhe on

$$\text{var}(\mathbf{a}_{PNS}) = \text{MSE}(\mathbf{a}_{PNS}) = \sigma_y^2 \text{tr}(\Lambda_r^{-1}) = \sigma_y^2 \sum_{j=1}^r \lambda_j^{-1} . \quad (2.27)$$

Harhattomilla estimaattoreilla \mathbf{b}_{PNS} ja \mathbf{a}_{PNS} on siis sama kokonaisvarianssi ja siten myös yhtäsuuret keskineliövirheet.

Kun matriisilla $\Lambda_p^{1/2}$, joka kuuluu havaintoaineistoon, ennustetaan tuntematonta \mathbf{u} , niin ennusteen $\hat{\mathbf{u}}_{PNS}$ keskineliövirhe on

$$\text{MSE}(\hat{\mathbf{u}}_{PNS}) = \sigma_y^2(n+1) + E((\mathbf{a}_{PNS} - \alpha)' \Lambda_r (\mathbf{a}_{PNS} - \alpha)) \quad (2.28)$$

$$= \sigma_y^2(n+1) + \sum_{j=1}^p \lambda_j \text{var}(a_{jPNS}) = (n+1+p) \sigma_y^2 ,$$

missä $\text{var}(a_{jPNS}) = \sigma_y^2 / \lambda_j$.

2.3 Pienimmän neliösumman menetelmä, kun selitettäviä muuttujia on useita

Oletetaan nyt, että estimoitavana on regressiomalli (2.5), missä selitettäviä muuttujia on useita eli $q > 1$. Pienimmän neliösumman menetelmässä valitaan parametrin \mathbf{B} estimaatti siten, että kaikkien havaintojen poikkeamien neliösumma

$$S(\mathbf{B}) = (\mathbf{Y} - \hat{\mathbf{Y}})'(\mathbf{Y} - \hat{\mathbf{Y}}) = \sum_{i=1}^n \sum_{j=1}^q (y_i^{(j)} - \beta_1^{(j)} x_{i1} - \dots - \beta_p^{(j)} x_{ip})^2$$

on mahdollisimman pieni. Tällöin parametrin \mathbf{B} pienimmän neliösumman estimaattori \mathbf{B}_{PNS} ratkaistaan normaaliyhtälöistä

$$(\mathbf{X}'\mathbf{X})\mathbf{B}_{PNS} = \mathbf{X}'\mathbf{Y} \quad (2.29)$$

(vrt. (2.14)). Kun neliöatriisi $\mathbf{X}'\mathbf{X}$ on epäsingulaarinen, niin pienimmän neliösumman estimaattori \mathbf{B}_{PNS} voidaan kirjoittaa muotoon

$$\mathbf{B}_{PNS} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}. \quad (2.30)$$

Havaintomatriisin \mathbf{Y} pienimmän neliösumman sovite on

$$\hat{\mathbf{Y}}_{PNS} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = \mathbf{H}\mathbf{Y}, \quad (2.31)$$

missä $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ on projektiomatriisi. Koska

$$\mathbf{E} = \mathbf{Y} - \hat{\mathbf{Y}}_{PNS} = (\mathbf{I} - \mathbf{H})\mathbf{Y}, \quad (2.32)$$

niin $E(\mathbf{E}'\mathbf{E}) = E(\mathbf{Y}'(\mathbf{I}_p - \mathbf{H})\mathbf{Y}) = \Sigma_Y \text{tr}(\mathbf{I}_p - \mathbf{H}) = \Sigma_Y (n - p)$ ja täten

$$\hat{\Sigma}_Y = \frac{\mathbf{Y}'(\mathbf{I}_n - \mathbf{H})\mathbf{Y}}{n - p} \quad (2.33)$$

(vrt. (2.20)) (Wang & Chow, 1994).

Estimaattorin vek(\mathbf{B}_{PNS}) kovarianssimatriisi on

$$\text{cov}(\text{vek}(\mathbf{B}_{PNS})) = \Sigma_Y \otimes (\mathbf{X}'\mathbf{X})^{-1}, \quad (2.34)$$

missä \otimes on Kroneckerin tulo, eli $\text{cov}(\mathbf{b}_{PNS}^{(i)} \mathbf{b}_{PNS}^{(j)}) = \sigma^{(ij)}(\mathbf{X}'\mathbf{X})^{-1}$, $i, j = 1, \dots, q$ (Wang & Chow, 1994). Selitettävien muuttujien $y^{(i)}$ ja $y^{(j)}$ korkeat korrelaatiot kasvattavat siis kovarianssia $\text{cov}(\mathbf{b}_{PNS}^{(i)} \mathbf{b}_{PNS}^{(j)})$, $i \neq j$, $i, j = 1, \dots, q$. Koska Brownin ja Paynen (1975) mukaan

$$\text{MSE}(\mathbf{B}) = E\left(\sum_{i=1}^p \sum_{j=1}^q (b_i^{(j)} - \beta_i^{(j)})^2\right), \quad (2.35)$$

missä $\mathbf{B} = (\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(q)})$ on jokin paramerin $\mathbf{B} = (\boldsymbol{\beta}^{(1)}, \dots, \boldsymbol{\beta}^{(q)})$ estimaattori, niin estimaattorin \mathbf{B}_{PNS} varianssi ja keskineliövirhe on varianssien $\text{var}(\mathbf{b}_{PNS}^{(j)}), j = 1, \dots, q$, summa.

Huomattavaa on, että koska monimuuttujainen pienimmän neliösumman estimaattori \mathbf{B}_{PNS} voidaan kirjoittaa muotoon $\mathbf{B}_{PNS} = (\mathbf{b}_{PNS}^{(1)}, \dots, \mathbf{b}_{PNS}^{(q)})$, missä

$$\mathbf{b}_{PNS}^{(j)} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}^{(j)}, \quad j = 1, \dots, q, \quad (2.36)$$

niin estimaattorit $\mathbf{b}_{PNS}^{(j)}, j = 1, \dots, q$, ovat q yksimuuttujaisen regressiomallin

$$\mathbf{y}^{(j)} = \mathbf{X}\boldsymbol{\beta}^{(j)} + \boldsymbol{\epsilon}^{(j)}, \quad j = 1, \dots, q \quad (2.37)$$

parametrien $\boldsymbol{\beta}^{(j)}, j = 1, \dots, q$, pienimmän neliösumman estimaattoreita. Koska ratkaistaessa pienimmän neliösumman estimaattoria \mathbf{B}_{PNS} ei käytetä matriisia $\boldsymbol{\Sigma}_Y$, niin määrättäessä estimaattoria \mathbf{B}_{PNS} ei hyödynnetä q normaaliyhtälöiden yhteisinformaatiota.

2.4 Multikollineaarisuus

Yleisesti multikollineaarisuus termiä käytetään, kun lineaarisen mallin selittävistä muuttujista $\mathbf{x}_j, j = 1, \dots, p$, kaikki tai osajoukko $\mathbf{x}_j, j = A + 1, \dots, p$, on lähes lineaarisesti riippuvia. Vektorit $\mathbf{x}_j, j = 1, \dots, p$, ovat multikollineaarisia, jos on olemassa vakiot $c_j \neq 0$, siten että

$$\sum_{j=1}^p c_j \mathbf{x}_j \approx \mathbf{0}, \quad (2.38)$$

missä vektorit $\mathbf{x}_j, j = 1, \dots, p$, on keskistetty. Täten selittävien muuttujien $\mathbf{x}_j, j = 1, \dots, p$, multikollineaarisuudella tarkoitetaan sitä, että selittävien muuttujien $\mathbf{x}_j, j = 1, \dots, p$, välillä on lähes suoria tai lineaarikombinaatioiden välisiä riippuvuuksia. Mitä lähempänä summa (2.38) on nollaa sitä suurempi on selittävien muuttujien $\mathbf{x}_j, j = 1, \dots, p$, multikollineaarisuuden aste. Harvoin selittävät muuttujat $\mathbf{x}_j, j = 1, \dots, p$, ovat täysin lineaarisesti riippuvia, jolloin summa (2.38) on nolla ja rakennematriisin \mathbf{X} aste $r < p$. Yleensä regressioanalyysissä ongelmana onkin selittävien muuttujien multikollineaarisuuden asteen arvioinen eikä täydellisen lineaarisen riippuvuuden tunnistaminen. (Gunst & Mason, 1980; Leskinen, 1977.)

Selittävien muuttujien multikollineaarisuus samaistetaan usein selittävien muuttujien voimakkaaseen korreloitumiseen ja multikollineaarisuutta arvioidaan selittävien muuttujien korrelaatiokertoimien avulla. Jos multikollineaarisuutta on useamman kuin kahden selittävän

muuttujan välillä, niin multikollineaarisuutta ei voida päätellä selittävien muuttujien korrelaatiomatriisista, koska selittävät muuttujat voivat olla lähes lineaarisesti riippuvia, vaikka niiden väliset korrelaatiokertoimet olisivatkin pieniä (Leskinen, 1977).

Seuraavassa on muutamia tapoja arvioida selittävien muuttujien $x_p, j = 1, \dots, p$, multikollineaarisuutta:

1) Jos selittävien muuttujien $x_p, j = 1, \dots, p$ korrelaatiomatriisissa on korrelaatiokertoimia, jotka ovat suurempia kuin 0.7, niin selittävien muuttujien multikollineaarisuuden astetta voidaan pitää korkeana (Gunst & Mason, 1980).

2) Oletetaan, että neliömatriisin XX pienimmät ominaisarvot $\lambda_p, j = A + 1, \dots, p$, ovat lähes nollia. Ominaisarvojen määritelmästä seuraa, että

$$X'X v_j = \lambda_j v_j \approx 0, \quad j = A + 1, \dots, p. \quad (2.39)$$

missä v_j on ominaisarvoa λ_j vastaava neliömatriisin XX ominaisvektori. Kun kerrotaan yllä oleva yhtälö ominaisvektorilla v_j , niin saadaan

$$(X v_j)'(X v_j) = \lambda_j v_j' v_j = \lambda_j \approx 0. \quad (2.40)$$

Kun $X v_j = z_j$, niin $z_j' z_j = (z_{j1}^2 + \dots + z_{jp}^2)$, missä alkiot $z_{ji}, i = 1, \dots, p$, eivät ole negatiivisia, niin yllä olevassa yhtälössä $z_j = X v_j \approx 0, j = A + 1, \dots, p$. Jos $\lambda_j \approx 0$, niin

$$X v_j = \sum_{i=1}^p v_{ij} x_i \approx 0 \quad (2.41)$$

(vrt. (2.38)). Täten kun neliömatriisilla XX on pieni ominaisarvo $\lambda_j \approx 0$, niin selittävät muuttujat $x_p, j = 1, \dots, p$, ovat multikollineaarisia. Kun $\lambda_j \approx 0$, niin ominaisvektorin v_j suuret alkiot määräävät, mitkä selittävät muuttujat $x_p, j = 1, \dots, p$, ovat eniten multikollineaarisia. (Wang & Chow, 1994.)

3) Neliömatriisin XX determinantti on

$$\det(X'X) = \prod_{j=1}^p \lambda_j. \quad (2.42)$$

Koska $\det(XX)$ on matriisin XX ominaisarvojen tulo, niin selittävien muuttujien multikollineaarisuuden asteen ollessa korkea $\det(XX) \approx 0$. Jos selittävien muuttujien välillä on täydellisiä

lineaarisia riippuvuuksia, niin $\lambda_j = 0$, jolloin $\det(\mathbf{X}\mathbf{X}) = 0$. (Leskinen, 1977.)

4) Matriisin \mathbf{X} sanotaan olevan kollineaarinen, jos

$$f_1 = \sqrt{\lambda_1/\lambda_p}, \quad (2.43)$$

missä λ_1 on suurin ja λ_p on pienin neliömatriisin $\mathbf{X}\mathbf{X}$ ominaisarvo, on suuri (Jackson, 1991). Kun rakennematriisi \mathbf{X} on kollineaarinen, niin se on myös multikollineaarinen.

5) Jos rakennematriisin \mathbf{X} korrelaatiomatriisi $\mathbf{X}\mathbf{X}$ on ortogonallinen eli selittävät muuttujat $\mathbf{x}_j, j = 1, \dots, p$, ovat riippumattomia, niin $\lambda_1 = \dots = \lambda_p$, jolloin $1/\lambda_1 = \dots = 1/\lambda_p = 1$. Mitä enemmän lukujen $1/\lambda_j, j = 1, \dots, p$, summa eroaa selittävien muuttujien lukumäärästä p sitä enemmän selittävien muuttujien $\mathbf{x}_j, j = 1, \dots, p$, välillä on multikollineaarisuutta.

6) Kun oletetaan, että selittävät muuttujat ovat satunnaismuuttujia ja ne noudattavat p -ulotteista normaalijakaumaa, niin multikollineaarisuuden havaitseminen voidaan perustaa testisuureeseen

$$\chi^2 = -(n-1-(1/6))(2p+5) \ln|\mathbf{X}\mathbf{X}|, \quad (2.44)$$

joka noudattaa suurilla otoskoon n arvoilla selittävien muuttujien ollessa korreloimattomia χ^2 -jakaumaa vapausastein $p(p-1)/2$ (Leskinen, 1977).

Huomattavaa on, että kaikkia edellä olevien testien tuloksia on käsiteltävä lähinnä suuntaa antavina.

Kun on todettu, että selittävien muuttujien välillä on multikollineaarisuutta, niin ratkaistaessa selittävien muuttujien multikollineaarisuuden ongelmaa on hyödyllistä yrittää selvittää selittävien muuttujien multikollineaarisuuden lähde, jotta tiedettäisiin, miten tätä ongelmaa tulisi käsitellä regressioanalyysissä. Seuraavassa on kolme selittävien muuttujien multikollineaarisuuden lähdettä Leskisen (1977) mukaan:

1) Multikollineaarisuutta saattaa esiintyä mallissa, jos aineisto on kerätty virheellisesti. Jos selittävien muuttujien multikollineaarisuus johtuu otannasta, niin on kerättävä joko uusi aineisto tai riittävä määrä uutta aineistoa. Multikollineaarisuutta saattaa aiheuttaa myös vieraat (outliers)

havainnot.

2) Regressiomallin rakennematriisi X on väärin määritelty. Yleisemmin regressiomalliin valitaan liian paljon selittäviä muuttujia, joissa 'ylimääräiset selittävät muuttujat' korreloivat muiden muuttujien kanssa.

3) Regressiomalliin sisäiset rajoitukset voivat aiheuttaa selittävien muuttujien välille multikollineaarisuutta. Tällä tarkoitetaan sitä, että selittävien muuttujien multikollineaarisuus on nyt tutkittavan ilmiön luontainen ominaisuus, jota ei voida muuttaa. Tällöin jos multikollineaarisuuden poistamiseksi poistettaisiin regressiomallista multikollineaarisia selittäviä muuttujia, niin tutkittavaa ilmiötä kuvattaisiin puutteellisesti. Jos selittävien muuttujien multikollineaarisuus ei ole vain otoksesta johtuvaa, vaan se osa tutkittavan ilmiön teoreettista ominaisuutta, niin tällöin on turvauduttava harhaisiin estimointimenetelmiin eli säännöllistämismenetelmiin.

2.5 Selittävien muuttujien multikollineaarisuuden vaikutus pienimmän neliösumman menetelmässä

Selittävien muuttujien $x_j, j = 1, \dots, p$, korkea multikollineaarisuuden aste aiheuttaa epätarkkuutta pienimmän neliösumman estimaattorissa b_{PNS} , koska siinä joudutaan kääntämään neliömatriisi $X'X$. Neliömatriisin $X'X$ käänteismatriisi $(X'X)^{-1}$ on

$$(X'X)^{-1} = V_p \Lambda_p^{-1} V_p' = \sum_{j=1}^p (1/\lambda_j) v_j v_j', \quad (2.45)$$

missä $\Lambda_p = \text{diag}(\lambda_1, \dots, \lambda_p)$, $V_p = (v_1, \dots, v_p)$, λ_j on neliömatriisin $X'X$ ominaisarvo ja v_j sitä vastaava ominaisvektori, $j = 1, \dots, p$. Jos selittävät muuttujat $x_j, j = 1, \dots, p$, ovat täysin lineaarisesti riippuvia, niin neliömatriisin $X'X$ pienin ominaisarvo $\lambda_p = 0$ ja $\det(X'X) = 0$. Jos käänteismatriisissa joudutaan jakamaan nolllalla, niin käänteismatriisia $(X'X)^{-1}$ ei ole olemassa ja matriisin $X'X$ sanotaan olevan singulaarinen. Jos selittävät muuttujat eivät ole täysin lineaarisesti riippuvia, mutta niiden multikollineaarisuuden aste on korkea, niin $\lambda_p \approx 0$ ja $\det(X'X) \approx 0$. Tällöin neliömatriisia $X'X$ on epäsingulaarinen eli on olemassa yksikäsitteinen käänteismatriisi $(X'X)^{-1}$, mutta käänteismatriisi $(X'X)^{-1}$ saattaa olla epävakaa ja sen alkiot saattavat olla kohtuuttoman suuria, koska käänteismatriisissa $(X'X)^{-1}$ joudutaan jakamaan hyvin pienellä positiivisella luvulla. Kun $\lambda_p \approx 0$, niin $\det(X'X)$

≈ 0 ja neliömatriisin $X'X$ sanotaan olevan lähes singulaarinen.

Oletetaan nyt, että selittävät muuttujat $x_j, j = 1, \dots, p$, ovat lineaariset riippuvia, jolloin rakennematriisin $X = (x_1, \dots, x_p)$ aste r on pienempi kuin p . Tällöin pienimmän neliösumman estimaattori b_{PNS} ei ole yksikäsitteinen, koska kääntematriisia $(X'X)^{-1}$ ei ole olemassa. Tällöin estimaattori b_{PNS} ei ole parametrin β harhaton estimaattori, sillä jos estimaattori b_{PNS} on parametrin β harhaton estimaattori, niin $E(b_{jPNS}) = \beta_j$ kaikilla $j = 1, \dots, p$, jolloin $(X'X)^{-1}(X'X) = I_p$, mikä on nyt mahdotonta, koska rakennematriisin X aste r on pienempi kuin p . (Wang & Chow, 1994.)

Lause 2.1 Jos rakennematriisin $X = (x_1, \dots, x_p)$ aste r on pienempi kuin p , niin parametrilla β ei ole olemassa harhatonta estimaattoria (Wang & Chow, 1994).

Todistus: Oletetaan, että C on $(p \times n)$ -matriisi ja Cy on parametrin β harhaton estimaattori, jolloin $CX = I_p$. Koska $p = \text{rank}(I_p) = \text{rank}(CX) \leq \text{rank}(X) < p$, niin matriisia C ei ole olemassa. Täten parametrilla β ei ole olemassa harhatonta estimaattoria. \square

Koska pienimmän neliösumman estimaattorin b_{PNS} kovarianssimatriisi on

$$\text{cov}(b_{PNS}) = \sigma_y^2 (X'X)^{-1},$$

niin selittävien muuttujien multikollineaarisuuden asteen ollessa korkea neliömatriisi $X'X$ on lähes singulaarinen, mikä kasvattaa estimaattoreiden $b_{jPNS}, j = 1, \dots, p$, välisiä kovariansseja. Täten, kun selittävät muuttujat ovat multikollineaarisia, niin kovarianssit $\text{cov}(b_{jPNS}, b_{iPNS})$ ja varianssit $\text{var}(b_{jPNS})$ saattavat olla kohtuuttoman suuria, $i, j = 1, \dots, p$.

Pienimmän neliösumman estimaattorin b_{PNS} neliöidun pituuden odotusarvo on

$$E(b_{PNS}' b_{PNS}) = \beta' \beta + \text{var}(b_{PNS}) = \beta' \beta + \sigma_y^2 \text{tr}(X'X)^{-1} \quad (2.46)$$

(Brown, 1993). Koska yllä olevassa yhtälössä on kääntematriisi $(X'X)^{-1}$, niin korkea selittävien muuttujien multikollineaarisuuden aste lisää estimaattorin b_{PNS} etäisyyttä sen teoreettisesta arvosta β , jolloin estimaattorit $b_{jPNS}, j = 1, \dots, p$, saattavat olla itseisarvoltaan kohtuuttoman suuria. Koska selittävien muuttujien korkea multikollineaarisuuden aste epävakauttaa matriisia $(X'X)^{-1}$, niin se epävakauttaa myös estimaattoria b_{PNS} ja estimaattoreilla $b_{jPNS}, j = 1, \dots, p$, saattaa olla väärät etumerkit.

Estimaattorin b_{PNS} varianssi ja samalla sen keskineliövirhe on

$$\text{var}(\mathbf{b}_{PNS}) = \text{MSE}(\mathbf{b}_{PNS}) = \sigma_y^2 \text{tr}(\mathbf{X}'\mathbf{X})^{-1},$$

missä selitettävien muuttujien ollessa multikollineaarisia joudutaan kääntämään lähes singulaarinen matriisi $\mathbf{X}'\mathbf{X}$. Tällöin estimaattorin \mathbf{b}_{PNS} keskineliövirhe saattaa olla kohtuuttoman suuri, koska selittävien muuttujien multikollineaarisuus on kasvattanut estimaattorin \mathbf{b}_{PNS} kokonaisvarianssia kohtuuttomasti.

Sovitteen $\hat{\mathbf{y}}_{PNS} = \mathbf{X}\mathbf{b}_{PNS}$ korrelaatiomatriisi on $\mathbf{b}_{PNS}'\mathbf{X}'\mathbf{X}\mathbf{b}_{PNS}$. Mitä enemmän selittävät muuttujat korreloivat sitä enemmän korreloivat myös sovitteet $\hat{y}_{jPNS}, j = 1, \dots, n$.

Kun y_0 on määritelty yhtälöllä (2.23), niin ennusteen \hat{y}_{0PNS} varianssi ja keskineliövirhe on

$$\text{var}(\hat{y}_{0PNS}) = \text{MSE}(\hat{y}_{0PNS}) = \sigma_y^2(1 + 1/n) + \sigma_y^2 \text{tr}(\mathbf{x}_0'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_0).$$

Jos annetun \mathbf{x}_0 multikollineaarisuusyhtälö (2.38) on vastaava kuin matriisin \mathbf{X} multikollineaarisuusyhtälö (2.38), niin annetulla \mathbf{x}_0 voidaan melko tarkasti ennustaa tuntematonta y_0 selittävien muuttujien multikollineaarisuudesta huolimatta. Kun annettujen $\mathbf{x}_{01}, \dots, \mathbf{x}_{0p}$ multikollineaarisuuden aste on pienempi kuin muuttujien $\mathbf{x}_1, \dots, \mathbf{x}_p$ multikollineaarisuuden aste, niin tämä kasvattaa ennusteen \hat{y}_{0PNS} keskineliövirhettä.

Selittävien muuttujien multikollineaarisuuden asteen ollessa korkea pienimmän neliösumman estimaattoreiden $b_{jPNS}, j = 1, \dots, p$, varianssit ja kovarianssit saattavat olla siis kohtuuttoman suuria. Tällöin regressiokertoimien $\beta_j, j = 1, \dots, p$, pienimmän neliösumman estimaatit ovat voineet liukua kauaksi pois niiden todellista arvoista ja estimaattorin \mathbf{b}_{PNS} harhattomuudella ei ole enää merkitystä. Täten kun selittävät muuttujat ovat multikollineaarisia, niin yksittäisten regressiokertoimen β_1, \dots, β_p pienimmän neliösumman estimaateilla ei voida tehdä mitään sisällöllisiä tulkintoja regressiomallin rakenteesta (Leskinen, 1977). Geometrisesti tulkittuna selittävien muuttujien multikollineaarisuuden vaikutus parametrien estimointiin nähdään selvästi: jos kaksi selittävää muuttujaa korreloivat voimakkaasti, niin niiden sirontakuvioiden on litteä ellipsi (Ranta, 1991). Regressiotason sovittaminen tapahtuu tällöin suppealle alueelle keskittyneiden havaintojen perusteella. Tällainen kovin kapealle pohjautuva estimointi on epätarkkaa. Lisäksi nämä korreloineet muuttujat liikkuvat yhdessä. Säännöllistämismenetelmillä voidaan saada parametrin β estimaattori, joka on keskimäärin tarkempi kuin estimaattori \mathbf{b}_{PNS} selittävien muuttujien ollessa multikollineaarisia. (Leskinen, 1977; Ranta, 1991.)

3 Harjaregressio

Harjaregressiolla (engl. ridge regression) approksimoidaan rakennematriisin X neliömatriisin käänteismatriisia $(X'X)^{-1}$ matriisilla

$$G_H = (X'X + kI_p)^{-1}, \quad (3.1)$$

missä harjaparametri (engl. ridge parameter) $k \geq 0$ ja neliömatriisi $X'X$ on korrelaatiomuodossa (Hoerl & Kennard, 1970a). Kun harjaregressiossa lisätään tarpeeksi pieni positiivinen luku k neliömatriisin $X'X$ diagonaalille, niin selittävien muuttujien ollessa multikollineaarisia ainoastaan hyvin pienet neliömatriisin $X'X$ ominaisarvot kasvavat ja neliömatriisin $X'X$ kohtuullisilla ja suurilla ominaisarvoilla tulee olemaan lähes sama approksimaatio $\lambda_j \approx \lambda_j + k$. Koska neliömatriisin $X'X + kI_p$ pienimmät ominaisarvot ovat suurempia kuin neliömatriisin $X'X$ pienimmät ominaisarvot, niin matriisi G_H on vakaampi kuin käänteismatriisi $(X'X)^{-1}$ (ks. (2.45)). Täten approksimaatiossa (3.1) harjaparametrin k avulla pyritään vähentämään selittävien muuttujien multikollineaarisuuden vaikutusta neliömatriisissa $X'X$.

Kun regressiomallissa selittävät muuttujat ovat multikollineaarisia ja harjaparametri k on valittu sopivasti approksimaatiossa (3.1), niin harjaregressiossa saadaan regressiokertoimien estimaatit, jotka eivät ole itseisarvoltaan liian suuria, niillä on oikeat etumerkit ja ne ovat vakaampia havaintoaineiston vaihtelun suhteen kuin vastaavat pienimmän neliösumman estimaatit (Leskinen, 1977). Harjaparametri k pyritään valitsemaan approksimaatiossa (3.1) siten, että selittävien muuttujien multikollineaarisuuden asteen ollessa korkea harjaregressiolla saadaan regressiokertoimien β harhainen estimaattori, jonka keskineliövirhe on pienempi kuin estimaattorin b_{PNS} keskineliövirhe.

3.1 Yksimuuttujaiset harjaestimaattorit $b_H(k)$ ja $a_H(k)$

Tarkastellaan regressiomallia (2.5), missä $q = 1$ ja rakennematriisin X neliömatriisi $X'X$ on korrelaatiomatriisi. Ratkaistaan parametrin β harjaestimaattori $b_H(k)$ (engl. ordinary ridge estimator) normaaliyhtälöistä

$$(X'X + kI_p)b_H(k) = X'y \quad (3.2)$$

missä harjaparametri $k \geq 0$. Harjaestimaattori $\mathbf{b}_H(k)$ on

$$\mathbf{b}_H(k) = (\mathbf{X}'\mathbf{X} + k\mathbf{I}_p)^{-1}\mathbf{X}'\mathbf{y} , \quad (3.3)$$

kun $(\mathbf{X}'\mathbf{X} + k\mathbf{I}_p)$ on epäsingulaarinen. Yleistetyksi harjaestimaattoriksi (engl. generalized ridge estimator) kutsutaan estimaattoria, joka ratkaistaan normaaliyhtälöistä

$$(\mathbf{X}'\mathbf{X} + \mathbf{K})\mathbf{b}_H(\mathbf{K}) = \mathbf{X}'\mathbf{y} , \quad (3.4)$$

missä $\mathbf{K} > \mathbf{0}$ ja $\mathbf{K} = \text{diag}(k_1, \dots, k_p)$. (Hoerl & Kennard, 1970a & 1970b.)

Harjaestimaattori esitetään usein kanonisessa muodossa. Oletetaan, että regressiomalli on määritelty kanonisessa muodossa (2.6), missä $q = 1$. Tällöin normaaliyhtälöt (3.2) ovat

$$(\mathbf{V}'\mathbf{X}'\mathbf{T}\mathbf{T}'\mathbf{X}\mathbf{V} + k\mathbf{I}_p)\mathbf{V}'\mathbf{b}_H(k) = \mathbf{V}'\mathbf{X}'\mathbf{T}\mathbf{T}'\mathbf{y} . \quad (3.5)$$

Kun käytetään matriisin \mathbf{X} singulaariarvohajotelmaa (2.8), missä $r = p$, ja lisäksi merkintätapaa $\mathbf{T}'_p\mathbf{y} = \mathbf{u}$, niin yllä olevat normaaliyhtälöt ovat

$$(\Lambda_p + k\mathbf{I}_p)\mathbf{V}'\mathbf{b}_H(k) = \sqrt{\Lambda_p}\mathbf{u} , \quad (3.6)$$

$$\mathbf{V}'\mathbf{b}_H(k) = \mathbf{a}_H(k) = \frac{\sqrt{\Lambda_p}\mathbf{u}}{\Lambda_p + k\mathbf{I}_p} . \quad (3.7)$$

Saatu kanoninen harjaestimaattori $\mathbf{a}_H(k)$ voidaan kirjoittaa muotoon

$$\mathbf{a}_{jH}(k) = \frac{u_j\sqrt{\lambda_j}}{(\lambda_j + k)} , \quad j = 1, \dots, p . \quad (3.8)$$

Kun on laskettu kanoninen harjaestimaattori $\mathbf{a}_H(k)$, niin aina voidaan tehdä muunnos

$$\mathbf{V}\mathbf{a}_H(k) = \mathbf{b}_H(k) , \quad (3.9)$$

jos on tarvetta esittää harjaestimaattori muodossa $\mathbf{b}_H(k)$.

Harjaestimaattorin $\mathbf{a}_H(k)$ ja pienimmän neliösumman estimaattorin \mathbf{a}_{PNS} välisestä yhteydestä saadaan harjaestimaattorille $\mathbf{a}_H(k)$ vaihtoehtoinen esitys

$$a_{jH}(k) = \frac{\sqrt{\lambda_j} \sqrt{\lambda_j} u_j}{\sqrt{\lambda_j} (\lambda_j + k)} = \frac{\lambda_j}{\lambda_j + k} a_{jPNS} = f_j a_{jPNS}, \quad j = 1, \dots, p, \quad (3.10)$$

missä $a_{jPNS} = u_j/\lambda_j^{1/2}$ ja $f_j = \lambda_j/(\lambda_j + k)$, kun $r = p$. Harjaestimaattori $a_{jH}(k)$ on eräänlainen estimaattorin a_{jPNS} konvekssi kombinaatio

$$a_{jH}(k) = w_j 0 + (1 - w_j) a_{jPNS}, \quad (3.11)$$

missä w_j on

$$w_j = \frac{k}{\lambda_j + k}. \quad (3.12)$$

(Brown, 1993). Jos $r < p$, niin $\lambda_{r+1} = \dots = \lambda_p = 0$ ja täten harjaestimaattoria $a_{jH}(k)$ ei voida esittää muodossa (3.10) ja (3.11), koska estimaattoreissa a_{jPNS} jouduttaisiin jakamaan nolllalla, $j = r + 1, \dots, p$. Kun harjaparametri $k > 0$ on tarpeeksi pieni, niin harjaregressiossa vain hyvin pienet neliömatriisin XX ominaisarvot kasvavat ja suurilla ominaisarvoilla tulee olemaan lähes sama approksimaatio $\lambda_j \approx \lambda_j + k$. Kun $k > 0$, niin $|a_{jH}(k)| < |a_{jPNS}|$, $j = 1, \dots, p$. Mitä pienempi on ominaisarvo λ_j , sitä pienempi on f_j ja sitä suurempi on w_j , missä f_j ja w_j ovat välillä $[0, 1]$, ja sitä enemmän tulee estimaattori a_{jPNS} kutistetuksi kohti origoa estimaattorissa $a_{jH}(k)$. Täten mitä suurempi on pienimmän neliösumman estimaattorin varianssi $\text{var}(a_{jPNS}) = 1/\lambda_j$, sitä enemmän harjaestimaattori $a_{jH}(k)$ kutistaa sitä. Niinpä harjaregressiota kutsutaan kutsitavaksi menetelmäksi (engl. shrinking method) ja harjaestimaattori $a_{jH}(k)$ on kutistavaksi estimaattoriksi (engl. shrinking estimator) (Brown, 1993).

3.2 Harjaestimaattoreiden $b_H(k)$ ja $a_H(k)$ keskineliövirhe

Harjaregressiossa tehtävänä on määrätä harhainen estimaattori $a_H(k)$, joka on MSE-kriteerin mukaan parempi estimaattori kuin harhaton estimaattori a_{PNS} , jolloin $\text{MSE}(a_H(k)) < \text{MSE}(a_{PNS})$. Harjaestimaattorin $a_{jH}(k)$ harha on

$$E(a_{jH}(k) - \alpha_j) = \frac{\lambda_j \alpha_j}{\lambda_j + k} - \alpha_j = \frac{\lambda_j \alpha_j - (\lambda_j \alpha_j + k \alpha_j)}{\lambda_j + k} = -\frac{k}{\lambda_j + k} \alpha_j \quad (3.13)$$

eli estimaattorin $a_{jH}(k)$ harha on painon w_j vastaluku yhtälössä (3.11). Harjaestimaattorin $a_{jH}(k)$ varianssi on

$$\text{var}(a_{jH}(k)) = \text{var}\left(\frac{\lambda_j a_{jPNS}}{\lambda_j + k}\right) = \frac{\lambda_j^2 \text{var}(a_{jPNS})}{(\lambda_j + k)^2} = \frac{\lambda_j^2 \sigma_y^2}{(\lambda_j + k)^2 \lambda_j} = \sigma_y^2 \frac{\lambda_j}{(\lambda_j + k)^2}, \quad (3.14)$$

missä $\text{var}(a_{jPNS}) = \sigma_y^2 / \lambda_j$ (vrt. yhtälö (2.27)). Harjaestimaattorin $a_{jH}(k)$ varianssi on pienempi kuin estimaattorin a_{jPNS} varianssi, kun harjaparametri $k > 0$ ja

$$\lambda_j(\lambda_j + k)^{-2} < \lambda_j^{-1} \quad (3.15)$$

(Gunst & Mason, 1980).

Harjaestimaattoreiden $b_H(k)$ ja $a_H(k)$ keskineliövirhe on

$$\begin{aligned} \text{MSE}(b_H(k)) &= \sigma_y^2 (\text{tr}(X'X + kI_p)^{-1} - k \text{tr}(X'X + kI_p)^{-2}) + k^2 \beta'(X'X + kI_p)^{-2} \beta \\ &= \gamma_1(k) + \gamma_2(k) = \end{aligned} \quad (3.16)$$

$$\text{MSE}(a_H(k)) = \sigma_y^2 \sum_{j=1}^p \lambda_j (\lambda_j + k)^{-2} + k^2 \sum_{j=1}^p \alpha_j^2 (\lambda_j + k)^{-2},$$

missä komponentti $\gamma_1(k)$ on estimaattoreiden $b_H(k)$ ja $a_H(k)$ kokonaisvarienssi ja komponentti $\gamma_2(k)$ on estimaattoreiden $b_H(k)$ ja $a_H(k)$ harhan neliö. Komponentti $\gamma_1(k)$ on harjaparametrin k funktiona monotonisesti jatkuva ja vähenevä funktio. Komponentti $\gamma_2(k)$ on harjaparametrin k funktiona monotonisesti jatkuva ja kasvava funktio. Komponenttien derivaatat harjaparametrin k lähestyessä nollaa ovat

$$\lim_{k \rightarrow 0^+} \left(\frac{d\gamma_1}{dk} \right) = -2\sigma_y^2 \sum_{j=1}^r \left(\frac{1}{\lambda_j} \right) \quad (3.17)$$

ja

$$\lim_{k \rightarrow 0^+} \left(\frac{d\gamma_2}{dk} \right) = 0 . \quad (3.18)$$

Kun selittävät muuttujat ovat riippumattomia ja harjaparametri k lähestyy nollaa, niin komponentilla γ_1 on negatiivinen derivaatta $-2p\sigma_y^2$. Kun selittävät muuttujat ovat multikollineaarisia ja harjaparametri k lähestyy nollaa, niin komponentin γ_1 derivaatta lähestyy negatiivista ääretöntä. Kun $k \rightarrow 0^+$, niin komponentti γ_2 on vakaa ja nolla origossa. Nämä edellä esitetyt komponenttien γ_1 ja γ_2 ominaisuudet mahdollistavat, että kun

$$0 < k < \frac{\sigma_y^2}{\max_{(1 \leq j \leq p)} \alpha_j^2} , \quad (3.19)$$

niin harjaregressiolla saadaan harhainen estimaattori $a_H(k)$, jolla on pienempi keskineliövirhe kuin harhattomalla estimaattorilla a_{PNS} . (Hoerl & Kennard, 1970a & 1970b.)

Kun

$$k_j = \sigma_y^2 / \alpha_j^2 , \quad j = 1, \dots, p , \quad (3.20)$$

niin harjaestimaattorilla $a_H(k)$ on pienin mahdollinen keskineliövirhe. Yleensä parametrit α ja σ_y^2 ovat tuntemattomia, joten harjaparametrin k yläraja on teoriassa tuntematon. Kuitenkin käytännössä harjaparametrin k yläraja on äärellinen, sillä on aina olemassa harjaparametri k , jolla

$$\text{MSE} (b_H(k)) = \text{MSE} (a_H(k)) < \text{MSE} (b_{PNS}) = \text{MSE} (a_{PNS}) ,$$

kun harjaparametri k on

$$0 < k < \frac{\hat{\sigma}_y^2}{\max_{(1 \leq j \leq p)} a_{jPNS}^2} . \quad (3.21)$$

(Hoerl & Kennard, 1970a & 1970b.)

3.3 Harjaregression geometrinen tulkinta

Olkoon a mikä tahansa parametrin α estimaateista ja tehtävänä on minimoida estimaattorin a neliöity pituus $a'a$ ehdolla

$$\sum_{j=1}^r (a_j - a_{jPNS})^2 \lambda_j = \text{vakio} . \quad (3.22)$$

Langrangen lauseke tälle minimointitehtävälle on

$$\sum_{j=1}^r a_j^2 + (1/k) \sum_{j=1}^r (a_j - a_{jPNS})^2 \lambda_j . \quad (3.23)$$

Kun derivoidaan tämä lauseke estimaattorin a_j suhteen ja asetetaan derivoitu lauseke nolllaksi, niin saadaan

$$a_j = \frac{\lambda_j a_{jPNS}}{\lambda_j + k} , \quad (3.24)$$

joka on harjaestimaattori $a_{jH}(k)$. Lisäksi

$$\sum_{j=1}^r a_{jH}^2(k) = \frac{\sum_{j=1}^r a_{jPNS}^2 \lambda_j^2}{(\lambda_j + k)^2} = f(k) , \quad (3.25)$$

missä $f(k)$ on harjaparametrin k funktiona monotoninen ja

$$f(0) = \sum_{j=1}^r a_{jPNS}^2 \quad \text{ja} \quad f(\infty) = 0 . \quad (3.26)$$

Harjaregression ratkaisu $\mathbf{a}_H(k)$ voidaan pelkistää säteeksi origosta pisteessä, jossa se kohtaa pienimmän neliösumman ratkaisun. Pienimmän neliösumman estimaattorin jäännösneliösumma on parametrin α funktiona ellipsoidi, jonka keskipiste on \mathbf{a}_{PNS} ja pääakseleina ovat neliömatriisin $\mathbf{X}\mathbf{X}$ ominaisarvojen neliöjuuret. Tämä ellipsoidi tulee olemaan hyvin kapea multikollineaarisessa tapauksessa, koska tällöin neliömatriisin $\mathbf{X}\mathbf{X}$ pienin ominaisarvo $\lambda_p \approx 0$. Kun harjaparametri k on nolla, niin estimaattorin $\mathbf{a}_H(k)$ etäisyys origosta on sama kuin estimaattorin \mathbf{a}_{PNS} etäisyys origosta ja harjaregression jäännösneliösumman ellipsoidi ei ole pyörästynyt pienimmän neliösumman jäännösneliösumman ellipsoidista, eli harjaregressio on sama kuin pienimmän neliösumman regressio. Kun harjaparametri k lähestyy ääretöntä, niin estimaattorin $\mathbf{a}_H(k)$ etäisyys origosta on rajattu hyvin pieneksi, jolloin estimaattori $\mathbf{a}_H(k) = 0$. Koska parametri α sijaitsee säteellä $[0, \mathbf{a}_{PNS}]$, niin harjaregressiossa tehtävänä on harjaparametrin k avulla kutistaa estimaattoria \mathbf{a}_{PNS} kohti origoa ja pyöristää tämän jäännösneliösumman ellipsoidia siten, että sen keskipiste, joka on nyt siis harjaestimaattori $\mathbf{a}_H(k)$, on mahdollisimman lähellä parametria α . Kun selittäviä muuttujia on kaksi, niin pienimmän neliösumman jäännösneliösumman muoto on harjanne, jonka mukaan harjaregressio on saanut nimensä.

3.4 Ennusteet harjaregressiossa

Oletetaan, että tuntematon y_0 voidaan kirjoittaa muotoon (2.23). Tällöin ennusteen $\hat{y}_{0H} = \mathbf{x}_0' \mathbf{b}_H(k)$ varianssi on

$$\text{var}(\hat{y}_{0H}) = \sigma_y^2(1 + 1/n + \mathbf{x}_0' \mathbf{G}_H \mathbf{x}_0) = \sigma_y^2(1 + 1/n + \mathbf{x}_0' (\mathbf{X}'\mathbf{X} + k\mathbf{I}_p)^{-1} \mathbf{x}_0), \quad (3.27)$$

missä $\mathbf{G}_H = (\mathbf{X}\mathbf{X} + k\mathbf{I}_p)^{-1}$ (Brown, 1993).

Oletetaan, että regressiomalli on määritelty kanonisen yhtälön (2.6) mukaan, ja tehtävänä on matriisilla $\Lambda_p^{1/2}$ ennustaa tuntematonta \mathbf{u} . Tällöin ennusteen $\hat{\mathbf{u}}_H$ keskineliövirhe on

$$\text{MSE}(\hat{\mathbf{u}}_H(k)) = \sigma_y^2(n + 1) + E \sum_{j=1}^p \lambda_j (\mathbf{a}_{jH}(k) - \alpha_j)^2 \quad (3.28)$$

$$\begin{aligned}
&= \sigma_y^2(n+1) + \sum_{j=1}^p \lambda_j (\text{var}(a_{jH}(k)) + B^2(a_{jH}(k))) \\
&= \sigma_y^2(n+1) + \sigma_y^2 \sum_{j=1}^p \lambda_j^2 / (\lambda_j + k)^2 + k^2 \sum_{j=1}^p \lambda_j \alpha_j^2 / (\lambda_j + k)^2,
\end{aligned}$$

missä $B^2(a_{jH}(k))$ on estimaattorin $a_{jH}(k)$ harhan neliö. Wang ja Chow (1994) ovat osoittaneet, että aina on olemassa harjaparametrin k arvo, jolla $\text{GMSE}(a_H(k)) < \text{GMSE}(a_{PNS})$, kun

$$0 < k < 2\sigma_y^2 / \alpha' \alpha. \quad (3.29)$$

Täten $\text{MSE}(a_H(k)) < \text{MSE}(a_{PNS})$ ja $\text{MSE}(\hat{u}_H) < \text{MSE}(\hat{u}_{PNS})$, kun harjaparametrille k on voimassa yllä oleva epäyhtälö.

3.5 Harjaparametrin k valintamenetelmiä yksimuuttujaisessa tapauksessa

Kun parametrit α ja σ_y^2 ovat tuntemattomia, niin ei ole mahdollista määrätä selvää automaattista estimointimenetelmää harjaparametrille k . Seuraavassa on esitetty muutamia menetelmiä harjaparametrin k valinnalle.

1) Harjajälki: Piirretään harjaestimaattori $a_H(k)$ harjaparametrin $k \geq 0$ funktiona. Kyseistä kuvaa kutsutaan harjajäljeksi (engl. ridge trace). Piirretään kuvaan myös sovitteen jäännöseliösumma harjaparametrin k funktiona. Kuvan perusteella valitaan se harjaparametrin k arvo, jolla estimaattoreiden $a_{jH}(k)$, $j = 1, \dots, r$, arvot näyttävät vakiintuneen ja sovitteen jäännöseliösumma ei ole kasvanut kohtuuttoman suureksi. Jos matriisin X aste r on pienempi kuin selittävien muuttujien lukumäärä p , niin $\lambda_{r+1} = \dots = \lambda_p = 0$, jolloin estimaattoria $a_{jH}(k)$, missä $k = 0$, ei ole määritelty, $j = r+1, \dots, p$. Kun harjaparametri k määrätään graafisesti, niin sopivan harjaparametrin k arvon määrääminen perustuu tutkijan subjektiiviseen käsitykseen sopivasta harjaparametrin k arvosta, mikä vaatii tämän menetelmän huolellista ja ammattitaitoista käyttöä. (Hoerl & Kennard, 1970a & 1970b.)

2) *Cross-validation*: Kullakin annetulla harjaparametrin k arvolla jätetään harjaestimaattorin $b_H(k)$ laskennassa vuorollaan kukin havainto $i = 1, \dots, n$ pois ja lasketaan jokaiselle poistetulle havainnolle i poistetun havainnon $y_{(i)}$ ennuste. Lasketaan lopuksi

$$\text{PRESS}(k) = \sum_{i=1}^n (y_i - \hat{y}_{(i)})^2, \quad (3.30)$$

missä siis ennuste $\hat{y}_{(i)}$ on laskettu harjaregressiolla ilman havaittua havainto i . Valitaan se harjaparametrin k estimaattori k_{cross} , jolla ennusteen virheneliösumma $\text{PRESS}(k)$ on pienin. (Brown, 1993.)

3) *Generalized cross-validation*: Valitaan harjaparametrin k estimaattoriksi k_{gcross} , joka minimoi lausekkeen

$$\|(I - H(k))y\|^2 / (\text{tr}(I - H(k)))^2, \quad (3.31)$$

missä $H(k) = X(X'X + kI_p)^{-1}X'$ (Brown, 1993).

4) Hoerlin, Kennardin ja Baldwin säännön mukaan:

$$k_{\text{HKB}} = (r - 2) \hat{\sigma}_y^2 / \mathbf{a}'_{\text{PNS}} \mathbf{a}_{\text{PNS}}. \quad (3.32)$$

Tämä harjaparametrin k valintamentelmä k_{HKB} perustuu epäyhtälöön (3.19), missä parametreja σ_y^2 ja α on estimoitu niiden pienimmän neliösuman estimaattoreilla. Estimaattori k_{HKB} saattaa olla liian pieni, kun selittävien muuttujien x_{pj} , $j = 1, \dots, p$, kollineaarisuuden indeksi (2.43) on suuri ja neliömatriisin $X'X$ pienimpiä ominaisarvoja vastaavat neliömatriisin $X'X$ ominaisvektorit ovat korreloimattomia havaintovektorin y kanssa. Jos

$$h_p^2(p + 2) < 2p \sum_{j=1}^p h_j^2/p, \quad (3.34)$$

missä $h_j = \lambda_j^{-1}$, niin $\text{MSE}(\mathbf{a}_H(k_{\text{HKB}})) < \text{MSE}(\mathbf{a}_{\text{PNS}})$ ja jos

$$h_p(p + 2) < 2p \sum_{j=1}^p h_j/p, \quad (3.35)$$

niin $\text{MSE}(\hat{\mathbf{a}}_H(k_{\text{HKB}})) < \text{MSE}(\hat{\mathbf{a}}_{\text{PNS}})$. (Brown, 1993.)

5) Brown (1993):

$$k_{\text{MLW}} = (r - 2) \hat{\sigma}_y^2 \text{tr}(X'X) / (r \mathbf{b}'_{\text{PNS}} X'X \mathbf{b}_{\text{PNS}}). \quad (3.36)$$

Kun selittävien muuttujien $x_j, j = 1, \dots, p$, kollineaarisuudenindeksi on suuri ja neliömatriisin XX pienimpiä ominaisarvoja vastaavat neliömatriisiin XX ominaisvektorit ovat korreloimattomia havaintovektorin y kanssa, niin k_{MLW} voidaan pitää luotettavampana kuin k_{HKB} . Jos ominaisarvo λ_j on pieni, niin estimaattoria a_{jPNS} , joka voi olla itseisarvoltaan kohtuuttoman suuri, painotetaan pienellä ominaisarvolla λ_j harjaestimaattorissa $a_H(k_{MLW})$. (Brown, 1993.)

6) Brown (1993): Lauseke

$$\hat{\sigma}_y^2 \sum_{j=1}^p (\lambda_j - k) / (\lambda_j (\lambda_j + k)) + k^2 \sum_{j=1}^p a_{jPNS}^2 / (\lambda_j + k)^2 \quad (3.37)$$

eli keskineliövirheen harhaton estimaattori minimoidaan harjaparametrin k suhteen, mistä saadaan harjaparametrin k valinta k_{MUR} .

7) Wang ja Chow (1994): Yleistetyyn harjaestimaattoriin $a_H(k)$ harjaparametrivektori k voidaan valinta

$$k_{jad} = \hat{\sigma}_y^2 / a_{jPNS}^2 \quad j = 1, \dots, r \quad (3.38)$$

Tämä harjaparametrivektorin k_{ad} valintamenetelmä perustuu epäyhtälöön (3.20), missä parametreja σ_y^2 ja α_j on estimoitu niiden pienimmän neliösumman estimaattoreilla. Jos tällä harjaparametrin k valintamenetelmällä halutaan valita harjaestimaattorin $b_H(k)$ harjaparametrivektori k , niin korvataan tässä valintamenetelmässä estimaattori a_{jPNS} estimaattorilla b_{jPNS} .

3.6 Monimuuttujainen harjaregressio

Olkoon monimuuttujainen regressiomalli määritelty yhtälön (2.5) mukaan. Kirjoitetaan regressiomatriisi B vektorimuotoon β^*

$$\beta^* = (\beta_1^*, \dots, \beta_p^*)' \quad (3.39)$$

missä β_j^* on $(1 \times q)$ -vektori ja se on regressioparametrimatriisin B j :nnes rivivektori. Brownin ja Zidekin (1980) ehdottama parametrin β^* harjaestimaattori $b_H^*(K)$ on

$$\mathbf{b}^*(\mathbf{K}) = (\mathbf{X}'\mathbf{X} \otimes \mathbf{I}_q + \mathbf{I}_p \otimes \mathbf{K})^{-1}(\mathbf{X}'\mathbf{X} \otimes \mathbf{I}_q)\mathbf{b}_{PNS}^* , \quad (3.40)$$

missä \otimes on Kroneckerin tulo, $\mathbf{K} > 0$ on $(q \times q)$ -harjaparametrimatriisi, \mathbf{b}_{PNS}^* on parametrin β^* pienimmän neliösumman estimaattori. Kun regressiomalli on määritelty kanonisen yhtälön (2.6) mukaan, niin monimuuttujaisen harjaestimaattorin $\mathbf{b}_{H}^*(\mathbf{K})$ kanoninen muoto $\mathbf{a}_{H}^*(\mathbf{K})$ on

$$\mathbf{a}^*(\mathbf{K}) = (\Lambda_p \otimes \mathbf{I}_q + \mathbf{I}_p \otimes \mathbf{K})^{-1}(\Lambda_p \otimes \mathbf{I}_q)\mathbf{a}_{PNS}^* , \quad (3.41)$$

missä \mathbf{a}_{PNS}^* on estimaattorin \mathbf{b}_{PNS}^* kanoninen muoto. Harjaestimaattori $\mathbf{a}_{H}^*(\mathbf{K})$ voidaan esittää myös muodossa

$$\mathbf{a}_{jH}^*(\mathbf{K}) = \mathbf{a}_{jPNS}^* \lambda_j (\lambda_j \mathbf{I}_q + \mathbf{K})^{-1} , \quad j = 1, \dots, p . \quad (3.42)$$

Kuten yksimuuttujainen harjaestimaattori $\mathbf{a}_{jH}(k)$, niin myös monimuuttujainen harjaestimaattori $\mathbf{a}_{jH}^*(\mathbf{K})$ on painotettu \mathbf{a}_{jPNS}^*

$$\mathbf{a}_{jH}^*(\mathbf{K}) = \mathbf{W}_j + (\mathbf{I}_q - \mathbf{W}_j(\mathbf{K}))\mathbf{a}_{jPNS}^* , \quad j = 1, \dots, p , \quad (3.43)$$

missä \mathbf{W}_j on

$$\mathbf{W}_j(\mathbf{K}) = \mathbf{K} (\lambda_j \mathbf{I}_q + \mathbf{K})^{-1} . \quad (3.44)$$

Monimuuttujaisessa harjaregressiossa kutistetaan estimaattoria \mathbf{a}_{PNS}^* kohti origoa ja sitä enemmän estimaattori \mathbf{a}_{jPNS}^* kutistuu, mitä suurempi on paino $\mathbf{W}(\mathbf{K})$ eli mitä pienempi on neliömatriisin $\mathbf{X}'\mathbf{X}$ ominaisarvo λ_j . Harjaestimaattorin $\mathbf{a}_H^*(\mathbf{K})$, missä $\mathbf{K} = k\mathbf{I}_q$, keskineliövirhe on

$$\text{MSE}(\mathbf{a}_H^*) = \text{tr}(\Sigma_Y) \sum_{j=1}^p \lambda_j / (\lambda_j + k)^2 + \sum_{j=1}^p k^2 / (\lambda_j + k)^2 \sum_{i=1}^q \alpha_j^{(i)2} , \quad (3.45)$$

missä $\Sigma_Y = (\sigma_Y^{(ij)})$, $i, j = 1, \dots, q$. Kun yksimuuttujaista harjaestimaattoria $\mathbf{a}_H(k)$ ja sen keskivirhettä $\text{MSE}(\mathbf{a}_H(k))$ verrataan monimuuttujaisen harjaestimaattoriin $\mathbf{a}_H^*(\mathbf{K})$ ja sen keskivirheeseen $\text{MSE}(\mathbf{a}_H^*(\mathbf{K}))$, niin estimaattoreilla $\mathbf{a}_H^*(\mathbf{K})$ ja $\mathbf{a}_H(k)$ ja täten myös keskivirheillä $\text{MSE}(\mathbf{a}_H^*(\mathbf{K}))$ ja

$MSE(\mathbf{a}_H(k))$ ovat samat matemaattiset ominaisuudet. Siis kun $q = 1$ ja $\mathbf{K} = k$, niin $\mathbf{a}_H^*(\mathbf{K}) = \mathbf{a}_H(k)$ ja $MSE(\mathbf{a}_H^*(\mathbf{K})) = MSE(\mathbf{a}_H(k))$. (Brown & Payne, 1975; Broewn & Zidek, 1980.)

3.7 Harjaparametrimatriisin \mathbf{K} valinta monimuuttujaisessa tapauksessa

Parametrin \mathbf{B} monimuuttujainen harjaestimaattori voidaan aina hajottaa q yksimuuttujaiseen harjaestimaattoriin

$$\mathbf{b}_H^{(j)}(k^{(j)}) = (\mathbf{X}'\mathbf{X} + k^{(j)})^{-1}\mathbf{X}'\mathbf{y}^{(j)}, \quad j = 1, \dots, q, \quad (3.46)$$

missä harjaparametrin $k^{(j)}$ ratkaisu riippuu ainoastaan yhdestä selitettävästä muuttujasta $y^{(j)}$. Tällöin harjaparametrit $k^{(j)}, j = 1, \dots, q$, voidaan ratkaista kappaleessa 3.5 esitettyjen harjaparametrien k valintamenetelmien mukaan. Tämä saattaa kuitenkin olla epätyydyttävä monimuuttujaisessa tapauksessa, koska siinä ei hyödynnetä kovarianssimatriisia Σ_Y eli kaikkien q normaaliyhtälöiden yhteisinformaatiota.

Monimuuttujaisessa harjaregressiossa voidaan hyödyntää kovarianssimatriisia Σ_Y valittaessa harjaparametrimatriisia \mathbf{K} . Seuraavassa on eri tapoja valita harjaparametrimatriisi \mathbf{K} monimuuttujaisessa harjaregressiossa.

3.7.1 Harjaparametrimatriisin \mathbf{K} valinta Brownin ja Paynen (1975) mukaan

Oletetaan, että regressiomallina on malli (2.6), missä Σ_Y on tunnettu ja $\alpha_i^{(j)}$ noudattaa jakaumaa $N(0, \sigma_\alpha^2)$, $i = 1, \dots, p, j = 1, \dots, q$. Kun harjaparametrimatriisi \mathbf{K} on valittu Brownin ja Paynen (1975) mukaan, niin harjaestimaattori $A_H(\mathbf{K}_{BP})$ on

$$a_{iH}^{(j)}(k_{BP}^{(j)}) = \frac{\lambda_i a_{iPNS}^{(j)}}{\lambda_i + h^{(j)} \sigma_\alpha^2} = \frac{\lambda_i a_{iPNS}^{(j)}}{\lambda_i + k_{BP}^{(j)}}, \quad i = 1, \dots, p, \quad j = 1, \dots, q, \quad (3.47)$$

missä $h^{(j)}$ on kovarianssimatriisin Σ_Y ominaisarvo ja harjaparametri $k^{(j)} = h^{(j)} \sigma_\alpha^2, j = 1, \dots, q$. Jos ominaisarvo λ_i on pieni, niin erityisesti tällöin estimaattori $a_{iH}^{(j)}(k^{(j)})$ on kutistanut enemmän estimaattoria $a_{iPNS}^{(j)}$, missä kovarianssimatriisin Σ_Y ominaisarvo $h^{(j)}$ on suuri, kohti origoa kuin estimaattoria $a_{iPNS}^{(k)}$, missä kovarianssimatriisin Σ_Y ominaisarvo $h^{(k)}$ on pieni. Kun kovarianssimatriisi Σ_Y on diagonaalinen, niin tällöin kovarianssimatriisilla Σ_Y on vain yksi nollasta eroava

ominaisarvo ja $\mathbf{K} = k\mathbf{I}_q$, missä harjaparametri k on

$$k = \frac{h}{\sigma_\beta^2} \quad (3.48)$$

ja h on ainoa matriisin Σ_Y nolasta eroava ominaisarvo. Jos oletetaan, että matriisi Σ_Y on diagonaalinen, niin tällöin monimuuttujaisessa harjaregressiossa ei hyödynnetä q normaaliyhtälöiden yhteisinformaatiota. (Brown & Payne, 1975.)

3.7.2 Harjaparametrimatriisin \mathbf{K} valinta Brownin ja Zidekin (1982) mukaan

Oletetaan, että regressiomalli on määritelty kanonisella yhtälöllä (2.6), missä $r = p$, $E(\epsilon^{(i)}) = 0$, $\text{cov}(\epsilon^{(i)}, \epsilon^{(j)}) = \sigma_Y^{(ij)}\mathbf{I}_p$, $i, j = 1, \dots, q$. Kun parametri \mathbf{A} noudattaa jakaumaa $N(0, \Sigma_\alpha)$, niin yleinen muoto harjaparametrimatriisille \mathbf{K} on

$$\mathbf{K} = \Sigma_Y \Sigma_\alpha^{-1}. \quad (3.49)$$

Jos parametrit Σ_Y ja Σ_α ovat tuntemattomia, niin käytetään estimaattoreita

$$\hat{\Sigma}_Y^{-1} = \mathbf{R} \quad \text{ja} \quad \hat{\Sigma}_\alpha = \sum_{j=1}^p w_j \mathbf{a}_{jPNS}' \mathbf{a}_{jPNS}, \quad (3.50)$$

jossa \mathbf{R}^{-1} on matriisin Σ_Y suurimman uskottavuuden estimaattori tai pienimmän neliösumman estimaattori.

Merkitään $v_j = \lambda_j^{-1}$, $\lambda_j, j = 1, \dots, p$, ovat matriisin $\mathbf{X}'\mathbf{X}$ ominaisarvoja, jotka ovat suuruusjärjestyksessä $\lambda_1 \geq \dots \geq \lambda_p, f = p - q - 1$,

$$\bar{v}^{-1} = \sum_{j=1}^p v_j^{-1}/p \quad \text{ja} \quad \bar{v}^2 = \sum_{j=1}^p v_j^2/p. \quad (3.51)$$

Kun harjaparametrimatriisi \mathbf{K} valitaan Scloven säännön mukaan, niin monimuuttujainen harjaestimaattori (3.42) on

$$\mathbf{a}_{jH}^*(\mathbf{K}_S) = \mathbf{a}_{jPNS}^* (\mathbf{I}_q - v_j f \bar{v}^{-1} (v_j f \bar{v}^{-1} \mathbf{I}_q + \mathbf{R} \sum_{i=1}^p v_i^{-1} \mathbf{a}_{iPNS}' \mathbf{a}_{iPNS})^{-1}), \quad (3.52)$$

$j = 1, \dots, p$. Monimuuttujaisella harjaestimaattorilla $\mathbf{a}_{H(K_S)}^*$ on pienempi keskineliövirhe kuin pienimmän neliösumman estimaattorilla \mathbf{a}_{PNS}^* kun

$$(n - p)(p - q - 1)\bar{v}^{-1}v_p^3 - 2(n - p - 2(q + 1))(p\bar{v}^2 - (q + 1)v_p^2) < 0 \quad (3.53)$$

ja

$$n - p - 2(q + 1) > 0. \quad (3.54)$$

Kun harjaparametrimatriisi \mathbf{K} valitaan Hoerlin, Kennardin ja Baldwinin säännön mukaan, niin monimuuttujainen harjaestimaattori (3.42) on muotoa

$$\mathbf{a}_{jH(\mathbf{K}_{HKB})}^* = \mathbf{a}_{jPNS}^*(\mathbf{I}_q - f\mathbf{v}_j(f\mathbf{v}_j\mathbf{I}_q + R\sum_{i=1}^p \mathbf{a}_{iPNS}^* \mathbf{a}_{iPNS}^{*'})^{-1}), \quad (3.55)$$

$j = 1, \dots, p$. Monimuuttujaisella harjaestimaattorilla $\mathbf{a}_{H(K_{HKB})}^*$ on pienempi keskineliövirhe kuin pienimmän neliösumman estimaattorilla \mathbf{a}_{PNS}^* kun

$$n - p - 2q - 2 > 0 \quad (3.56)$$

ja

$$(n - p)(p - q - 1)v_p^2 - 2(n - p - 2q - 2)(p\bar{v}^2 - (q + 1)v_p^2) < 0. \quad (3.57)$$

3.7.3 Harjaparametrimatriisin \mathbf{K} valinta Frankin ja Friedmanin (1993) mukaan

Kun harjaparametrimatriisi \mathbf{K} valitaan Frankin ja Friedmanin (1993) mukaan, niin kanonisen monimuuttujaisen regressiomallin (2.12) harjaestimaattori on

$$\mathbf{a}_{iH(k_{FF}^{(j)})}^{(j)} = \mathbf{a}_{iPNS}^{(j)} \lambda_i / (\lambda_i + k_{FF}^{(j)}), \quad i = 1, \dots, p \text{ ja } j = 1, \dots, q, \quad (3.58)$$

missä $\lambda_i, i = 1, \dots, p$, ovat neliömatriisin $\mathbf{X}\mathbf{X}$ ominaisarvoja ja $k_{FF}^{(j)}$ on

$$k_{PP}^{(j)} = p \sigma_Y^2 / n (h^{(j)} - \sigma_Y^2), \quad (3.59)$$

$h^{(j)}, j = 1, \dots, q$, ovat neliömatriisin $Y'Y$ ominaisarvoja, n on otoskoko ja $\Sigma_Y = \sigma_Y^2 I_q$. Jos σ_Y^2 on tuntematon, niin käytetään tässä sen estimaattoria

$$\hat{\sigma}_Y^2 = \frac{\sum_{j=1}^q (y^{(j)} - b_{PNS}^{(j)})'(y^{(j)} - b_{PNS}^{(j)})}{q(n-p)}. \quad (3.60)$$

Mitä pienemmät ovat ominaisarvot λ_i ja $h^{(j)}$ sitä enemmän kutistetaan estimaattoria $\alpha_{i, PNS}^{(j)}$ harjaestimaattorissa $\alpha_{i, H}^{(j)}(k_{PP}^{(j)})$ kohti origoa. (Frank & Friedman, 1993.)

3.7.4 Harjaparametrimatriisin K valinta *cross-validation* -menetelmällä

Sovelletaan yksimuuttujaisen harjaestimaattorin $a_H(k)$ harjaparametrin k *cross-validation* valintamenetelmää monimuuttujaisen harjaestimaattorin $a_H^*(K)$ harjaparametrimatriisin K valintamenetelmäksi. Oletetaan, että $\Sigma_Y = \sigma_Y^2 I_q$, jolloin $K = kI_q$, missä $k \geq 0$. Lasketaan ennusteen virhevariassi

$$\text{PRESS}(k) = \sum_{i=1}^n \sum_{j=1}^q (u_i^{(j)} - \hat{u}_{(i)}^{(j)})^2, \quad (3.61)$$

missä $\hat{u}_{(i)}^{(j)}$ on ennuste, joka on saatu harjaregressiolla ilman havaintoa $i, i = 1, \dots, n$. Valitaan se harjaparametrin k estimaattori k_{cross} , jolla $\text{PRESS}(k) = \text{PRESS}(k_{cross})$ on pienin.

3.7.5 Harjaparametrimatriisin K valinta Fülen (1995) mukaan

Kirjoitetaan regressiomalli (2.5) tässä muotoon

$$\text{vek}(Y) = (I_q \otimes X) \text{vek}(B) + \text{vek}(E), \quad (3.62)$$

missä \otimes on Kroneckerin tulo. Fülen (1995) mukaan parametrin $\text{vek}(B)$ harjaestimaattori on

$$\begin{aligned} \text{vek}(B(K)) &= (D'D + K)^{-1} (D' \text{vek}(Y) + K \text{vek}(B_0)) \\ &= (D'D + K)^{-1} (D'D \text{vek}(B_{PNS}) + K \text{vek}(B_0)), \end{aligned} \quad (3.63)$$

missä $\text{vek}(\mathbf{B}_0)$ on valittu $(pq \times 1)$ -vektori, \mathbf{K} on $(pq \times pq)$ -harjaparametrimatriisi ja $\mathbf{D} = \mathbf{I}_q \otimes \mathbf{X}$. Tämä harjaestimaattori $\text{vek}(\mathbf{B}_H(\mathbf{K}))$ kutistaa pienimmän neliösumman estimaattoria $\text{vek}(\mathbf{B}_{PNS})$ kohti valittua vektoria $\text{vek}(\mathbf{B}_0)$, joka on oletettu estimoitavan parametrin $\text{vek}(\mathbf{B})$ todellinen arvo. Koska estimoitava parametri $\text{vek}(\mathbf{B})$ on välillä $[0, \text{vek}(\mathbf{B}_{PNS})]$, niin

$$\text{vek}(\mathbf{B}_{j0}) = c \text{vek}(\mathbf{B}_{jPNS}), \quad j = 1, \dots, pq, \quad (3.64)$$

missä $0 < c < 1$ eli täten valitaan $\text{vek}(\mathbf{B}_0)$ väliltä $[0, \text{vek}(\mathbf{B}_{PNS})]$. Fülen (1995) ehdottama harjaparametrimatriisin \mathbf{K} valinta on $\mathbf{K}_F = \text{diag}(k_{jF})$, missä

$$k_{jF} = \frac{pq \hat{\sigma}_Y^2}{(\text{vek}(\mathbf{B}_{jPNS}) - \text{vek}(\mathbf{B}_{j0}))^2}, \quad j = 1, \dots, pq, \quad (3.65)$$

ja $\hat{\sigma}_Y$ on (3.60). Täten valittuun \mathbf{K}_F vaikuttaa valittu $\text{vek}(\mathbf{B}_0)$.

Tässä harjaparametrimatriisin \mathbf{K} valintamenetelmässä on ongelmana valita sopiva $\text{vek}(\mathbf{B}_0)$. Koska tehtävänä oli muodostaa harhainen estimaattori $\text{vek}(\mathbf{B}_H(\mathbf{K}_{HMSE}^F))$, jonka keskineliövirhe on pienempi kuin estimaattorin $\text{vek}(\mathbf{B}_{PNS})$ keskineliövirhe, niin nollahypoteesina on

$$\text{MSE}(\text{vek}(\mathbf{B}_H(\mathbf{K}_F))) < \text{MSE}(\text{vek}(\mathbf{B}_{PNS}))$$

ja sitä voidaan testata testisuurella

$$F = \frac{(\text{vek}(\mathbf{B}_{PNS}) - \text{vek}(\mathbf{B}_0))'((\mathbf{D}'\mathbf{D})^{-1} + \mathbf{K}_F)^{-1}(\text{vek}(\mathbf{B}_{PNS}) - \text{vek}(\mathbf{B}_0))/v_1}{(\text{vek}(\mathbf{Y}) - \mathbf{D} \text{vek}(\mathbf{B}_{PNS}))'(\text{vek}(\mathbf{Y}) - \mathbf{D} \text{vek}(\mathbf{B}_{PNS}))/v_2}, \quad (3.66)$$

joka noudattaa F-jakumaa vapausastein $v_1 = qp$ ja $v_2 = q(n-p)$. (Füle, 1995.)

4 Pääkomponenttiregressio

Neliömatriisin $X'X$ käänteismatriisi on

$$(X'X)^{-1} = (V_p \Lambda_p V_p')^{-1} = \sum_{j \in S_p} (1/\lambda_j) v_j v_j', \quad (4.1)$$

missä $\Lambda_p = \text{diag}(\lambda_1, \dots, \lambda_p)$, $V_p = (v_1, \dots, v_p)$, $S_p = \{1, \dots, p\}$, λ_j , $j \in S_p$, ovat neliömatriisin $X'X$ ominaisarvot ja v_j , $j \in S_p$, ovat vastaavat ominaisvektorit. Pääkomponenttiregressiossa (engl. principal components regression) tätä approksimoidaan matriisilla

$$G_{PK} = (V_A \Lambda_A V_A')^{-1} = \sum_{j \in S_A} (1/\lambda_j) v_j v_j', \quad (4.2)$$

missä $\Lambda_A = \text{diag}(\lambda_j)$, $V = (v_j)$, $j \in S_A$, S_A on valittujen pääkomponenttien indeksijoukko, A on valittujen pääkomponenttien lukumäärä ja $S_A \in S_p$ (vrt. (3.1)) (Brown, 1993). Kun selittävien muuttujien multikollineaarisuuden aste on korkea, $\lambda_j \neq 0$ ja $j \in S_A$, niin erityisesti tällöin matriisi $V_A \Lambda_A V_A'$ on paremmin kääntyvä kuin matriisimatriisi $V_p \Lambda_p V_p'$.

Kirjoitetaan regressiomalli (2.5), missä selitettäviä muuttujia on yksi eli $q = 1$, kanoniseen muotoon

$$y = X\beta + \epsilon = XVV'\beta + \epsilon = Z\alpha + \epsilon = \alpha_1 z_1 + \dots + \alpha_p z_p + \epsilon, \quad (4.3)$$

missä $V'\beta = \alpha = (\alpha_1, \dots, \alpha_p)'$, $XV = Z = (z_1, \dots, z_p)$, matriisi $V = (v_1, \dots, v_p)$ on ortogonallinen ja v_j , $j \in S_p$, ovat neliömatriisin $X'X$ ominaisvektorit, joiden pituus on skaalattu ykköseksi. Vektoreita z_1, \dots, z_p kutsutaan rakennematriisin X pääkomponenteiksi (engl. principal components) ja ne ovat selittäviä muuttujia pääkomponenttiregressiossa (Jackson, 1991).

Pääkomponenttiregressiossa tehdään siis rakennematriisille X ortogonaalinen muunnos $XV = Z = (z_1, \dots, z_p)$. Pääkomponentit z_j , $j = 1, \dots, p$, ovat selittävien muuttujien x_j , $j = 1, \dots, p$, lineaarikombinaatiota

$$z_j = v_{1j}x_1 + v_{2j}x_2 + \dots + v_{pj}x_p, \quad j = 1, \dots, p. \quad (4.4)$$

Niiden kovarianssimatriisi $\text{cov}(Z, Z)$ on

$$\text{cov}(Z'Z) = E((XV)'(XV)) = V'E(X'X)V = V'\Sigma_X V = \Lambda_p = \text{diag}(\lambda_1, \dots, \lambda_p), \quad (4.5)$$

missä siis $\text{cov}(z_i, z_j) = 0$, kun $i \neq j$, ja $\text{var}(z_j) = \lambda_j$, $i, j = 1, \dots, p$, eli pääkomponentit $z_j, j = 1, \dots, p$, ovat korreloimattomia. Pääkomponentit $z_j, j = 1, \dots, p$, selittävät rakennematriisin X kokonaisvaihtelun, sillä

$$\text{var}(Z) = \text{tr}(Z'Z) = \sum_{j=1}^p \lambda_j = r = \text{tr}(X'X) = \text{var}(X), \quad (4.6)$$

missä r on rakennematriisin X aste. Nämä kaikki edellä esitetyt pääkomponenttien $z_j, j = 1, \dots, p$, ominaisuudet ovat seurausta ortogonaalisesta muunnoksesta $XV = Z$.

Pääkomponenttiregressiossa selitettävä muuttuja y regressoidaan pääkomponentteja $z_j, j \in S_A$, vasten pienimmän neliösumman menetelmällä. Koska matriisi Z on ortogonaalinen, niin neliömatriisi $Z_A'Z_A$ on diagonaalinen ja sen käänteismatriisi on $(Z_A'Z_A)^{-1} = \text{diag}(1/\lambda_j), j \in S_A$, joka ei ole yhtä herkkä mahdollisille aineiston pyöristys- ja mittausvirheelle kuin käänteismatriisi $(X'X)^{-1}$. Koska pääkomponentit $z_j, j = 1, \dots, p$, ovat korreloimattomia, niin pääkomponentteja $z_j, j \in S_A$, vastaavat regressiokertoimet ovat korreloimattomia (Jackson, 1991).

Pääkomponenttiregressiossa valitaan kaikista rakennematriisin X pääkomponenteista $z_j, j \in S_p$ sopiva osajoukko $z_j, j \in S_A$, selittäviksi muuttujiksi pääkomponenttiregressiomalliin. Usein pääkomponenttiregressiomalliin valitaan selittäviksi muuttujiksi pääkomponentit $z_j, j \in S_A$, jotka selittävät suurimman osan rakennematriisin X vaihtelusta, jolloin pääkomponenttiregressiossa joukon S_A valinta perustuu pelkästään rakennematriisin X pääkomponenttianalyysiin. Tällöin yritetään selittää mahdollisimman paljon rakennematriisin X kokonaisvaihtelusta mahdollisimman pienellä selittävien pääkomponenttien lukumäärällä, jolloin indeksijoukkoon S_A eivät tule valituksi pääkomponentit, joiden $\text{var}(z_j) = \lambda_j \approx 0$. Indeksijoukko S_A voidaan valita myös siten, että valitaan pääkomponentit $z_j, j \in S_A$, jotka korreloivat voimakkaasti selitettävän muuttujan y kanssa. Täten pääkomponenttiregressiossa samassa tapauksessa indeksijoukon S_A alkioit saattavat vaihdella sen mukaan, mitä valintamenetelmää on käytetty joukkon S_A määräämiseen ja mikä on tämän valintamenetelmän tilastollinen peruste. (Brown, 1991.)

Kun joukko S_A on valittu, niin pääkomponenttiregressiossa rakennematriisin X approksimaatio on matriisi $Z_A V_A'$, missä $Z_A = (z_j)$, $V_A = (v_j)$ ja $j \in S_A$, jolloin neliömatriisin $X'X$ approksimaatio on matriisi (4.2). Pääkomponenttiregressiossa rakennematriisia X voidaan approksimoida myös singulaariarvohajotelmalla (2.10), missä nyt $T_r = (t_j)$, $V_r = (v_j)$, $\Lambda_r = \text{diag}(\lambda_j)$, $j \in S_A$, jolloin neliömatriisin $X'X$ approksimaatio on matriisi (2.10). Pääkomponenttiregressiossa

rakennematriisiin X approksimaatioiden $Z_A V_A'$ ja $T_A \Lambda_A V_A'$ aste on $A \leq p$.

4.1 Pääkomponenttiestimaattori a_{APK}

Oletetaan, että regressiomalli on määritelty yhtälön (4.3) mukaan. Oletetaan nyt, että lopulliseen pääkomponenttiregressiomalliin on valittu selittäviksi muuttujiksi pääkomponentit z_j , missä $j \in S_A$ ja $S_A \in S_p$. Kun käytetään merkintätapaa $Z_A = (z_j)$, $V_A = (v_j)$ ja $\Lambda_A = \text{diag}(\lambda_j)$, $j \in S_A$, niin pääkomponenttiregressiossa estimoidaan regressiomallia (4.3) mallilla $y = Z_A a_{APK} + e$, missä pääkomponenttiestimaattori $a_{APK} = (a_{jAPK})$, $j \in S_A$, on normaaliyhtälöiden

$$(Z_A' Z_A)^{-1} Z_A' y = Z_A' a_{APK} \quad (4.7)$$

ratkaisu, joka on

$$\begin{aligned} a_{APK} &= (Z_A' Z_A)^{-1} Z_A' y = (V_A' X' X V_A)^{-1} V_A' X' y = (V_A' V_p \Lambda_p V_p' V_A)^{-1} V_A' X' y \\ &= ((I_A \ 0) \Lambda_p (I_A \ 0)')^{-1} V_A' X' y = \Lambda_A^{-1} V_A' X' y = \Lambda_A^{-1} Z_A' y. \end{aligned} \quad (4.8)$$

Kun regressiomalli on esitetty kanonisessa muodossa (2.6), missä $q = 1$, niin pääkomponenttiestimaattori a_{APK} on

$$a_{jAPK} = u_j / \sqrt{\lambda_j}, \quad j \in S_A, \quad \text{ja} \quad a_{jAPK} = 0, \quad j \notin S_A. \quad (4.9)$$

Pääkomponenttiestimaattorin a_{APK} ja pienimmän neliösumman estimaattorin a_{PNS} välisestä yhteydestä saadaan pääkomponenttiestimaattorille a_{APK} vaihtoehtoinen esitys

$$a_{jAPK} = f_j a_{jPNS},$$

missä $f_j = 1$, kun $j \in S_A$, ja $f_j = 0$, kun $j \notin S_A$, eli

$$a_{jAPK} = a_{jPNS}, \quad j \in S_A, \quad \text{ja} \quad a_{jAPK} = 0, \quad j \notin S_A \quad (4.10)$$

(vrt. (3.10)). Pääkomponenttiestimaattori a_{APK} ja harjaestimaattori $a_H(k)$ ovat estimaattorin a_{PNS} eräänlaisia konveksikombinaatiota. Pääkomponenttiestimaattori a_{APK} on

$$a_{jAPK} = w_j 0 + (1 - w_j) a_{jPNS}, \quad j = 1, \dots, p, \quad (4.11)$$

missä $w_j = 0$, kun $j \in S_A$, ja $w_j = 1$, kun $j \notin S_A$ (vrt. (3.11)). Jos $S_A \neq S_p$, niin pääkomponentt-
iestimaattorin a_{APK} etäisyys origosta on pienempi kuin estimaattorin a_{PNS} etäisyys origosta ts.

$$\|a_{APK}\| < \|a_{PNS}\|. \quad (4.12)$$

Tämän vuoksi pääkomponenttiregressiota kutsutaan kutistavaksi menetelmäksi (engl. *shrinking method*) ja pääkomponenttiestimaattoria kutistavaksi estimaattoriksi (engl. *shrinking estimator*) (Brown, 1993).

4.2 Pääkomponenttiestimaattorin a_{APK} keskineliövirhe

Pääkomponenttiregressioestimaattorin a_{jAPK} harha on

$$E(a_{jAPK} - \alpha_j) = E(a_{jPNS}) - \alpha_j = 0, \quad j \in S_A \quad (4.13)$$

tai

$$E(a_{jAPK} - \alpha) = 0 - \alpha_j = -\alpha_j, \quad j \notin S_A.$$

Pääkomponenttiestimaattorin a_{APK} varianssi on

$$\text{var}(a_{jAPK}) = \text{var}(a_{jPNS}) = \sigma_y^2 / \lambda_j, \quad j \in S_A \quad (4.14)$$

tai

$$\text{var}(a_{jAPK}) = 0, \quad j \notin S_A.$$

Koska estimaattorin a_{APK} kokonaisvarianssi on

$$\text{var}(a_{APK}) = \sum_{j \in S_A} \text{var}(a_{jAPK}) = \sigma_y^2 \sum_{j \in S_A} 1 / \lambda_j, \quad (4.15)$$

niin mitä pienempi on $\text{var}(z_j) = \lambda_j, j \in S_A$, sitä pienempi on $\text{var}(\mathbf{a}_{APK})$ kuin $\text{var}(\mathbf{a}_{PNS})$.

Koska estimaattorin keskineliövirhe on sen harhan neliön ja kokonaisvariانسin summa, niin pääkomponenttiestimaattorin \mathbf{a}_{APK} keskineliövirhe on

$$\begin{aligned} \text{MSE}(\mathbf{a}_{APK}) &= \sum_{j \in S_A} \alpha_j^2 + \sigma_y^2 \sum_{j \in S_A} 1/\lambda_j \\ &- \sum_{j \in S_A} \alpha_j^2 + \text{MSE}(\mathbf{a}_{PNS}) - \sigma_y^2 \sum_{j \in S_A} 1/\lambda_j. \end{aligned} \quad (4.16)$$

Pääkomponenttiestimaattorilla \mathbf{a}_{APK} on pienempi keskineliövirhe kuin pienimmän neliösumman estimaattorilla \mathbf{a}_{PNS} eli $\text{MSE}(\mathbf{a}_{APK}) < \text{MSE}(\mathbf{a}_{PNS})$, kun

$$\sigma_y^2 \sum_{j \in S_A} 1/\lambda_j + \sum_{j \in S_A} \alpha_j^2 < \sigma_y^2 \sum_{j \in S_p} 1/\lambda_j \quad (4.17)$$

eli

$$\sum_{j \in S_A} \alpha_j^2 < \sigma_y^2 \sum_{j \in S_A} 1/\lambda_j$$

eli

$$\sum_{j \in S_A} \lambda_j \alpha_j^2 / \sigma_y^2 < 1.$$

Estimaattorin \mathbf{a}_{APK} keskineliövirheeseen vaikuttavat siis neliömatriisin $\mathbf{X}'\mathbf{X}$ ominaisarvot ja parametrit α ja σ_y^2 . (Leskinen, 1981.)

4.3 Ennusteet pääkomponenttiregressiossa

Ennustetaan estimoidulla pääkomponenttiregressiomallilla tuntematonta selitettävän muuttujan arvoa y_0 , joka voidaan kirjoittaa muotoon (2.23). Ennusteen $\hat{y}_{0APK} = \mathbf{x}_0' \mathbf{b}_{APK}$ varianssi on

$$\text{var}(\hat{y}_{0APK}) = \sigma_y^2 (1 + 1/n + \mathbf{x}_0' \mathbf{G}_{PK} \mathbf{x}_0) \quad (4.18)$$

$$- \sigma_y^2 (1 + 1/n + \mathbf{x}_0' (\sum_{j \in S_A} (1/\lambda_j) \mathbf{v}_j \mathbf{v}_j') \mathbf{x}_0),$$

missä G_{PK} on (4.2) (Brown, 1993). Oletetaan, että estimointiaineiston rakennematriisin X pääkomponenteilla $z_p, j \in S_A$, ennustetaan tuntematonta havaintovektoria y . Tällöin ennusteen $\hat{y}_{APK} = Z a_{APK}$ keskineliövirhe $MSE(\hat{y}_{APK})$ on

$$\begin{aligned} MSE(\hat{y}_{APK}) &= \sigma_y^2(n+1) + \sigma_y^2 \sum_{j \in S_p} \lambda_j (\text{var}(a_{jAPK}) + B^2(a_{jAPK})) \\ &= \sigma_y^2(n+1+A) + \sum_{j \in S_A} \lambda_j \alpha_j^2, \end{aligned} \quad (4.19)$$

missä $B^2(a_{jAPK})$ on estimaattorin a_{jAPK} harhan neliö.

4.4 Ykismuuttujaisessa tapauksessa valittujen pääkomponenttien indeksijoukon S_A valintamenetelmiä

Pääkomponenttiregressiossa on ongelmana valita kaikista pääkomponenteista $z_p, j \in S_p$ sopiva osajoukko $z_p, j \in S_A$, selittäväksi pääkomponenteiksi pääkomponenttiregressiomalliin. Seuraavassa on muutamia valittujen pääkomponenttien indeksijoukon S_A valintamenetelmiä.

1) Oletetaan, että $\epsilon \sim N(0, \sigma_y^2 I_n)$. Pääkomponenttiestimaattori a_{APK} voidaan nyt muodostaa seuraavasti

$$a_{jAPK} = \begin{cases} 0, & \text{jos } \lambda_j a_{jPNS}^2 / \hat{\sigma}_y < c \\ a_{jPNS}, & \text{jos } \lambda_j a_{jPNS}^2 / \hat{\sigma}_y \geq c \end{cases} \quad j = 1, \dots, p, \quad (4.20)$$

missä c on 4, 2 tai 1, jolloin se vastaa likimain merkitsevyytstasoja 5%, 16% ja 35%. Jos epäyhtälössä (4.17) parametreja α ja σ_y^2 approksimoidaan niiden pienimmän neliösumman estimaattoreilla, niin $c = 1$. (Leskinen, 1981.)

2) Rakennematriisin X pääkomponentit ovat seuraavassa suuruusjärjestyksessä $\text{var}(z_1) = \lambda_1 \geq \dots \geq \text{var}(z_p) = \lambda_p$, missä $\lambda_p, j = 1, \dots, p$, ovat neliömatriisin XX ominaisarvot. Valitaan indeksijoukkoon S_A ne A ensimmäistä rakennematriisin X pääkomponenttia, joiden varianssien summa on selittänyt

esimerkiksi 75% neliömatriisin XX kokonaisvarianssista. Jos neliömatriisi XX on korrelaatiomatriisi, niin tällöin tulee valituksi ne A ensimmäistä pääkomponenttia, jotka toteuttavat ensimmäisenä ehdon

$$\sum_{j=1}^A \text{var}(z_j)/r = \sum_{j=1}^A \lambda_j/r \geq 0.75 \quad (4.21)$$

(Draper & Smith, 1966). Tämän valintamenetelmän mukaan valittujen pääkomponenttien indeksijoukkoon S_A eivät kuulu ne pääkomponentit, joiden varianssi on pieni eli $\text{var}(z_j) = \lambda_j \approx 0$, sillä nämä pääkomponentit selittävät heikosti rakennematriisin X kokonaisvaihtelua. Tämä valintamenetelmä takaa, että pääkomponentit, jotka selittävät hyvin vähintään yhden selittävän muuttujan varianssin, tulevat mukaan joukkoon S_A .

3) Oletetaan, että rakennematriisi X on täysiasteinen eli $r = p$, mutta sen pienimmät ominaisarvot ovat lähes nollia. Jos oletetaan, että matriisin X osittaisaste (engl. fractional rank) s on välillä $A - 1 < s < A$, missä niin tällöin voidaan muodostaa osittainen pääkomponenttiestimaattori (engl. fractional principal component estimator)

$$a_{jPK} = \begin{cases} 0, & j = A + 1, \dots, p \\ a_A a_{jPNS}, & j = A \\ a_{jPNS}, & j = 1, \dots, A - 1, \end{cases} \quad (4.22)$$

missä $0 \leq a_A \leq 1$. Kun

$$\hat{a}_A = \hat{\tau}_A^2 / (1 + \hat{\tau}_A^2) \quad \text{ja} \quad \hat{\tau}_A^2 = \lambda_j a_{jPNS}^2 / \sigma_y^2, \quad (4.23)$$

niin saadaan seuraava osittainen pääkomponenttiestimaattori (engl. modified fractional rank estimator)

$$a_{jPK} = \begin{cases} 0, & \text{jos } \hat{\tau}_j^2 \leq 1, j = A + 1, \dots, p \\ \hat{a}_A a_{jPNS}, & j = A \\ a_{jPNS}, & j = 1, \dots, A - 1, \end{cases} \quad (4.24)$$

missä osittaisaste $s = A + \hat{\alpha}_A - 1$. (Leskinen, 1981.)

4) Estimoidaan indeksijoukko S_A *cross-validation* -menetelmällä. Muodostetaan selityksasteen mielessä paras yhden selittävän pääkomponentin regressiomalli ja tähän malliin lisätään uusia selittäviä pääkomponentteja, kunnes ennusteen virheneliösumma

$$PRESS(S_A) = \sum_{i=1}^n (y_i - \hat{y}_{(i)})^2, \quad (4.25)$$

missä y_i on selitettävän muuttujan havaittu arvo ja $\hat{y}_{(i)}$ on tämän ennuste, joka on saatu pääkomponenttiregressiolla ilman havaittua havaintoa i , on saavuttanut miniminsä. Indeksijoukkoon S_A lisätään uusia pääkomponentteja kuten etenevässä (engl. forward) valinnassa. (Stone & Brown, 1990.)

5) Lasketaan testisuure

$$W = \frac{(PRESS(S_A) - PRESS(S_{A+1}))/D_m}{PRESS(S_{A+1})/D_r}, \quad (4.26)$$

missä $PRESS(S_A)$ on ennusteen virheneliösumma (4.25), kun valittujen indeksien joukko on S_A , $PRESS(S_{A+1})$ on ennusteen virheneliösumma (4.25), kun valittujen indeksien joukkoon S_A on lisätty uusi selittävä pääkomponentti kuten etenevässä (engl. forward) valinnassa, $D_m = (n + p - 2A)$ ja

$$D_r = p(n - 1) - \sum_{j=1}^A (n + p - 2j). \quad (4.27)$$

Oletetaan aluksi, että $S_A = \{\}$, jolloin $\hat{y}_{(i)} = \bar{y}_{(i)}$, missä $\bar{y}_{(i)}$ selitettävän muuttujan y keskiarvo, joka on laskettu ilman havaintoa i . Jos testisuure $W > 1$, niin valittujen indeksien joukko S_A kasvaa yhdellä indeksillä, jolloin $S_A = S_{A+1}$, minkä jälkeen testataan uutta indeksijoukkoa S_{A+1} , kunnes testisuure $W \leq 1$. (Jackson, 1991.)

4.5 Monimuuttujainen pääkomponenttiestimaattori A_{PK}

Oletetaan, että regressiomalli on määritelty yhtälön (2.5) mukaan, missä selitettäviä muuttujia on useita. Kun pääkomponenttiregressiomalliin valitut pääkomponentit ovat z_j , $j \in S_A$, niin monimuuttujainen pääkomponenttiestimaattori A_{PK} on

$$A_{PK} = \Lambda_A^{-1} Z_A' Y = \Lambda_A^{-1} V_A' X' Y, \quad (4.28)$$

missä $S_A \in S_p$, $\Lambda_A^{-1} = \text{diag}(1/\lambda_j)$, $Z_A = XV_A$ ja $V_A = (v_j)$, v_j on neliömatriisin XX ominaisarvoa λ_j vastaava ominaisvektori ja S_A on valittujen pääkomponenttien indeksijoukko, $j \in S_A$ (vrt. kaava (4.8)).

Kun regressiomalli on määritelty kanonisen yhtälön (2.6) mukaan, niin monimuuttujainen pääkomponenttiestimaattori A_{PK} on

$$A_{PK} = (\sqrt{\Lambda_A^{-1}} \mathbf{0})U, \quad (4.29)$$

missä $\Lambda_A^{-1} = \text{diag}(1/\lambda_A)$, $j \in S_A$, S_A on valittujen pääkomponenttien indeksijoukko, nollamatriisin $\mathbf{0}$ ulottuvuus on $(n - A \times A)$ ja $U = TY$ (vrt. (4.9)).

Monimuuttujainen pääkomponenttiestimaattori A_{PK} voidaan kirjoittaa muotoon

$$a_{iPK}^{(j)} = u_i^{(j)} / \sqrt{\lambda_i}, \quad i \in S_A^{(j)}, \quad \text{ja} \quad a_{iPK}^{(j)} = 0, \quad i \notin S_A^{(j)}, \quad j = 1, \dots, q, \quad (4.30)$$

eli pääkomponenttiestimaattori A_{PK} voidaan hajottaa q yksimuuttujaiseen pääkomponenttiestimaattoriin $a_{APK}^{(j)}$, $j = 1, \dots, q$ (vrt. (4.9)). Jos pääkomponenttiregressiossa indeksijoukkojen $S_A^{(j)}$, $j = 1, \dots, q$, estimointi perustuu pelkästään selittävien muuttujien x_1, \dots, x_p , pääkomponenttianalyysiin eikä siinä huomioida pääkomponenttien z_p , $j = 1, \dots, p$ korrelaatiota selitettävien muuttujien $y^{(j)}$, $j = 1, \dots, q$, kanssa, niin tällöin indeksijoukot $S_A^{(j)}$, $j = 1, \dots, q$, ovat samoja ja monimuuttujainen pääkomponenttiestimaattori voidaan esittää muodossa (4.29). Jos indeksijoukkojen $S_A^{(j)}$, $j = 1, \dots, q$, yhtenä valintaperustana on pääkomponenttien z_p , $j = 1, \dots, p$, korrelaatiot selitettävien muuttujien $y^{(j)}$, $j = 1, \dots, q$, kanssa, niin tällöin indeksijoukot $S_A^{(j)}$, $j = 1, \dots, q$, saattavat erota toisistaan ja monimuuttujaista pääkomponenttiestimaattoria ei voida esittää muodossa (4.29), vaan se esitetään muodossa (4.30).

Monimuuttujainen pääkomponenttiestimaattori A_{PK} on harhainen estimaattori, jonka harha on yksimuuttujaisten estimaattoreiden $a_{APK}^{(j)}$, $j = 1, \dots, q$, harhojen summa ja jonka kokonaisvarianssi on yksimuuttujaisten estimaattoreiden $a_{APK}^{(j)}$, $j = 1, \dots, q$, varianssien summa. Täten monimuuttujaisen pääkomponenttiestimaattorin A_{PK} keskineliövirhe on yksimuuttujaisten estimaattoreiden $a_{APK}^{(j)}$, $j = 1, \dots, q$, keskineliövirheiden summa.

4.6 Monimuuttujaisessa tapauksessa valittujen pääkomponenttien indeksijoukon S_A valintamenetelmiä

Monimuuttujaiselle pääkomponenttiestimaattorille $A_{PK} = (\mathbf{a}_{APK}^{(1)}, \dots, \mathbf{a}_{APK}^{(q)})$ voidaan valita indeksijoukot $S_A^{(j)}, j = 1, \dots, q$, toisistaan riippumatta kappaleessa 4.4. esitettyin valintamenetelmin. Koska kappaleessa 4.4 indeksijoukkojen $S_A^{(j)}, j = 1, \dots, q$, valintamenetelmä (4.21) perustuu pelkästään selittävien muuttujien pääkomponenttianalyysiin, niin niissä indeksijoukot $S_A^{(j)}, j = 1, \dots, q$, tulevat olemaan samoja. Jos indeksijoukot $S_A^{(j)}, j = 1, \dots, q$, valitaan toisistaan riippumatta, niin tällöin ei hyödynnetä q normaaliyhtälöiden yhteisinformaatiota.

4.6.1 Cross-validation

Valitaan indeksijoukko $S_A = S_A^{(j)}, j = 1, \dots, q$, *cross-validation* -menetelmällä. Muodostetaan selityksasteen mielessä paras yhden selittävän pääkomponentin regressiomalli ja lasketaan ennusteen virheneliösumma

$$\text{PRESS}(S_A) = \sum_{i=1}^n \sum_{j=1}^q (y_i^{(j)} - \hat{y}_{(i)}^{(j)})^2, \quad (4.31)$$

missä $y_i^{(j)}$ on selitettävän muuttujan havaittu arvo ja $\hat{y}_{(i)}^{(j)}$ on tämän ennuste, joka on saatu pääkomponenttiregressiolla ilman havaintoa i . Valitaan indeksijoukkoon S_A uusia pääkomponentteja kuten etenevässä menetelmässä (vrt. (4.25)). Kun joukkoon S_A on lisätty uusi pääkomponentti, niin lasketaan tämän jälkeen aina ennusteen virheneliösumma. Valitaan se valittujen pääkomponenttien indeksijoukko S_A , joka antaa pienimmän ennusteen virheneliösumman.

5 Osittainen pienimmän neliösumman regressio

Osittaista pienimmän neliösumman regressiota (engl. partial least squares regression) eli PLS-regressiota käytetään, kun joko selittävät tai selitettävät muuttujat tai kumpikin ovat multikollineaarisia. PLS-regressiota voidaan käyttää jopa silloin, kun havaintoja on vähemmän kuin selittäviä muuttujia tai kun selitettäviä muuttujia on enemmän kuin selittäviä muuttujia. PLS-algoritmi on laskuteknisesti helpohko menetelmä, jolla pystytään melko vaivattomasti mallittamaan suuriakin aineistoja. Selittävien muuttujien multikollineaarisuuden ongelma ratkaistaan pääkomponenttiregressiossa ja PLS-regressiossa samalla periaatteella, mutta lisäksi PLS-regressiossa mallitetaan havaintomatriisia Y . Koska PLS-regressiossa regressoidaan havaintomatriisi Y rakennematriisia X vasten ja päin vastoin, niin siinä voidaan hylätä selittävän ja selitettävän muuttujan roolit. PLS-regressiolla on saatu tarkkoja selitettävien muuttujien ennusteita ja sitä on käytetty usein kalibrointimenetelmänä esim. Garthwaite (1994), Otto & Wegscheider (1985) ja Haaland & Thomas (1988). PLS-regressio on tässä tutkielmassa esitetyistä säännöllistämismenetelmistä nuorin. Sen teoriaa ovat kehittäneet mm. Helland (1988) ja Höskuldsson (1988).

PLS-algoritmissa matriisi X , jonka ulottuvuus on $n \times p$, kirjoitetaan rekursiiviseen muotoon

$$X = TP' + X_{A+1} = t_1 p'_1 + t_2 p'_2 + \dots + t_A p'_A + X_{A+1}, \quad (5.1)$$

missä matriisin T ulottuvuus on $n \times A$, matriisin P ulottuvuus on $p \times A$ ja X_{A+1} on $(n \times p)$ -jäännösmatriisi. PLS-regressiossa matriisi Y , jonka ulottuvuus on $n \times q$, kirjoitetaan myös rekursiiviseen muotoon

$$Y = TC' + Y_{A+1} = t_1 c'_1 + t_2 c'_2 + \dots + t_A c'_A + Y_{A+1}, \quad (5.2)$$

missä matriisin T ulottuvuus on $n \times A$, matriisin C ulottuvuus on $q \times A$ ja jäännösmatriisi Y_{A+1} ulottuvuus on $n \times q$. PLS-regressiossa matriisin X ja matriisin Y yhteyden ajatellaan välittyvän PLS-komponenttien kautta. PLS-regressiossa muodostetaan PLS-komponentit siten, että niitä on mahdollisimman vähän, mutta ne selittävät suurimman osan rakennematriisin X vaihtelusta, jolla voidaan parhaiten ennustaa havaintomatriisia Y .

Kun matriisi X kirjoitetaan rekursiiviseen muotoon (5.1), niin joko matriisin T sarakevekt-

reiden tai matriisin C sarakevektoreiden tulee olla keskenään ortogonaalisia. Tässä tutkielmassa esitetyssä PLS-algoritmissa, joka on Höskuldssonin (1988) mukaan, on vaadittu matriisin T sarakevektoreiden keskinäinen ortogonaalisuus. Riippumatta siitä, onko PLS-algoritmissa vaadittu matriisin T tai matriisin C sarakevektoreiden keskinäinen ortogonaalisuus, se antaa samat sovitteet (Helland, 1988).

Kun pääkomponenttianalyysiä voidaan pitää menetelmänä, jossa mallitetaan havaintoaineiston selitettävien muuttujien kokonaisvarianssin rakennetta, niin PLS-menetelmää voidaan pitää menetelmänä, jossa mallitetaan havaintoaineiston selitettävien ja selittävien muuttujien kovarianssirakennetta. Tällöin yksi havainto ei ole oleellinen osa mallin teoriaa vaan yksittäisellä havainnolla estimoidaan ainoastaan mallin parametreja.

Tässä tutkielmassa on esitetty vain monimuuttujainen PLS-regressio. Yksimuuttujaisen PLS-algoritmin on esittänyt mm. Helland (1988). Tässä tutkielmassa osittaisen pienimmän neliösumman regression teoria ja algoritmi on esitetty Höskuldssonin (1988) mukaan.

5.1 Osittaisen pienimmän neliösumman algoritmi ja sen geometrinen tulkinta

Asetetaan vektorin u_0 arvoksi matriisin Y ensimmäinen sarake. Asetetaan $X_1 = X$ ja $Y_1 = Y$. PLS-algoritmin i :nnes kierros on

- 1) $w_i = X_i' u_{i-1} / (u_{i-1}' u_{i-1})$
- 2) Skaalataan vektorin w_i pituudeksi yksi.
- 3) $t_i = X_i w_i$
- 4) $c_i = Y_i' t_i / (t_i' t_i)$
- 5) Skaalataan vektorin c_i pituudeksi yksi
- 6) $u_i = Y_i c_i / (c_i' c_i)$
- 7) Jos u_i ei konvergoi askeleessa 6, niin siirrytään askeleeseen 2, muuten siirrytään askeleeseen 8.
- 8) $p_i = X_i' t_i / (t_i' t_i)$
- 9) $q_i = Y_i' u_i / (u_i' u_i)$
- 10) $b_i = u_i' t_i / (t_i' t_i)$
- 11) $X_{i+1} = X_i - t p_i'$
- 12) $Y_{i+1} = Y_i - b t_i c_i'$

4/

13) Jos matriisi $X \neq \mathbf{0}$ tai käytetyn pysähdyssäännön kriteerin mukaan jatketaan iterointia, niin siirrytään askeleeseen 1, muuten lopetetaan iterointi.

PLS-algoritmissa on useita pienimmän neliösumman estimaattoreita. Alla ovat regressiomallit, joissa numero regressiomallin edessä viittaa PLS-algoritmin askeleeseen, jossa on esiteyn regressiomallin pienimmän neliösumman estimaattori:

$$1) X_i = u_{i-1} w_i' + \epsilon$$

$$4) Y_i = t_i c_i' + \epsilon$$

$$8) X_i = t_i p_i' + \epsilon$$

$$9) Y_i = u_i q_i' + \epsilon$$

$$10) u_i = t_i b_i' + \epsilon$$

PLS-regressiomallin ulottuvuuden i vektorit voidaan kirjoittaa edellisen ulottuvuuden $i-1$ vektoreiden funktiona. Vektori u_i on vektorin u_{i-1} funktiona

$$\begin{aligned} u_i &= Y c_i / (c_i' c_i) \\ &= Y Y' t_i / (c_i' c_i) (t_i' t_i) \\ &= Y Y' X w_i / ((c_i' c_i) (t_i' t_i) (w_i' w_i)) \\ &= Y Y' X X' u_{i-1} / ((c_i' c_i) (t_i' t_i) (w_i' w_i) (u_{i-1}' u_{i-1})) . \end{aligned}$$

Vektorit c_i , t_i ja w_i ovat

$$\begin{aligned} c_i &= Y' X X Y c_{i-1} / ((t_i' t_i) (w_i' w_i) (u_{i-1}' u_{i-1}) (c_{i-1}' c_{i-1})) \\ t_i &= X X' Y Y' t_{i-1} / ((w_i' w_i) (u_{i-1}' u_{i-1}) (c_{i-1}' c_{i-1}) (t_{i-1}' t_{i-1})) \\ w_i &= X' Y Y' X w_{i-1} / ((u_{i-1}' u_{i-1}) (c_{i-1}' c_{i-1}) (t_{i-1}' t_{i-1}) (w_{i-1}' w_{i-1})) . \end{aligned}$$

$$XX'YY't = at \quad (5.3)$$

Jos PLS-algoritmin vektorit u , t , w ja c konvergoivat, niin

$$YY'XX'u = au \quad (5.4)$$

$$Y'XX'Yc = ac \quad (5.5)$$

$$X'YY'Xw = aw, \quad (5.6)$$

missä a on kyseisen ominaisarvottehtävän suurin ominaisarvo. Täten vektorit u , c , t ja w ovat approksimoitujen matriisien suurinta ominaisarvoa vastaavat ominaisvektorit.

Tässä tutkielmassa matriisien U ja T sarakevektoreita kutsutaan PLS-komponenteiksi ja vektorit u , ja t , muodostavat ns. PLS-komponenttiparin.

Tarkastellaan PLS-algoritmissa olevien vektoreiden ja matriisien geometrisia ominaisuuksia. Merkitään

$$W = (w_1, w_2, \dots, w_A), \quad (5.7)$$

missä W on $(n \times A)$ -matriisi, jonka sarakkeina ovat vektorit w . Käytetään tätä merkintätapaa kaikille PLS-algoritmin tuottamille matriiseille.

PLS-algoritmin ideana on, että jäännösmatriisi X_i lasketaan aina edellisestä jäännösmatriisista X_{i-1} :

$$\begin{aligned} X_i &= X_{i-1} - t_{i-1}p'_{i-1} \\ &= X_{i-1} - t_{i-1}t'_{i-1}X_{i-1}/(t'_{i-1}t_{i-1}) \\ &= (I - t_{i-1}t'_{i-1}/(t'_{i-1}t_{i-1}))X_{i-1} \end{aligned} \quad (5.8)$$

$$= (\mathbf{I} - \mathbf{t}_{i-1}\mathbf{t}'_{i-1}/(\mathbf{t}'_{i-1}\mathbf{t}_{i-1}))(\mathbf{X}_{i-2} - \mathbf{t}_{i-2}\mathbf{t}'_{i-2}\mathbf{X}_{i-2}/(\mathbf{t}'_{i-2}\mathbf{t}_{i-2})) .$$

PLS-algoritmissa olevilla vektoreilla \mathbf{w} , \mathbf{t} , \mathbf{u} ja \mathbf{p} ovat seuraavat ortogonalisuusominaisuudet:

Ominaisuus 1

Painovektorit \mathbf{w} ovat keskenään ortogonalisia:

$$\mathbf{w}'_i\mathbf{w}_j = 0 , \quad \text{kun } i \neq j . \quad (5.9)$$

Todistus : Oletetaan, että $i < j$. Kirjoitetaan jäännösmatriisi \mathbf{X}_j muotoon

$$\mathbf{X}_j = \mathbf{Z}(\mathbf{X}_i - \mathbf{t}_i\mathbf{t}'_i/(\mathbf{t}'_i\mathbf{t}_i)\mathbf{X}_i) , \quad (5.10)$$

missä \mathbf{Z} on jokin matriisi. Osoitetaan, että

$$\mathbf{X}_j\mathbf{w}_i = 0 , \quad (5.11)$$

kun $j > i$. Kun määritellään jäännösmatriisi \mathbf{X}_j yhtälön (5.11) mukaan, niin saadaan

$$\mathbf{X}_j\mathbf{w}_i = (\mathbf{X}_i - \mathbf{t}_i\mathbf{t}'_i/(\mathbf{t}'_i\mathbf{t}_i)\mathbf{X}_i)\mathbf{w}_i - \mathbf{t}_i - \mathbf{t}_i\mathbf{t}'_i/(\mathbf{t}'_i\mathbf{t}_i)\mathbf{t}_i - \mathbf{t}_i - \mathbf{t}_i = 0 . \quad (5.12)$$

Kun \mathbf{w}_j on määritelty yhtälön (5.7) mukaan ja käytetään apuna yhtälöä (5.12), niin saadaan

$$\mathbf{w}'_j\mathbf{w}_i = \mathbf{w}'_j\mathbf{X}'_j\mathbf{Y}_j\mathbf{Y}'_j\mathbf{X}_j\mathbf{w}_i/a_j = 0 . \quad \square \quad (5.13)$$

Ominaisuus 2

Vektorit \mathbf{t}_i ovat keskenään ortogonalisia:

$$\mathbf{t}'_i\mathbf{t}_j = 0 , \quad \text{kun } i \neq j . \quad (5.14)$$

Todistus: Oletetaan $i < j$. Jäännösmatriisi X_j on jäännösmatriisin X_i funktiona

$$\begin{aligned}
 X_j &= X_{j-1} - X_{j-1} w_{j-1} t'_{j-1} X_{j-1} / (t'_{j-1} t_{j-1}) \\
 &= X_{j-1} (I - w_{j-1} t'_{j-1} X_{j-1} / (t'_{j-1} t_{j-1})) \\
 &= X_{i+1} Z \\
 &= (X_i - t_i t'_i X_i / (t'_i t_i)) Z,
 \end{aligned} \tag{5.15}$$

missä Z on jokin matriisi. Nyt saadaan

$$t'_i X_j - (t'_i X_i - t'_i t_i t'_i X_i / (t'_i t_i)) = (t'_i X_i - t'_i X_i) = 0 \tag{5.16}$$

kun $i < j$ ja tällöin

$$t'_i t_j - t'_i X_j w_j = 0, \quad \text{kun } i < j. \quad \square \tag{5.17}$$

Täten avaruuden $C(X)$, joka on matriisin X sarakkeiden virittämä avaruus, ortogonaalisena kantana ovat vektorit t . Kun oletetaan, että matriisin X aste on r , niin matriisi X voidaan kirjoittaa rekursiiviseen muotoon

$$X = \sum_{i=1}^r t_i p'_i + X_{r+1}, \tag{5.18}$$

missä jäännösmatriisi X_{r+1} on ortogonaliinen matriisin Y kanssa. Jos oletetaan, että matriisi X_{r+1} on nollamatriisi, niin matriisi X_0 ei sisällä informaatiota matriisista Y .

Ominaisuus 3

Vektorit w_i ja p_j ovat keskenään ortogonaalisia:

$$w'_i p_j = 0, \quad \text{kun } i < j. \tag{5.19}$$

Todistus: Yhtälön (5.11) mukaan saadaan

$$\mathbf{w}'_i \mathbf{p}_j = \mathbf{w}'_i \mathbf{X}'_j \mathbf{t}_j (\mathbf{t}'_j \mathbf{t}_j)^{-1} = 0, \quad \text{kun } i < j. \quad \square \quad (5.20)$$

Ominaisuus 4

Vektorit \mathbf{p}_i ovat ortogonalisia avaruudessa $\ker(\mathbf{X}) = \{\mathbf{p}_i \mid \mathbf{X}\mathbf{p}_i = 0\}$:

$$(\mathbf{p}_i, \mathbf{p}_j)_X = \mathbf{p}'_i (\mathbf{X}'\mathbf{X})^{-1} \mathbf{p}_j = 0 \quad (5.21)$$

kaikilla $i \neq j$, missä $(\mathbf{X}'\mathbf{X})^{-1}$ on matriisin $\mathbf{X}'\mathbf{X}$ yleistetty käänteismatriisi.

Todistus: Yhtälöstä (5.18) saadaan neliömatriisille $\mathbf{X}'\mathbf{X}$ esitys

$$\mathbf{X}'\mathbf{X} = \sum (\mathbf{t}'_i \mathbf{t}_i) \mathbf{p}_i \mathbf{p}'_i = \mathbf{P}\mathbf{L}\mathbf{P}' \quad (5.22)$$

missä \mathbf{L} diagonaalimatriisi, jossa on kertoimet $\mathbf{t}'_i \mathbf{t}_i$. \square

Edellä on esitetty vektoreiden \mathbf{w} , \mathbf{t} , \mathbf{u} ja \mathbf{p} ortogonaaliset ominaisuudet. Yleisesti kuitenkin on

$$\mathbf{p}'_i \mathbf{w}_j \neq 0 \quad (5.23)$$

kaikilla $i < j$. Neljästä edellä mainitusta PLS-algoritmin ominaisuudesta seuraa, että jäännösmatriisi \mathbf{X}_i ei riipu tavasta, jolla uusi vektori \mathbf{t}_i saadaan, vaan vektorit \mathbf{t} määritellään siten, että ne muodostavat avaruuden $C(\mathbf{X})$ aliavaruuden $C(\mathbf{T})$ ortogonaalisen kannan.

Kun vektorin \mathbf{c}_i pituus on skaalattu ykköseksi, niin jäännösmatriisit voidaan kirjoittaa muotoon

$$\mathbf{X}_{i+1} = \mathbf{X}_i - \mathbf{t}_i \mathbf{p}'_i \quad \text{ja} \quad \mathbf{Y}_{i+1} = \mathbf{Y}_i - b_i \mathbf{t}_i \mathbf{c}'_i \quad (5.24)$$

$$\mathbf{Y}'_{i+1} \mathbf{Y}_{i+1} = \mathbf{Y}'_i \mathbf{Y}_i - b_i^2 (\mathbf{t}'_i \mathbf{t}_i) \mathbf{c}_i \mathbf{c}'_i. \quad (5.25)$$

$$Y'_{i+1}X_{i+1} = Y'_iX_i - b_i(t'_i t_i)c_i p'_i. \quad (5.26)$$

$$X'_{i+1}X_{i+1} = X'_iX_i - (t'_i t_i)p_i p'_i \quad (5.27)$$

Oletetaan nyt, että vektorin c pituutta ei ole skaalattu ykköseksi PLS-algoritmissa. Tällöin

$$u't = c'Y't/(c'c) = c'(Y't)/(c'c) = c'c(t't)/(c'c) = t't. \quad (5.28)$$

Yllä olevasta yhtälöstä seuraa, että

$$b_i = u'_i t_i / (t'_i t_i) = 1 \quad (5.29)$$

kaikilla i . Tästä lähtien oletetaan, että vektorin c pituus ei ole skaalattu ykköseksi. Jos kuitenkin halutaan saada kerroin b_i , niin se saadaan yhtälöstä

$$b_i = |c_i|. \quad (5.30)$$

5.2 Muita osittaisen pienimmän neliösumman tulkintoja

Rotatoidaan matriisi X eli tehdään ortogonaalinen muunnos

$$S = X O_X, \quad (5.31)$$

missä O_X on ortogonaalinen matriisi ja S on rotatoitu matriisi X . Koska

$$\sum s_{ij}^2 = \text{tr}(S'S) = \text{tr}(X'X) = \sum x_{ij}^2, \quad (5.32)$$

niin matriisin X rotatointi ei muuta matriisin X kokonaisvaihtelua (vrt. (4.6)). Rotatoidaan myös matriisi Y

$$Z = YO_Y, \quad (5.33)$$

missä Z on rotatoitu matriisi Y ja O_Y ortogonalinen matriisi. Merkitään matriisin S sarakkeita merkinnällä s_i ja matriisin Z sarakkeita merkinnällä z_i ja oletetaan lisäksi, että selittäviä muuttujia on enemmän kuin selitettäviä muuttujia eli $p > q$. Mitataan vektorin s etäisyyttä vektorista z niiden neliöidyllä etäisyydellä I . Vektorin s neliöity etäisyys I vektorista z on

$$I = \sum_{i=1}^q |s_i - z_i|^2 + \sum_{i=q+1}^p |s_i|^2. \quad (5.34)$$

Tulkinta 1

Oletetaan, että ortogonaaliset matriisit O_X ja O_Y ovat määritelty siten, että neliöity etäisyys I minimoituu, mikä intuitiivisesti merkitsee, että rotatoidaan matriisin X sarakkeiden virittämä avaruus $C(X)$ ja matriisin Y sarakkeiden virittämä avaruus $C(Y)$ siten, että komponentit, jotka kuuluvat avaruuteen $C(Y)$, ovat mahdollisimman lähellä komponenttejaan, jotka kuuluvat avaruuteen $C(X)$.

Oletetaan, että kaksi ortogonaalista matriisia on määritelty siten, että neliöity etäisyys I minimoituu, jolloin PLS-algoritmin vektorit t ja u eli PLS-komponenttipari saa arvon

$$t = s_1 \quad \text{ja} \quad u = z_1. \quad (5.35)$$

Todistus: Yhtälö (5.35) voidaan kirjoittaa muotoon

$$I = \text{tr}(X'X) + \text{tr}(Y'Y) - 2\text{tr}(X'YO_YO_X'). \quad (5.36)$$

Kuten myöhemmin tullaan osoittamaan, niin ortogonaaliset matriisit O_Y ja O_X , jotka ovat matriisin XY singulaariarvohajotelmasta, minimoivat neliöidyn etäisyyden I

$$X'Y = \sum e_i g_i f_i' = GEF', \quad (5.37)$$

missä matriisin G ensimmäisenä sarakeena on vektori t ja matriisin F ensimmäisenä sarakeena on vektori u . □

Neliöidyn etäisyyden I minimi on

(5.38)

$$\min I = \text{tr}(X'X) + \text{tr}(Y'Y) - 2\text{tr}((X'Y Y'X)^{1/2})$$

$$= \text{tr}(X'X) + \text{tr}(Y'Y) - 2 \sum e_i .$$

PLS-algoritmissa valitaan yksi PLS-komponenttipari kerallaan. Etäisyys I vähenee jokaisella iterointikierröksellä i $2e_{(i)}$ verran, missä $e_{(i)}$ on suurin matriisin $X_i'Y_i$ singulaariarvo. Koska iterointikierröksellä i matriisin $X_i'Y_i$ toiseksi suurin singulaariarvo on pienempi kuin iterointikierröksellä $i + 1$ matriisin $X_{i+1}'Y_{i+1}$ suurin singulaariarvo, niin vain yksi PLS-komponenttipari tulee valituksi jokaisella PLS-algoritmin iterointikierröksellä. Täten etäisyyttä I voidaan pienentää eniten, kun valitaan yksi PLS-komponenttipari kerrallaan, lasketaan matriisin X ja Y jäännösmatriisit ja lasketaan seuraavan ulottuvuuden PLS-komponenttipari saaduista jäännösmatriiseista. PLS-algoritmia voidaan pitää siis askeltavana menetelmänä (engl. stepwise procedure), jossa valitaan yksitellen PLS-komponenttipari, joiden välinen etäisyys on pienin. PLS-algoritmia jatketaan niin kauan kuin avaruudessa $C(X)$ tai avaruudessa $C(Y)$ ei ole PLS-komponenttipareja, jotka ovat riittävän lähellä toisiaan, ts. tilastollisesti merkitseviä PLS-komponenttipareja ei ole jäljellä. Kun avaruudessa $C(X_i)$ ei ole jäljellä tilastollisesti merkitseviä PLS-komponentteja, niin avaruudessa $C(X_i)$ ei ole mitään, millä voitaisiin ennustaa havaintomatriisia Y .

Olkoon komponentit $f \in C(X)$ ja $g \in C(Y)$:

$$f = Xd \quad \text{ja} \quad g = Ye , \quad (5.39)$$

missä $|d| = 1$ ja $|e| = 1$. Komponenttien f ja g otoskovarianssi on

$$\text{cov}(f, g) = f'g/n . \quad (5.40)$$

PLS-algoritmi valitsee komponentit f ja g siten, että niiden välinen otoskovarianssi on suurin mahdollinen otoskovarianssi komponenteilla $f \in C(X)$ ja $g \in C(Y)$.

Tulkinta 2

Vektoreilla w ja c saadaan komponentit $f \in C(X)$ ja $g \in C(Y)$ siten, että

$$(\text{cov}(t, u))^2 = (\text{cov}(Xw, Yc))^2 = \max(\text{cov}(f, q))^2 , \quad (5.41)$$

kun $|d| = |e| = 1$.

Todistus: Oletetaan, että matriisin $X'Y$ singulaariarvohajotelma on

$$X'Y = \sum a_i f_i g_i' . \quad (5.42)$$

Matriisin $X'Y$ suurin singulaariarvon a_1 neliö toteuttaa ehdon

$$(a_1)^2 = \max (d'X'Ye)^2 , \quad (5.43)$$

kun $|d| = |e| = 1$, $d = f_1$ ja $e = g_1$. Vektorit w ja c saavat PLS-algoritmissa arvot

$$w = f_1 \quad \text{ja} \quad c = g_1 , \quad (5.44)$$

jolloin väite (5.41) on tosi. \square

Tulkinta 3

Monissa käytännön sovelluksissa korvataan kovarianssimatriisi $X'X$ matriisilla $X'VX$, missä matriisi V on positiivisesti definiitimatriisi. PLS-algoritmissa matriisi V määritellään

$$V = YY' , \quad (5.45)$$

jolloin

$$X'VX = X'YY'X . \quad (5.46)$$

Tällöin PLS-algoritmissa havainnot, joita vastaavat selitettävien muuttujien arvot ovat lähellä nollaa, saa alhaisen painon ja havainnot, joita vastaavat selitettävien muuttujien arvot ovat itseisarvoltaan suuria, saa vastaavasti suuremman painon. Koska matriisi Y on keskitetty, niin matriisin Y keskiarvoa lähellä olevat arvot ovat lähes nollia. PLS-algoritmissa lähellä keskiarvoon olevat havaintomatriisin Y arvot katsotaan olevan kohinaa, jota ei PLS-regressiolla hyödytä mallitaa. PLS-algoritmissa painotetaan havaintoja, joita vastaavilla $y^{(j)}$, $j = 1, \dots, q$, on suurimmat arvot, koska niissä katsotaan olevan eniten informaatiota käytetystä datasta. Tämä kuitenkin edellyttää, että matriisissa Y arvot, jotka eivät ole lähellä nollaa, ovat virheettömästi mitattu.

PLS-algoritmissa lasketaan ensiksi ominaisarvot ja -vektorit painotetusta kovarianssimatriisista

$$X'YY'X = ODO', \quad (5.47)$$

jonka jälkeen rotatoidaan matriisi X eli tehdään muunnos

$$F = XO. \quad (5.48)$$

PLS-algoritmissa regressoidaan Y aina yhtä muunnoksen (5.48) komponenttia vasten iterointi-
kierroksittain. Huomattavaa on, että pääkomponenttiregressiossa ominaisarvot ja -vektorit
lasketaan kovarianssimatriisista

$$X'X = ODO', \quad (5.49)$$

jonka jälkeen rotatoidaan matriisi X eli tehdään muunnos

$$F = XO$$

ja regressiodaan havaintomatriisi Y valittuja $f_j, j \in S_A$, vasten.

Oletetaan, että iteroidaan PLS-algoritmin iterointikierrosta i . Tällöin w_i on suurinta
ominaisarvoa vastaava ominaisarvovektori yhtälössä (5.48), missä X on korvattu jäännösmatriisil-
la X_i ja Y on korvattu jäännösmatriisilla Y_i . Vektori $t_i = X_i w_i$ eli vektori t_i vastaa ortogonalista
muunnosta (5.49). Regressoidaan matriisi Y_i muunnosta t_i vasten. Jäännösmatriisin Y_i projektio
on

$$t_i c_i' \quad (5.50)$$

ja täten jäännösmatriisi Y_{i+1} voidaan kirjoittaa muotoon

$$Y_{i+1} = Y_i - t_i c_i'. \quad (5.51)$$

Jäännösmatriisin Y_{i+1} kovarianssimatriisi on

$$Y_{i+1}' Y_{i+1} = Y_i' Y_i - c_i c_i' (t_i' t_i). \quad (5.52)$$

Projektio (5.50) ei muutu, jos jäännösmatriisi Y_i korvataan matriisilla Y .

Todistus: Olkoon w_i matriisin $X'YY'X$ suurinta ominaisarvoa vastaava ominaisarvovektori ja $t_i =$

$x_i w_i$. Matriisin Y_i projektiio vektorille t_i on

$$t_i' Y_i / (t_i' t_i) t_i = t_i c_i' \quad (5.53)$$

ja jäännöskovarianssimatriisi $Y_{i+1}' Y_{i+1}$ on

$$\begin{aligned} Y_{i+1}' Y_{i+1} &= (Y_i - t_i c_i') (Y_i - t_i c_i') \\ &= Y_i' Y_i - c_i t_i' Y_i - Y_i' t_i c_i' + c_i c_i' (t_i' t_i) \\ &= Y_i' Y_i - c_i c_i' (t_i' t_i) . \end{aligned} \quad (5.54)$$

Kun $Y_i' t_i = Y' t_i$, niin jäännösmatriisin Y_i projektiio vektorille t_i on sama kuin matriisin Y projektiio vektorille t_i . \square

PLS-algoritmia voidaan täten pitää ortogonaalisten komponenttien askeltavana regressiomenetelmänä, jossa jokaisessa askeleessa määrätään PLS-komponenttipari, jolla saadaan painotetun kovarianssimatriisin $X'VX = X'Y'YX$ suurin 'varianssi'. Tämä kovarianssimatriisi on painotettu siten, että x_{ip} , jota vastaava $y_i^{(k)}$ on lähellä kohinaa eli $y_i^{(k)} \approx 0$, saa pienen painon ja x_{ip} , jota vastaava $y_i^{(k)}$ ei ole lähellä kohinaa, saa suuremman painon $i = 1, \dots, n, j = 1, \dots, p, k = 1, \dots, q$.

Matriisin Y regressio matriisille $T = (t_1, t_2, \dots, t_i)$ voidaan kirjoittaa muotoon

$$Y = \sum_{j=1}^i t_j c_j' + Y_{i+1} , \quad (5.55)$$

Täten jäännösmatriisin Y_{i+1} kovarianssimatriisi voidaan kirjoittaa muotoon

$$Y_{i+1}' Y_{i+1} = Y' Y - \sum_{j=1}^i c_j c_j' (t_j' t_j) . \quad (5.56)$$

5.3 PLS-algoritmin matemaattinen perusta

PLS-menetelmä perustuu matriisin $X'Y$ singulaariarvohajotelmaan

$$X'Y = \sum a_i f_i g_i' , \quad (5.57)$$

missä f_i ja g_i ovat approksimoidun ulottuvuuden i ortonormaalit vektorit ja a_i on matriisin $X'Y$ i :nneksi suurin singulaariarvo. Matriisin $X'Y$ suurimmalle singulaariarvolle a_1 pätee

$$a_1^2 = \max (f'X'Yg)^2 , \quad (5.58)$$

missä $|f| = |g| = 1$. Tämä on PLS-algoritmin tulkinta 2. Yhtälön (5.58) todistus perustuu Cauchy-Schwartz -epäyhtälöön.

Tulkintaa 1 todistettaessa oletettiin, että matriisilla Y on vähemmän sarakkeita kuin matriisilla X . Matriisit Y ja O_Y saadaan samankokoiseksi kuin matriisi X , kun korvataan matriisien Y ja O_Y puuttuvat sarakkeet nollavektoreilla. Nyt neliöity etäisyys I voidaan kirjoittaa muotoon

$$I = |X O_X - Y O_Y|^2 = \text{tr}(X'X) + \text{tr}(Y'Y) - 2 \text{tr}(X'Y O_Y O_X') . \quad (5.59)$$

Neliöidyn etäisyyden I minimi on sama kuin

$$\max \text{tr}(X'YV) , \quad (5.60)$$

missä $X'YV$ maksimoidaan ortogonaalisten matriisien V suhteen. Tällöin matriisin V ratkaisu on

$$V = GF' , \quad (5.61)$$

missä G ja F ovat matriisin $X'Y$ singulaariarvohajotelman matriiseja. Täten yksi neliöidyn etäisyyden I mahdollinen minimi on neliöity etäisyys I (5.59), missä

$$O_X = F \quad \text{ja} \quad O_Y = G , \quad (5.62)$$

minkä todistaa tulkinta 1, eli

$$I_{\min} = \text{tr}(X'X) + \text{tr}(Y'Y) - 2 \text{tr}(X'YF'X)^{1/2} \quad (5.63)$$

$$= \text{tr}(X'X) + \text{tr}(Y'Y) - 2 \sum a_i .$$

Pääperiaate, jolla PLS-algoritmi toimii, saadaan epäyhtälöstä

$$s_i^2(A - B) \geq s_{i+k}^2(A) - a_{i+k}^2, \quad (5.64)$$

missä $s_i(A)$ on matriisin A i :nneksi suurin singulaariarvo ja k on matriisin B aste. Kun tätä epäyhtälöä sovelletaan PLS-algoritmiin, niin saadaan

$$s_1^2(X'_{i+1}Y) - s_1^2(X'_iY - p_i t'_i Y) \geq s_2^2(X'_iY). \quad (5.65)$$

Täten matriisin $X_i Y$ suurin singulaariarvo iterointikierröksellä $i + 1$ on suurempi kuin toiseksi suurin singulaariarvo iterointikierröksellä i . Tämän vuoksi vain yksi PLS-komponenttipari kerrallaan tulee valituksi PLS-algoritmissa.

Matriisin $X_i Y$ suurimmalla singulaariarvolla s_1 on myös tulkinta

$$(s_1(X'_i Y))^2 = \max(f' X'_i Y Y' X_i f), \quad (5.66)$$

kun $|f| = 1$. Tarkastellaan seuraavaa epäyhtälöä

$$(s_1(X'_{i+1} Y))^2 - \max(f' X'_{i+1} Y Y' X_{i+1} f) \quad (5.67)$$

$$= \max(f' (I - p_i w'_i) X'_i Y Y' X_i (I - w_i p'_i) f)$$

$$\leq \max(f' X'_i Y Y' X_i f)$$

$$= ((s_1(X'_i Y))^2),$$

kun $|f| = 1$. Yllä olevasta epäyhtälöstä seuraa, että koska valittu ehto on $X'_i Y_i = X'_i Y$, niin jokaisella PLS-algoritmin kierroksella matriisin $X'_i Y_i$ arvot ovat pienempiä kuin edellisen kierroksen vastaavat arvot. Tämän vuoksi PLS-algoritmissa voidaan valita PLS-komponenttipareja niin kauan, kunnes matriisin $X'_i Y$ suurin singulaariarvo $s_1(X'_i Y)$ on kyllin pieni PLS-regressiossa käytetyn kriteerin mukaan.

5.4 Ennusteet PLS-regressiossa

Höskuldsson (1988) on myös esittänyt, kuinka PLS-regressiossa annetulla x_0 voidaan ennustaa tuntematonta y_0 . Koska kappaleessa 6.6 annetulla testiaineiston rakennematriisilla X_t on ennustettu 'tuntematonta' testiaineiston havaintomatriisia Y_t , Jacksonin (1991) mukaan, niin tässä esitetään PLS-regressiolla ennustaminen Höskuldssonin (1988) ja Jacksonin (1991) mukaan.

Ennustettaessa annetulla X_0 tuntematonta Y_0 annettu X_0 pitää esittää PLS-regressiossa vastaavassa rekursiivisessa muodossa kuin rakennematriisin X approksimaatio on PLS-algoritmin askeleessa A , jotta voidaan laskea ennuste

$$Y_{0(A)} = t_{01}c'_1 + \dots + t_{0A}c'_A, \quad (6.68)$$

missä $Y_{0(A)}$ on tuntemattoman havaintomatriisin Y_0 ennuste, A on valittujen PLS-komponenttiparien lukumäärä, matriisi C on saatu jo PLS-regressiomallia sovitettaessa ja vektoreiden c pituutta ei ole skaalattu ykköseksi. Ennusteessa (6.67) vektorit t_{0i} , $i = 1, \dots, A$, saadaan Jacksonin (1991) mukaan

$$t_{0i} = X_{0i-1}w_i, \quad i = 1, \dots, A, \quad (6.69)$$

missä

$$X_{0i} = X_{0i-1} - t_{0i}p'_i, \quad i = 2, \dots, A, \quad (6.70)$$

vektorit w ja p on jo saatu PLS-regressiomallia sovitettaessa ja $X_{01} = X_0$.

Kun PLS-regressiossa rakennematriisi X aproksimoidaan matriisilla

$$X_{(A)} = \sum_{i=1}^A t_i p'_i, \quad (5.71)$$

missä A on valittujen PLS-komponenttiparien lukumäärä, niin PLS-regressiolla saatu havaintomatriisin Y sovite $Y_{(A)}$ voidaan kirjoittaa muotoon

$$Y_{(A)} = H_T Y = \left(\sum_{i=1}^A t_i t'_i / (t'_i t_i) \right) Y, \quad (5.72)$$

missä projektiomatriisi H_T on

$$H_T = \sum_{i=1}^A t_i t_i' / (t_i' t_i) = \sum_{i=1}^A X_i w_i t_i' / (t_i' t_i). \quad (5.73)$$

Projektiomatriisi H_T on symmetrinen ja idempotentti ja siten

$$\sum_i h_{ij}^2 = h_{ij}, \quad 0 \leq h_{ij} \leq 1 \quad \text{ja} \quad \sum h_{ij} = A, \quad (5.74)$$

missä h_{ij} on matriisin H_T alkio ja A on PLS-komponenttiparien lukumäärä. Projektiomatriisilla H_T on siis vastaavat ominaisuudet kuin matriisilla $H = X(X'X)^{-1}X'$, joten projektiomatriisia H_T voidaan analysoida vastaavasti kuin matriisia H lineaarisessa regressioitehtävässä.

5.5 PLS-regressiomallin ulottuvuuden valinta

Seuraavassa on muutamia valintamenetelmiä, joilla voidaan valita PLS-regressiomallin PLS-komponenttiparien lukumäärä A , joka on PLS-regressiomallin ulottuvuus.

5.5.1 PLS-regressiomallin jäännösten graafinen tarkastelu

Piirretään jäännösmatriisin Y_{A+1} normi $\|Y_{A+1}\|$ PLS-regressiomallin ulottuvuuden A funktiona. Kuvan perusteella yritetään paikallistaa mallin luonnollinen PLS-komponenttiparien lukumäärä. (Jackson, 1991.)

5.5.2 Cross-validation

Lasketaan PLS-algoritmin jokaisella kierroksella

$$\text{PRESS}(A) = \sum_{i=1}^n \sum_{j=1}^q (y_i^{(j)} - \hat{y}_{(i)}^{(j)})^2, \quad (5.75)$$

missä A on mallissa olevien PLS-komponenttiparien lukumäärä, $y_i^{(j)}$ on selitettävän muuttujan havaittu arvo ja $\hat{y}_{(i)}^{(j)}$ on tämän ennuste, joka on laskettu ilman havaintoa i . Valitaan se PLS-komponenttien lukumäärä A , jolla saadaan pienin $\text{PRESS}(A)$. Yleensä *cross-validation* kriteeri

'pysäyttää' ennen kuin kaikki mahdolliset PLS-komponentit ovat mukana mallisissa, sillä 'loput' PLS-komponenttiparit eivät paranna mallin kykyä toimia ennusteena. (Höskuldsson, 1988.)

6 Teräksen karkenevuuden mallittamisesta

Tehtävänä on muodostaa monimuuttujainen regressiomalli teräksen karkenevuudelle käyttäen estimointimenetelminä säännöllistämismenetelmiä ja vertailla käytettyjen menetelmien "luotettavuutta" keskenään. Copas (1983) ehdottaa kokonaisaineiston jaettavaksi retrospektiiviseen aineistoon (engl. retrospective sample), josta estimoidaan regressiomallin parametrit, ja prospektiiviseen tai vahvistavaan aineistoon (engl. prospective or validation sample), josta ei estimoida regressiomallin parametreja. Tässä kokonaisaineiston eli teräksen karkenevuusaineiston otoskoko on 220 havaintoa. Jotta estimodun regressiomallin sopivuutta voidaan testata myös havaitsemattomiin otospisteisiin, niin kokonaiaineisto on jaettu prospektiiviseen aineistoon, jota kutsutaan tässä testiaineistoksi, ja retrospektiiviseen aineistoon, jota kutsutaan tässä estimointiaineistoksi. Testiaineiston otoskoko on 44 ja se on valittu kokonaisaineistosta systemaattisella otannalla. Loput kokonaisaineiston havainnot, joita on 176, muodostavat estimointiaineiston.

Mallitettaessa teräksen karkenevuutta on tässä valittu selitettäviksi muuttujiksi teräksen kovuudet 1.5, 13 ja 25 mm etäisyyksiltä koesauvan alapäästä. Merkitään näitä mitattuja kovuuksia seuraavasti: $y^{*(j)}$ ja $y_i^{*(j)}$, $j = 1.5, 13, 25$ mm, missä muuttujat $y^{*(j)}$ ja $y_i^{*(j)}$ ovat mitatut teräksen kovuudet etäisyyksiltä j mm koesauvan alapäästä, muuttuja $y^{*(j)}$ kuuluu estimointiaineistoon ja muuttuja $y_i^{*(j)}$ kuuluu testiaineistoon. Muuttujien $y^{*(j)}$ ja $y_i^{*(j)}$, $j = 1.5, 13, 25$, keskiarvot ja -hajonnat ovat taulukossa 6.1.

Taulukko 6.1: Selitettävien muuttujien $y^{*(j)}$ ja $y_i^{*(j)}$, $j = 1.5, 13, 25$, keskiarvot ja -hajonnat

Tunnusluku	Estimointiaineisto			Testiaineisto		
	$y^{*(1.5)}$	$y^{*(13)}$	$y^{*(25)}$	$y_i^{*(1.5)}$	$y_i^{*(13)}$	$y_i^{*(25)}$
Keskiarvo	48.05	40.27	29.27	48.43	40.76	29.27
Keskihajonta	6.0259	9.5228	9.1767	5.9621	9.6699	8.8026

Kuten taulukosta 6.1 nähdään, niin testiaineiston selitettävien muuttujien keskiarvot ja -hajonnat eivät juurikaan eroa estimointiaineiston selitettävien muuttujien keskiarvoista ja -hajonnoista.

Tehdään selitettäville muuttujille $y^{*(j)}$ ja $y_i^{*(j)}$ muunnos (2.4), missä $\text{var}(y^{*(j)})^{1/2}$, $\bar{y}^{*(j)}$ ja n on laskettu estimointiaineistosta, $j = 1.5, 13, 25$. Tämän muunnoksen jälkeen estimointiaineiston havaintomatriisin $Y = (y^{(1.5)} y^{(13)} y^{(25)})$ neliömatriisi YY on korrelaatiomatriisi

$$Y'Y = \begin{pmatrix} 1.000 & .922 & .810 \\ .922 & 1.000 & .924 \\ .810 & .924 & 1.000 \end{pmatrix}. \quad (6.1)$$

Kuten yllä olevasta neliömatriisistä $Y'Y$ nähdään, niin selitettävät muuttujat korreloivat keskenään erittäin voimakkaasti. Neliömatriisin $Y'Y$ determinantti on 0.020 ja neliömatriisin $Y'Y$ ominaisarvot ovat 2.771, 0.190 ja 0.038. Jotta monimuuttujainen regressiomalli olisi hyödyllinen, niin selitettävien muuttujien tulee korreloida melko voimakkaasti keskenään, mikä pätee tässä aineistossa.

Teräksen karkenevuuteen vaikuttavat teräksen seosaineiden C, Si, Mn, P, S, Cr, Ni, Mo, V, Ti, Cu, Als, B, Nb ja N pitoisuudet. Edellä mainittujen seosaineiden pitoisuuksien lisäksi selittävänä muuttujana on ns. vapaa boori $B_{jab} = B + 0.212Ti - 0.715N$. Koska seosaine B_{jab} on seosaineiden B, Ti ja N lineaarikombinaatio, niin selittävät muuttujat ovat multikollineaarisia. Estimointi- ja testiaineiston selittävien muuttujien keskiarvot ja -hajonnat ovat liitteessä A taulukossa A.1. Testiaineiston selittävien muuttujien keskiarvot ja -hajonnat eivät juurikaan eroa estimointiaineiston selittävien muuttujien keskiarvoista ja -hajonnoista. Tehdään estimointiaineiston selittäville muuttujille x_j^* ja testiaineiston selittäville muuttujille x_j^* , $j = C, \dots, B_{jab}$, muunnos (2.3), missä \bar{x}_j^* , $\text{var}(x_j^*)$ ja n on laskettu estimointiaineistosta. Näiden muunnosten jälkeen rakennematriisin $X = (x_C, \dots, x_{B_{jab}})$ neliömatriisi $X'X$ on korrelaatiomatriisi, joka on liitteessä A taulukossa A.2.

Mitä enemmän selitettävien ja selittävien muuttujien välillä on lineaarista riippuvuutta, sitä luotettavampaa on muuttujilla x_j , $j = C, \dots, B_{jab}$, selittää muuttujien $y^{(j)}$, $j = 1, 5, 13, 25$, vaihtelua. Estimointiaineiston selitettävien ja selittävien muuttujien korrelaatiokertoimet ovat liitteessä A taulukossa A.3. Kuten estimointiaineiston selitettävien ja selittävien muuttujien korrelaatiomatriisista nähdään, varsinkin muuttuja C korreloi voimakkaasti selitettävien muuttujien kanssa.

Tässä estimoitava monimuuttujainen regressiomalli teräksen karkenevuudelle on

$$Y = XB + E, \quad (6.2)$$

missä selitettävien muuttujien (176×3) -havaintomatriisi Y ja selittävien muuttujien (176×16) -rakennematriisi X ovat estimointiaineistosta, B on (16×3) -regressiomatriisi ja E on (176×3) -jäännösmatriisi. Regressiomalli (6.2) voidaan kirjoittaa myös muotoon

$$y^{(j)} = X \beta^{(j)} + \epsilon^{(j)} = Z \alpha^{(j)} + \epsilon^{(j)}, \quad j = 1, 5, 13, 25, \quad (6.3)$$

missä vektori $\beta^{(j)}$ on selitettävää muuttujaa $y^{(j)}$ vastaava (16×1) -regressiovektori, $Z = XV$, $\alpha^{(j)} = V\beta^{(j)}$, V on neliömatriisin XX' ominaisvektorimatriisi ja $\epsilon^{(j)}$ on (176×1) -jäännösvektori. Vastaava kanoninen muoto (2.12) on

$$u_i^{(j)} = \sqrt{\lambda_i} \alpha_i^{(j)} + \epsilon_i^{(j)}, \quad i = 1, \dots, 16, \quad \text{ja} \quad u_i^{(j)} = \epsilon_i^{(j)}, \quad i = 17, \dots, 176, \quad (6.4)$$

$j = 1, 5, 13, 25$.

Kun regressiomalli on muodossa (6.3), niin alkuperäisten havaintojen $y^{*(j)}$, $j = 1, 5, 13, 25$ sovitteet, saadaan muunnoksella

$$\hat{y}^{*(j)} = \sqrt{n-1} \sqrt{\text{var}(y^{*(j)})} \hat{y}^{(j)} + \bar{y}^{*(j)}, \quad j = 1, 5, 13, 25, \quad (6.5)$$

missä n on estimointiaineiston otoskoko, $(\text{var}(y^{*(j)}))^{1/2}$ ja $\bar{y}^{*(j)}$ ovat estimointiaineiston muuttujan $y^{*(j)}$ keskihajonta ja keskiarvo, $\hat{y}^{(j)} = Xb^{(j)}$, $b^{(j)}$ on parametrin $\beta^{(j)}$ estimaattori tai $\hat{y}^{(j)} = Za^{(j)}$ ja $a^{(j)}$ on parametrin $\alpha^{(j)}$ estimaattori, $j = 1, 5, 13, 25$. Jos regressiomalli on määritelty kanonisen yhtälön (6.4) mukaan, niin alkuperäisten havaintojen $y^{*(j)}$, $j = 1, 5, 13, 25$, sovitteet saadaan muunnoksella

$$\hat{y}^{*(j)} = \sqrt{n-1} \sqrt{\text{var}(y^{*(j)})} T \hat{a}^{(j)} + \bar{y}^{*(j)}, \quad j = 1, 5, 13, 25, \quad (6.6)$$

missä n , $\text{var}(y^{*(j)})$ ja $\bar{y}^{*(j)}$ ovat kuten yhtälössä (6.5), matriisi T on estimointiaineiston neliömatriisin XX' ominaisvektorimatriisi, $\hat{a}^{(j)} = \Lambda_p^{-1/2} a^{(j)}$, $\Lambda_p = \text{diag}(\lambda_1, \dots, \lambda_p)$, λ_i , $i = 1, \dots, p$, ovat estimointiaineiston neliömatriisin XX' ominaisarvot ja $a^{(j)}$ on estimoitu parametrin $\alpha^{(j)}$ estimaattori, $j = 1, 5, 13, 25$.

Jos oletetaan, että regressiomalli on muotoa (6.2), niin testiaineiston havaintomatriisin Y_t ennuste on

$$\hat{Y}_t = X_t B = Z_t A, \quad (6.7)$$

missä X_t on testiaineiston rakennematriisi, $Z_t = X_t V$, matriisin V sarakevektoreina ovat estimointiaineiston neliömatriisin XX' ominaisvektorit, B ja A ovat parametrien B ja A estimaattorit, jotka on laskettu estimointiaineistosta. Ennuste (6.7) voidaan esittää myös muodossa

$$\hat{y}_t^{(j)} = X_t b^{(j)} + Z_t a^{(j)}, \quad j = 1, 5, 13, 25, \quad (6.8)$$

missä X_t ja $Z_t = X_t V$ ovat kuten yhtälössä (6.7) ja $b^{(j)}$ ja $a^{(j)}$ ovat parametrien $\beta^{(j)}$ ja $\alpha^{(j)}$ estimaattorit, jotka on laskettu estimointiaineistosta, $j = 1, 5, 13, 25$, (vrt. (6.3)). Jotta testiaineistosta saatua ennustetta voidaan verrata alkuperäiseen havaintomatriisiin Y , niin ennusteelle (6.7) tai (6.8) tehdään muunnos (6.5), missä $\hat{y}^{(j)}$ on ennuste $\hat{y}_t^{(j)}$ ja muut tarvittavat tunnusluvut ja matriisit saadaan estimointiaineistosta.

6.1 Sovitteen ja ennusteen selityskertoimet

Koska regressiomalleissa (6.3) ja (6.4) ovat parametrit α ja β tuntemattomia, niin parametrien α ja β säännöllistämismenetelmien estimaattoreiden keskineliövirheitä ei voida laskea eikä niitä siten voida verrata vastaavaan pienimmän neliösumman estimaattorin keskineliövirheeseen.

Koska tavoitteena on, että muodostettu teräksen karkenevuuden regressiomalli sopisi hyvin myös havaitsemattomiin otospisteisiin ts. regressiomallin validiteetti olisi hyvä, niin valittaessa lopullista regressiomallia sekä regressiomallin sovitteen että ennusteen tulisi olla riittäviä. Jos muodostetun regressiomallin parametrien estimaattoreilla ennustetaisiin vain estimointiaineiston havaintoja, niin mallin ennustamiskykyä olisi vaikea arvioida täsmällisesti, koska estimointiaineisto olisi käytetty kahdesti (Helland, 1988). Ennusteen ja sovitteen riittävyttä arvioidaan tässä tutkielmassa selityskertoimien avulla.

Yksimuuttujaisten regressiomallin (6.3) selityskertoimet, joita tässä kutsutaan sovitteiden selityskertoimeksi, saadaan Heikan, Minkkisen ja Taavitsaisen (1994) mukaan

$$R^{2(j)} = 1 - \frac{\sum_{i=1}^n (y_i^{(j)} - \hat{y}_i^{(j)})^2 / n}{\sum_{i=1}^{n_{total}} (y_{i_{total}}^{(j)} - \bar{y}_{total}^{(j)})^2 / n_{total}} \quad j = 1, 5, 13, 25, \quad (6.9)$$

missä n_{total} on kokonaisaineiston otoskoko, $y_{total}^{(j)}$ on kokonaisaineiston selitettävä muuttuja ja $\bar{y}_{total}^{(j)}$ sen keskiarvo, $j = 1, 5, 13, 25$. Selityskertoimessa (6.9) huomioidaan myös testiaineistossa olevien selitettävien muuttujien kokonaisvaihtelu. Koska sovitteen keskineliövirhe on

$$MSE(\hat{y}) = \sqrt{\sum_{i=1}^n (y_i - \hat{y}_i)^2 / n}, \quad (6.10)$$

niin mitä pienempi on sovituksen keskineliövirhe, sitä suurempi on sovituksen selityskerroin. Monimuuttujaisen regressiomallin selityskerroin, jota tässä kutsutaan sovituksen selityskertoimeksi, on Höskuldssonin (1988) mukaan

$$R^2 = 1 - \text{tr}((\hat{Y} - Y)'(\hat{Y} - Y)) / \text{tr}(Y'Y), \quad (6.11)$$

missä siis selitettävät muuttujat on keskistetty ja niiden varianssit ovat yhtäsuuret. Regressiomallin selityskerroin ilmaisee, kuinka suuren osan estimointiaineiston selitettävän muuttujan vaihtelusta estimoitu sovite selittää.

Koska ennusteen \hat{y}_t keskineliövirheessä on lauseke

$$E((\beta - b)'X_t'X_t(\beta - b)), \quad (6.12)$$

missä parametri β on tuntematon, niin ennusteen keskineliövirhettä ei voida laskea, kun yllä oleva lauseke on ennusteen keskineliövirheessä. Koska tässä ennusteen \hat{y}_t todellinen arvo y_t on tunnettu, niin ennusteiden (6.8) keskineliövirhe saadaan yhtälöstä

$$\text{MSE}(\hat{y}_t^{(j)}) = \sqrt{\sum_{i=1}^{n_t} (y_{it}^{(j)} - \hat{y}_{it}^{(j)})^2 / n_t}, \quad j = 1.5, 13, 25, \quad (6.13)$$

ja täten ennusteiden (6.8) selityskertoimet saadaan Heikan ym. (1994) mukaan

$$Q^{2(j)} = 1 - \frac{\sum_{i=1}^{n_t} (y_{it}^{(j)} - \hat{y}_{it}^{(j)})^2 / n_t}{\sum_{i=1}^{n_{total}} (y_{itotal}^{(j)} - \bar{y}_{total}^{(j)})^2 / n_{total}} \quad j = 1.5, 13, 25, \quad (6.14)$$

missä n_t on testiaineiston otoskoko, $y_t^{(j)}$ on testiaineiston selitettävä muuttuja ja $\hat{y}_t^{(j)}$ on tämän ennuste, n_{total} on kokonaisaineiston otoskoko, $y_{total}^{(j)}$ on kokonaisaineiston selitettävä muuttuja ja $\bar{y}_{total}^{(j)}$ on tämän keskiarvo, $j = 1.5, 13, 25$. Kun testiaineiston selitettäville muuttujille $y_t^{(j)}$, $j = 1.5, 13, 25$, on tehty muunnokset (2.3) ja (2.4), niin $\bar{y}_t^{(1.5)} = 0.0048$, $\bar{y}_t^{(13)} = 0.0042$, $\bar{y}_t^{(25)} = 0.0006$, $\text{var}(y^{(1.5)}) = 0.0748$, $\text{var}(y^{(13)}) = 0.0768$ ja $\text{var}(y^{(25)}) = 0.0725$. Laskettaessa sovituksen selityskertointa (6.11) Höskuldssonin (1988) mukaan oletetaan, että selitettävät muuttujat on keskistetty ja niiden varianssit ovat yhtäsuuria. Tämän vuoksi tehdään vielä testiaineiston selitettävälle muuttujalle $y_t^{(j)}$ ja sen ennusteelle $\hat{y}_t^{(j)}$ muunnos (2.4), missä nyt tarvittavat tunnusluvut ja matriisit saadaan testiaineistosta, $j = 1.5, 13, 25$. Tämän jälkeen lasketaan monimuuttujainen ennusteen

roin

$$Q^2 = 1 - \frac{\text{tr}((Y_t - \hat{Y}_t)'(Y_t - \hat{Y}_t))}{\text{tr}(Y_t'Y_t)}, \quad (6.15)$$

joka siis ilmaisee, kuinka suuri osa havaintomatriisin Y_t vaihtelusta on onnistuttu estimointiaineistosta estimoidulla regressiomallilla ennustamaan.

6.2 Selittävien muuttujien multikollinearisuus

Kuten aiemmin on jo todettu, niin säännöllistämismenetelmiä käytetään, kun selittävät muuttujat ovat multikollineaarisia. Vaikka ennalta jo tiedetään, että selittävien muuttujien välillä on lineaarikombinaatioiden välisiä riippuvuuksia, niin tutkitaan estimointiaineiston selittävien muuttujien $x_p, j = C, \dots, B_{jab}$, multikollinearisuuden astetta kappaleessa 2.4 esitetyn menetelmin.

1) Selittävien muuttujien neliömatriisi $X'X$, joka on siis korrelaatiomatriisi, on liitteessä A taulukossa A.2. Neliömatriisista $X'X$ nähdään, että selittävien muuttujien välillä on melko voimakkaita kahden muuttujan välisiä korrelaatioita. Näistä korrelaatiokertoimista kaksi eri korrelaatiokerrointa ovat suurempi kuin 0.7. Täten rakennematriisissa $X = (x_C, \dots, x_{B_{jab}})$ on selittäviä muuttujia, joiden multikollinearisuuden aste on korkea. Huomattavaa on, että vaikka selittävä muuttuja $x_{B_{jab}}$ on muuttujien x_B, x_{Ti} ja x_N lineaarikombinaatio, niin niiden korrelaatiokerroimet ovat $\text{cor}(x_{B_{jab}}, x_B) = 0.0271$, $\text{cor}(x_{Ti}, x_{B_{jab}}) = 0.7346$ ja $\text{cor}(x_{B_{jab}}, x_N) = -0.6724$. Täten selittävät muuttujat voivat olla multikollineaarisia, vaikka niiden korrelaatiokertoimet ovat pieniä.

2) Neliömatriisin $X'X$ ominaisarvot ovat liitteessä A taulukossa A.4. Vertailun vuoksi on myös laskettu vastaavat luvut testiaineistosta ja ne ovat myös taulukossa A.4. Vaikka selittävien muuttujien välillä on lineaarikombinaatioiden välisiä riippuvuuksia, niin tässä neliömatriisin $X'X$ pienin ominaisarvo ei ole nolla, vaan se on 0.005, mikä johtuu neliömatriisin $X'X$ ominaisarvoja approksimoitaessa käytetystä laskentatarkkuudesta. Kun tarkastellaan neliömatriisin $X'X$ pienintä ominaisarvoa $\lambda_p = 0.005$ vastaavan ominaisvektorin v_p alkioita, niin alkioit $v_{pp}, j = B_{jab}, B, Ti$ ja N ovat suurimmat, jolloin Wangin ja Chowin (1994) mukaan selittävät muuttujat $x_{B_{jab}}, x_{B_{jab}}, x_{Ti}$ ja x_N , ovat eniten multikollineaarisia, kuten olla pitääkin.

3) Neliömatriisin $X'X$ suurin ominaisarvo on 4.484 ja sen pienin ominaisarvo on 0.005. Rakennematriisin X kollineaarisuuden indeksi (2.43) on 954. Tämän perusteella voidaan sanoa, että selittävät muuttujat x_j , $j = C, \dots, Bjab$, ovat kollineaarisia. Koska selittävät muuttujat ovat multikollineaarisia, niin ne ovat myös kollineaarisia.

4) Neliömatriisin $X'X$ ominaisarvojen käänteislukujen summa on 243.63, joka on 15 kertaa suurempi kuin ortogonaalisen neliömatriisin $X'X$ käänteislukujen summa. Myös tämän perusteella voidaan sanoa, että selittävien muuttujien multikollineaarisuuden aste on korkea.

5) Neliömatriisin $X'X$ determinantti on 2.855×10^{-5} . Koska käytetystä laskentatarkkuudesta johtuen tämä determinantti ei ole nolla, vaikka todellisuudessa selittävien muuttujien välillä on lineaarikombinaatioiden välisiä riippuvuuksia, niin neliömatriisi $X'X$ on epäsingulaarinen eli sille on olemassa käänteismatriisi $(X'X)^{-1}$. Koska neliömatriisin $X'X$ determinantti on melko pieni, niin neliömatriisi $X'X$ on lähes singulaarinen.

6) Testisuure χ^2 (2.44) on 238.5 ja sen vapausasteet ovat 120. Koska tässä testisuureen χ^2 vapausasteet > 100 , niin tehdään testisuurelle χ^2 muunnos $Z = (2\chi^2)^{1/2} - (2df - 1)^{1/2}$, missä df on testisuureen χ^2 vapausasteet (Jackson, 1991). Tämä muunnos Z noudattaa normaalijakaumaa $N(0,1)$, jos nollahypoteesi eli selittävien muuttujien korreloimattomuus on tosi. Koska tässä $Z = 6.38$, niin nollahypoteesi hylätään 5% riskitasolla. Täten selittävät muuttujat x_j , $j = C, \dots, Bjab$, eivät ole korreloimattomia, vaan niiden välillä on korreloituneisuutta.

Vaikka selittävien muuttujien x_j , $j = C, \dots, Bjab$, välillä on lineaarisia kombinaatioiden välisiä riippuvuuksia, niin ominaisarvoja laskennassa käytetystä laskentatarkkuudesta johtuen neliömatriisi $X'X$ on kääntövä. Tällöin voidaan laskea regressiomalleissa (6.3) ja (6.4) parametrien $\alpha^{(j)}$ ja $\beta^{(j)}$ pienimmän neliösumman estimaatit, mutta ne eivät ole luotettavia, koska neliömatriisin $X'X$ käänteismatriisia ei todellisesti ole olemassa, $j = 1.5, 13, 25$.

6.3 Pienimmän neliösumman regressio

Regressiomallin (6.3) regressiomatriisin B pienimmän neliösumman estimaattori $B_{PNS} = (b_{PNS}^{(1.5)}, b_{PNS}^{(13)}, b_{PNS}^{(25)})$ ja sen varianssi $\text{var}(B_{PNS})$ ovat liitteessä A taulukossa A.5. Regressiomalli-

en (6.3) ja (6.4) parametrin $A = (\alpha^{(1.5)}, \alpha^{(13)}, \alpha^{(25)})$ pienimmän neliösumman estimaattori $A_{PNS} = (a_{PNS}^{(1.5)}, a_{PNS}^{(13)}, a_{PNS}^{(25)})$ ja sen varianssi $\text{var}(A_{PNS})$ ovat liitteessä A taulukossa A.6.

Kun regressiomalli (6.3) estimoidaan pienimmän neliösumman menetelmällä, niin virhevarianssin Σ_Y estimaattori (2.31) on

$$\begin{pmatrix} 0.000162 & 0.000118 & 0.000140 \\ 0.000118 & 0.000507 & 0.000363 \\ 0.000140 & 0.000363 & 0.000506 \end{pmatrix}, \quad (6.16)$$

joka on korrelaatiomuodossa

$$\begin{pmatrix} 1.0000 & 0.4123 & 0.4909 \\ 0.4123 & 1.0000 & 0.7170 \\ 0.4909 & 0.7170 & 1.0000 \end{pmatrix}.$$

Kuten yllä olevasta matriisista nähdään, niin regressiomallien $\hat{y}_{PNS}^{(j)} = Xa_{PNS}^{(j)} + e^{(j)}, j = 1.5, 13, 25$, jäännökset korreloivat melko voimakkaasti. Pienimmän neliösumman menetelmällä saatujen estimointiaineiston sovitteiden ja testiaineiston ennusteiden selityskertoimien arvot ovat taulukossa 6.2.

Taulukko 6.2: Pienimmän neliösumman menetelmällä saatujen sovitteiden ja ennusteiden selityskertoimien arvot

$R^{2(1.5)}$	$R^{2(13)}$	$R^{2(25)}$	R^2	$Q^{2(1.5)}$	$Q^{2(13)}$	$Q^{2(25)}$	Q^2
0.974	0.919	0.918	0.937	0.970	0.921	0.910	0.931

Kuten taulukosta 6.2 nähdään, niin pienimmän neliösumman menetelmällä saadaan melko korkeat selityskertoimien arvot sekä estimointiaineiston sovitteelle että testiaineiston ennusteelle. Erikoista on, että $Q^{2(13)} > R^{2(13)}$.

Koska selittävien muuttujien $x_p, j = C, \dots, B_j$, välillä on lineaarikombinaatioiden välisiä riippuvuuksia, niin todellisuudessa parametreilla $\alpha^{(j)}$ ja $\beta^{(j)}$ ei ole olemassa yksikäsitteisiä pienimmän neliösumman estimaatteja, $j = 1.5, 13, 25$. Koska käytännössä johtuen laskennassa käytettyä laskentatarkkuudesta matriisi XX on epäsingulaarinen, niin parametreille $\alpha^{(j)}$ ja $\beta^{(j)}$ voidaan laskea yksikäsitteiset pienimmän neliösumman estimaatit, $j = 1.5, 13, 25$. Koska todellisesti

käänteismatriisia $(XX)^{-1}$ ei ole olemassa, niin parametrien $\alpha^{(j)}$ ja $\beta^{(j)}$, $j = 1, 5, 13, 25$, pienimmän neliösumman estimaattien perusteella ei voida tehdä luotettavia sisällöllisiä tulkintoja estimoitavan regressiomallin rakenteesta.

Vaikka selittävien muuttujien ollessa multikollineaarisia ennusteen keskineliövirheessä $\text{var}(\mathbf{a}_{PNS})$ saattaa olla kasvanut kohtuuttoman suureksi, niin ennustaminen voidaan suorittaa kohtalaisen tarkasti, jos selittävien muuttujien multikollineaarisuusyhtälö pysyy ennustejakson aikana voimassa (Leskinen, 1977). Myös tässä esimerkkiaineistossa on saatu melko hyviä ennusteen selityskertoimia selittävien muuttujien multikollineaarisuudesta huolimatta.

6.4 Harjaregressio

Mallitettaessa teräksen karkenevuutta harjaregressiolla estimoidaan tässä regressiomalli (6.3) harjaregressiomallilla

$$Y = Z A_H(K) + E, \quad (6.17)$$

missä matriisit Y , E ja Z ovat kuten yhtälössä (6.3) ja harjaestimaattori $A_H(K) = (\mathbf{a}_H^{(1.5)}(k^{(1.5)}), \mathbf{a}_H^{(13)}(k^{(13)}), \mathbf{a}_H^{(25)}(k^{(25)}))$ on

$$\mathbf{a}_H^{(j)}(k^{(j)}) = \frac{\mathbf{u}^{(j)} \sqrt{\Lambda_p}}{\Lambda_p + k^{(j)} I_p}, \quad j = 1.5, 13, 25, \quad (6.18)$$

ja $\mathbf{u}^{(j)}$ ja Λ_p on määritelty kuten kanonisessa yhtälössä (6.4). Jotta harjaestimaattori (6.18) voidaan laskea, niin on valittava harjaparametrien $k^{(j)}$, $j = 1.5, 13, 25$, arvot. Valitaan harjaestimaattoriin (6.18) harjaparametrien $k^{(j)}$, $j = 1.5, 13, 25$, arvot kappaleissa 3.5 ja 3.7 esitettyjen valintamenetelmien mukaan. Valitut harjaparametrien arvot ovat taulukossa 6.3. Harjaestimaattori $A_H(K_{ad})$ ja sen varianssi ovat liitteessä A taulukossa A.6. Harjaestimaattorit $A_H(K_{FF})$ ja $A_H(K_{cross})$ ja niiden varianssit ovat liitteessä A taulukossa A.7.

Taulukossa 6.3 kaikissa muissa harjaparametrien $k^{(j)}$, $j = 1.5, 13, 25$, valintamenetelmissä kuin valintamenetelmässä K_{cross} valitaan harjaparametrien $k^{(j)}$, $j = 1.5, 13, 25$, arvot toisistaan riippumatta. Tällöin ei hydynnetä selitettävien muuttujien $\mathbf{y}^{(1.5)}$, $\mathbf{y}^{(13)}$ ja $\mathbf{y}^{(25)}$ yhteisinformaatiota, vaan on estimoitu kolme yksimuuttujaista harjaestimaattoria $\mathbf{a}_H^{(j)}(k^{(j)})$, $j = 1.5, 13, 25$.

Harjaregressiolla saadut estimointiaineiston soviteen ja testiaineiston ennusteen selityskertoimien arvot ovat taulukossa 6.4.

Taulukko 6.3: Harjaestimaattoriin (6.18) valitut harjaparametrien $k^{(j)}, j = 1.5, 13, 25$, arvot

k	$k^{(1.5)}$	$k^{(13)}$	$k^{(25)}$
k_{jalki}	0.03	0.07	0.09
k_{cross}	0.0046	0.0106	0.0068
k_{gross}	0.0037	0.0045	0.0037
k_{HKB}	0.0023	0.0050	0.0035
k_{MUR}	0.0013	0.0050	0.0103
k_{MLW}	0.0023	0.0077	0.0077
K_{cross}	0.0072	0.0072	0.0072

Taulukko 6.4: Harjaregressiolla saatujen sovitteiden ja ennusteiden selityskertoimien arvot

k	$R^2(1.5)$	$R^2(13)$	$R^2(25)$	R^2	$Q^2(1.5)$	$Q^2(13)$	$Q^2(25)$	Q^2
k_{jalki}	0.971	0.900	0.883	0.920	0.971	0.901	0.880	0.916
k_{cross}	0.974	0.919	0.917	0.937	0.971	0.922	0.911	0.932
k_{gross}	0.974	0.919	0.917	0.937	0.971	0.923	0.911	0.932
k_{HKB}	0.974	0.919	0.917	0.937	0.971	0.923	0.911	0.932
k_{MUR}	0.974	0.919	0.916	0.937	0.971	0.923	0.911	0.932
k_{MLW}	0.974	0.919	0.917	0.937	0.971	0.922	0.911	0.932
k_{ad}	0.974	0.918	0.917	0.937	0.971	0.921	0.908	0.931
K_{FF}	0.974	0.919	0.918	0.937	0.970	0.922	0.910	0.931
K_{cross}	0.974	0.919	0.917	0.937	0.971	0.923	0.911	0.932
$K = \theta^a$	0.974	0.919	0.918	0.937	0.970	0.921	0.910	0.931

^a Valinta $K = \theta$ vastaa regressiomallin (6.3) pienimmän neliösumman menetelmän ratkaisua.

Kun harjaparametrien $k^{(j)}, j = 1.5, 13, 25$, arvot ovat $k^{(1.5)} = 0.03$, $k^{(13)} = 0.07$ ja $k^{(25)} = 0.09$, niin harjaestimaattoreiden $a_H^{(j)}(k^{(j)}), j = 1.5, 13, 25$, arvot näyttivät vakiintuneen ilman, että sovitteiden jäännöseliösummat olisivat kasvaneet kohtuuttomasti. Kun verrataan valintoja $k_{jalki}^{(j)}$,

$j = 1.5, 13, 25$, muihin harjaparametrien $k^{(j)}$, $j = 1.5, 13, 25$, valintoihin, niin harjajäljellä valitut harjaparmetrien arvot ovat suuria. Kun verrataan taulukossa 6.4 olevien sovitteiden ja ennusteiden selityskertoimien arvoja, niin harjajäljellä valitut harjaparametrien $k^{(13)}$ ja $k^{(25)}$ arvot antavat huonoimmat sovitteiden ja ennusteiden selityskertoimien arvot, minkä perusteella voidaan olettaa, että harjajäljellä valitut harjaparmetrien $k^{(13)}$ ja $k^{(25)}$ arvot ovat epäonnistuneet.

Kun harjaparametrien $k^{(j)}$, $j = 1.5, 13, 25$, arvot valitaan *cross-validation* -menetelmällä, niin laskettaessa ennusteen virheneliösummaa PRESS käytetään apuna Wangin ja Chowin (1994) esittämää lauseketta

$$\text{PRESS} = \sum_{i=1}^n \frac{e_i^2}{(1 - h_{ii})^2}, \quad (6.19)$$

missä $\text{diag}(\mathbf{H}) = \text{diag}(\mathbf{X}(\mathbf{X}\mathbf{X}^{-1}\mathbf{X})) = \text{diag}(h_{ii})$, $i = 1, \dots, n$, jolloin ennusteen virheneliösumma $\text{PRESS}(k^{(j)})$ voidaan kirjoittaa muotoon

$$\text{PRESS}(k^{(j)}) = \sum_{i=1}^n \frac{(e_i^{(j)})^2}{(1 - h_{ii}^{(j)})^2}, \quad j = 1.5, 13, 25, \quad (6.20)$$

missä $e^{(j)} = \mathbf{y}^{(j)} - \mathbf{Z}\mathbf{a}_H^{(j)}(k^{(j)})$, $\text{diag}(\mathbf{H}) = \text{diag}(\mathbf{Z}(\mathbf{\Lambda}_p + k\mathbf{I}_p)^{-1}\mathbf{Z}) = \text{diag}(h_{ii})$, $i = 1, \dots, n$, ja matriisit \mathbf{Z} ja $\mathbf{\Lambda}_p$ on määritelty kuten yhtälössä (6.3). Vaikka valittaessa harjaparametrin k arvoa *cross-validation* -menetelmällä valitaan harjaparametrin k arvo siten, että sillä saadaan pienin ennusteen virheneliösumma (6.20), niin harjaparametrin $k^{(13)}$ valinta $k_{\text{cross}}^{(13)} = 0.0106$ ei anna parasta ennusteen selityskertoimen $Q^{2(13)}$ arvoa. Tämä saattaa johtua siitä, että kun $0.0068 \leq k^{(13)} \leq 0.0106$, niin $\text{PRESS}(k^{(13)}) \approx 0.1035$ ja täten valinta $k^{(13)} = 0.0106$ saattaa olla liian suuri.

Estimointiaineistossa neliömatriisin $\mathbf{X}\mathbf{X}$ ominaisarvojen välillä on suuria eroja ja täten kollineaarisuuden indeksi on suuri. Koska epäyhtälö (3.34) on epätosi, niin $\text{MSE}(\mathbf{a}_H^{(j)}(k_{HKB}^{(j)})) \geq \text{MSE}(\mathbf{a}_{PNS}^{(j)})$, $j = 1.5, 13, 25$. Koska epäyhtälö (3.35) on epätosi, niin $\text{MSE}(\hat{\mathbf{y}}_{0HKB}^{(j)}) \geq \text{MSE}(\hat{\mathbf{y}}_{0PNS}^{(j)})$, $j = 1.5, 13, 25$. Koska kollineaarisuuden indeksi on suuri, niin voidaan olettaa, että valinnat $k_{HKB}^{(j)}$, $j = 1.5, 13, 25$, ovat liian pieniä ja harjaparametrin k valintamenetelmä k_{HKB} ei sovi käytettyyn estimointiaineistoon. Koska $k_{MUR}^{(j)} \leq k_{HKB}^{(j)}$, $j = 1.5, 13$ ja $k_{\text{gross}}^{(13)} < k_{HKB}^{(13)}$, niin voidaan olettaa, että myös valinnat $k_{MUR}^{(1.5)}$, $k_{MUR}^{(13)}$ ja $k_{\text{gross}}^{(13)}$ ovat liian pieniä. Koska tässä aineistossa kollineaarisuuden indeksi on suuri, niin Brownin (1993) mukaan harjaparametrin k valintaa k_{MLW} tulee pitää parempana kuin valintaa k_{HKB} . Valinnat $k_{MLW}^{(j)}$, $j = 13, 25$, ovatkin suurempia kuin $k_{HKB}^{(j)}$, $j = 13, 25$, ja $k_{HKB}^{(1.5)} = k_{MLW}^{(1.5)}$, mutta $Q_{MLW}^{2(13)} = 0.922 < Q_{HKB}^{2(13)} = 0.923$.

parempana kuin valintaa k_{HKB} . Valinnat $k_{MLW}^{(j)}$, $j = 13, 25$, ovatkin suurempia kuin $k_{HKB}^{(j)}$, $j = 13, 25$, ja $k_{HKB}^{(1.5)} = k_{MLW}^{(1.5)}$, mutta $Q^{2(13)}_{MLW} = 0.922 < Q^{2(13)}_{HKB} = 0.923$.

Kun harjaparametrimatriisin K estimaattori valitaan Brownin ja Paynen (1975) mukaan, niin ongelmana on, että parametri σ_β^2 on tässä tuntematon. Parametrin σ_β^2 estimaattori on σ_Y^2/k (Brown & Payne, 1975). Koska tässä on myös harjaparametri k tuntematon, niin tässä tutkielmassa tätä harjaparametrimatriisin K valintamenetelmää ei käytetä.

Kun harjaparametrimatriisin K valitaan Scloven säännön (3.52) mukaan, niin tässä epäyhtälö (3.53) on epätosi. Tällöin $MSE(A_{PNS}) < MSE(A_H(K_{Sclove}))$. Täten harjaparametrimatriisin K valintamenetelmä Scloven säännön mukaan ei sovi tässä käytettyyn estimointiaineistoon.

Kun harjaparametrimatriisin K valitaan Hoerlin, Kennardin ja Baldwinin säännön mukaan, niin tässä epäyhtälö (3.57) on epätosi, minkä mukaan $MSE(A_{PNS}) < MSE(A_H(K_{HKB}))$. Täten harjaparametrimatriisin K valintamenetelmä Hoerlin, Kennardin ja Baldwinin säännön mukaan ei myöskään sovi käytettyyn estimointiaineistoon.

Kun verrataan pienimmän neliösumman menetelmällä ja harjaregressiolla saatuja sovitteiden selityskertoimien (6.9) ja (6.11) arvoja keskenään, niin harjaregressiolla saadut sovitteiden selityskertoimien arvot eivät ole juurikaan pienempiä kuin pienimmän neliösumman menetelmällä saadut sovitteiden selityskertoimien arvot. Täten harjaregressiossa pienimmän neliösumman estimaattoreissa $b_{PNS}^{(j)}$, $j = 1.5, 13, 25$, harjaparametrin k lisääminen neliömatriisin XX diagonaalille ei juurikaan kasvata sovitteiden jäännöseliösummia, sillä valitut harjaparametrit olivat hyvin pieniä. Toisaalta harjaregressiolla on saatu keskimäärin hiukan parempia ennusteen selityskertoimien (6.14) ja (6.15) arvoja kuin pienimmän neliösumman menetelmällä, joten kun harjaparametri k on valittu sopivasti, niin harjaregressiolla saadaan hieman tarkempia ennusteita kuin pienimmän neliösumman menetelmällä.

Valitaan harjaestimaattoriin (6.17) sopivimpana harjaparametrimatriisin K valintamenetelmä valintamenetelmä K_{CROSS} . Valintamenetelmällä K_{CROSS} valittu harjaparametrin $k^{(j)}$ arvo on suurempi kuin $k_{HKB}^{(j)}$, joka saattoi olla liian pieni, koska selittävien muuttujien kollineaarisuuden indeksi on suuri, $j = 1.5, 13, 25$. Kun valitaan harjaestimaattoriin (6.17) harjaparametrimatriisi K valintamenetelmällä K_{CROSS} , niin tällöin saadut sovitteiden ja ennusteiden selityskertoimien arvot ovat hyviä, kun niitä verrataan muihin sovitteiden ja ennusteiden selityskertoimien arvoihin taulukossa 6.4

6.5 Pääkomponenttiregressio

Kun teräksen karkenevuutta mallitetaan pääkomponenttiregressiolla, niin approksimoidaan regressiomalleissa (6.3) ja (6.4) parametreja $\alpha^{(j)}$, $j = 1.5, 13, 25$, pääkomponenttiestimaattoreilla

$$a_{iAPK}^{(j)} = f_i^{(j)} a_{iPNS}^{(j)}, \quad i = 1, \dots, 16, \quad j = 1.5, 13, 25, \quad (6.21)$$

missä $f_i^{(j)} = 1$, kun $i \in S_A^{(j)}$, $f_i^{(j)} = 0$, kun $i \notin S_A^{(j)}$, $j = 1.5, 13, 25$. Pienimmän neliösumman estimaattorit $a_{iPNS}^{(j)}$, $j = 1.5, 13, 25$, ovat liitteessä A taulukossa A.6. Pääkomponenttiestimaattorin $a_{iAPK}^{(j)}$ kertoimet $f_i^{(j)}$, $j = 1.5, 13, 25$, kun valittujen pääkomponenttien indeksijoukot $S_A^{(j)}$, $j = 1.5, 13, 25$, on valittu kappaleissa 4.4 ja 4.6 olevien valintamenetelmien mukaan, ovat taulukossa 6.5

Kun pääkomponenttiregressio perustuu rakennematriisin X korrelaatiomatriisiin, niin rakennematriisin X selityskerroin R_X^2 saadaan yhtälöstä

$$R_X^2 = \sum_{j \in S_A} \lambda_j / r, \quad (6.22)$$

missä r on rakennematriisin X aste. Pääkomponenttiregressiossa rakennematriisin selityskerroin R_X^2 kertoo, kuinka paljon rakennematriisin X kokonaisvaihtelusta valitut pääkomponentit z_j , $j \in S_A$, selittävät eli kuinka paljon pääkomponenttiregressiossa rakennematriisin X vaihtelusta on käytetty selittämään selitettävää muuttujaa y . Rakennematriisin X selityskertoimien arvot ovat taulukossa 6.6.

Kun valittujen pääkomponenttien indeksijoukko S_A valitaan *cross-validation* menetelmällä, niin laskettaessa ennusteen virheneliösummia (4.25) ja (4.31) käytetään apuna Wangin ja Chowin (1994) lauseketta (6.19), jolloin ennusteen virheneliösumma $\text{PRESS}(S_A^{(j)})$ voidaan kirjoittaa muotoon

$$\text{PRESS}(S_A^{(j)}) = \sum_{i=1}^n \frac{(e_i^{(j)})^2}{(1 - h_{ii}^{(j)})^2}, \quad j = 1.5, 13, 25, \quad (6.23)$$

missä $e^{(j)} = \mathbf{y}^{(j)} - \mathbf{Z}_A^{(j)} \mathbf{a}_{APK}^{(j)}$, $\mathbf{H}^{(j)} = \text{diag}(\mathbf{Z}_A^{(j)} (\mathbf{Z}_A^{(j)\top} \mathbf{Z}_A^{(j)})^{-1} \mathbf{Z}_A^{(j)}) = \text{diag}(h_{11}^{(j)}, \dots, h_{nn}^{(j)})$, $\mathbf{Z}_A^{(j)} = (\mathbf{z}_i^{(j)})$ ja $i \in S_A^{(j)}$, $j = 1.5, 13, 25$. Kuvassa 6.1 on lasketut ennusteen virheneliösummat $\text{PRESS}(S_A^{(j)})$, $j = 1.5, 13, 25$, ja $\text{PRESS}(S_A^{(j)})/3$, missä $S_A^{(j)} = S_A^{(1.5)} = S_A^{(13)} = S_A^{(25)}$, valittujen pääkomponenttien lukumäärän funktiona.

Taulukko 6.5: Pääkomponenttimestimaattorin (6.12) kertoimet $f^{(j)}$, $j = 1, 5, 13, 25$

S_A	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9	f_{10}	f_{11}	f_{12}	f_{13}	f_{14}	f_{15}	f_{16}
$S_{A(4.20)}^{(1.5)}$	1	1	1	1	1	1	0	1	1	1	0	1	1	1	1	0
$S_{A(4.20)}^{(13)}$	1	1	1	1	1	1	0	1	1	1	0	1	1	1	0	1
$S_{A(4.20)}^{(25)}$	1	1	1	1	1	1	1	0	0	1	0	0	1	1	1	1
$S_{A(4.21)}^{(1.5)}$	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
$S_{A(4.21)}^{(13)}$	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
$S_{A(4.21)}^{(25)}$	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
$S_{A(4.24)}^{(1.5)}$	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
$S_{A(4.24)}^{(13)}$	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
$S_{A(4.24)}^{(25)}$	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
$S_{A(4.25)}^{(1.5)}$	1	1	1	1	1	1	0	1	1	1	0	0	1	1	1	0
$S_{A(4.25)}^{(13)}$	1	1	1	1	1	1	0	1	1	1	0	1	1	1	0	1
$S_{A(4.25)}^{(25)}$	1	1	1	1	1	1	1	0	0	1	0	0	1	1	1	1
$S_{A(4.26)}^{(1.5)}$	1	1	1	1	1	1	0	0	0	1	0	0	1	1	1	0
$S_{A(4.26)}^{(13)}$	1	1	1	0	1	0	0	0	0	1	0	0	1	1	0	0
$S_{A(4.26)}^{(25)}$	1	1	1	0	1	0	0	0	0	1	0	0	1	1	1	0
$S_{A(4.31)}$	1	1	1	1	1	1	0	0	1	1	0	0	1	1	1	0

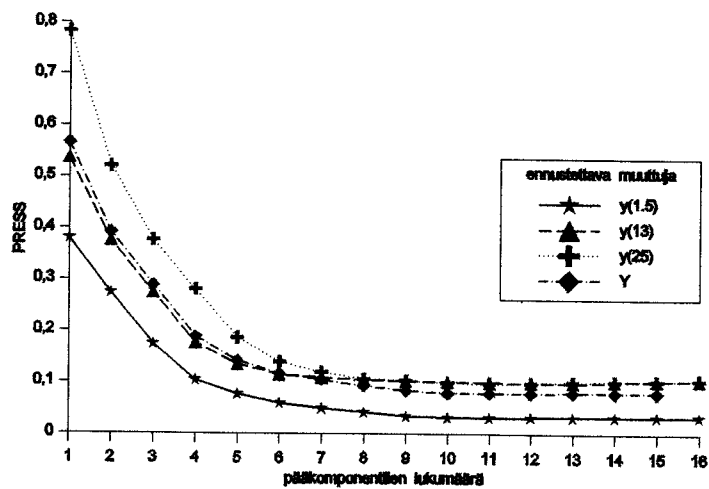
Valittaessa pääkomponenttien indeksijoukkoa S_A *cross-validation* -menetelmällä vaarana voi olla pääkomponenttiregressiomallin ylisovittaminen ts. liian monen pääkomponentin valitseminen indeksijoukkoon S_A (Jackson, 1991). Tämän vuoksi indeksijoukko S_A on valittu myös valintamenetelmän (4.36) mukaan. Kun verrataan indeksijoukkoja $S_{A(4.35)}$ ja $S_{A(4.36)}$, niin valittujen pääkomponenttien lukumäärä on indeksijoukoissa $S_{A(4.25)}^{(1.5)} = 12$, $S_{A(4.26)}^{(1.5)} = 10$, $S_{A(4.25)}^{(13)} = 13$, $S_{A(4.26)}^{(13)} = 7$, $S_{A(4.25)}^{(25)} = 12$ ja $S_{A(4.26)}^{(25)} = 8$.

Pääkomponenttiregressiolla saatujen estimointiaineiston sovitteiden ja testiaineiston ennusteiden selityskertoimien arvot ovat taulukossa 6.7

Taulukko 6.6: Pääkomponenttiregressiossa rakennematriisin X selityskertoimien arvot ja valittujen pääkomponenttien lukumäärä A

S_A	$R_X^2 (1.5)$	$A^{(1.5)}$	$R_X^2 (13)$	$A^{(13)}$	$R_X^2 (25)$	$A^{(25)}$
$S_{A(4.20)}$	0.913	13	0.906	13	0.852	12
$S_{A(4.21)}$	0.779	7	0.779	7	0.779	7
$S_{A(4.24)}$	1.000	16	1.000	16	1.000	16
$S_{A(4.25)}$	0.885	12	0.913	13	0.852	12
$S_{A(4.26)}$	0.800	10	0.636	7	0.646	8
$S_{A(4.31)}$	0.842	11	0.842	11	0.842	11
S_p^A	1.000	16	1.000	16	1.000	16

^a Pienimmän neliösumman menetelmä



Kuva 6.1: *Cross-validation* -menetelmässä lasketut ennusteen virheneliösummat $PRESS(S_A^{(j)})$, $j = 1.5, 13, 25$, ja $PRESS(S_A^{(j)})/3$ valittujen pääkomponenttien $z_p, j \in S_A^{(j)}$, lukumäärän funktiona

Taulukko 6.7: Pääkomponenttiregressiolla saatujen sovitteiden ja ennusteiden selityskertoimien arvot

S_A	$R^{2(1.5)}$	$R^{2(13)}$	$R^{2(25)}$	R^2	$Q^{2(1.5)}$	$Q^{2(13)}$	$Q^{2(25)}$	Q^2
$S_{A(4.20)}$	0.974	0.918	0.917	0.936	0.970	0.921	0.910	0.931
$S_{A(4.21)}$	0.781	0.665	0.553	0.669	0.830	0.679	0.593	0.689
$S_{A(4.24)}$	0.974	0.919	0.918	0.937	0.970	0.921	0.910	0.931
$S_{A(4.25)}$	0.973	0.918	0.917	0.935	0.972	0.915	0.905	0.928
$S_{A(4.26)}$	0.973	0.902	0.905	0.927	0.972	0.909	0.899	0.924
$S_{A(4.31)}$	0.973	0.912	0.914	0.934	0.973	0.923	0.904	0.931
S_P^a	0.974	0.919	0.918	0.937	0.970	0.921	0.910	0.931

^a Pienimmän neliösumman menetelmä

Kun verrataan pääkomponenttiregressiossa saatuja sovitteiden selityskertoimien $R^{2(j)}$, $j = 1.5, 13, 25$, ja R^2 ja ennusteiden selityskertoimien $Q^{2(j)}$, $j = 1.5, 13, 25$, ja Q^2 arvoja taulukossa 6.7, niin näiden selityskertoimien arvot erovat melkoisesti eri indeksijoukon S_A valintamenetelmien mukaan. Teräksen karkenevuusaineistossa saadaan pääkomponenttiregressiolla parhaimmat selityskertoimien R^2 , Q^2 , $R^{2(j)}$ ja $Q^{2(j)}$ $j = 1.5, 13, 25$, arvot, kun pääkomponenttiestimaattorissa (6.21) valitaan indeksijoukko S_A valintamenetelmillä $S_{A(4.20)}$, $S_{A(4.24)}$, $S_{A(4.31)}$ ja pienimmän neliösumman menetelmällä eli indeksijoukolla S_P . Kun näistä valittujen pääkomponenttien indeksijoukon S_A valintamenetelmistä valitaan sopivin pääkomponenttiestimaattoriin (6.21), niin ennusteiden selityskertoimien arvojen perusteella valitaan valintamenetelmä $S_{A(4.20)}$. Valintamenetelmä $S_{A(4.31)}$ ei valittu, koska ennusteen selityskertoimen $Q^{2(25)}$ arvo on pieni verrattuna valintamenetelmällä $S_{A(4.20)}$ saatuun ennusteen selityskertoimen $Q^{2(25)}$ arvoon.

6.6 Osittainen pienimmän neliösumman regressio

Kun teräksen karkenevuutta mallitetaan osittaisella pienimmän neliösumman regressiolla eli PLS-regressiolla, niin malli on

$$X = TP' + X_{A+1} \quad \text{ja} \quad Y = TC' + Y_{A+1}, \quad (6.24)$$

missä matriisit X ja Y ovat kuten regressiomallissa (6.3), matriisit $T = (t_1, \dots, t_A)$, $C = (c_1, \dots, c_A)$ ja $P = (p_1, \dots, p_A)$ saadaan PLS-algoritmin mukaan, matriisit X_{A+1} ja Y_{A+1} ovat jäännösmatriiseja ja valittujen PLS-komponenttiparien lukumäärä on A (ks. kappale 5.1). Mallissa (6.24) kovarianssimatriisin $X'Y$ selityskerroin R^2_{XY} saadaan yhtälöstä

$$R^2_{XY} = 1 - \text{tr}(X'_{A+1}Y_{A+1}Y'_{A+1}X_{A+1}) / \text{tr}(X'Y'X) \quad (6.25)$$

(Höskuldsson, 1988). Tämä selityskerroin R^2_{XY} kertoo, kuinka paljon malli (6.24) on selittänyt kovarianssimatriista $X'Y$. Rakennematriisin X selityskerroin R^2_X saadaan yhtälöstä

$$R^2_X = 1 - \text{tr}(X'_{A+1}X_{A+1}) / \text{tr}(X'X) \quad (6.26)$$

(vrt. (6.22)) (Höskuldsson, 1988). Tämä selityskerroin R^2_X kertoo, kuinka paljon rakennematriisin X varianssista on mallissa (6.24) käytetty selittämään havaintomatriisia Y . PLS-komponenttien lukumäärä A kappaleen 5.5 valintamenetelmin sekä vastaavat kovarianssi- ja rakennematriisin selityskertoimet R^2_{XY} ja R^2_X ovat taulukossa 6.8.

Taulukko 6.8: PLS-regressiomallin PLS-komponenttien lukumäärä A sekä vastaavat kovarianssi- ja rakennematriisin selityskertoimien arvot

Valintamenetelmä	PLS-komp. lkm. A	R^2_X	R^2_{XY}
A_{cross}	9	0.773	1.000
$A_{\ Y\ }$	12	0.875	1.000
S_p^a	16	1.000	1.000

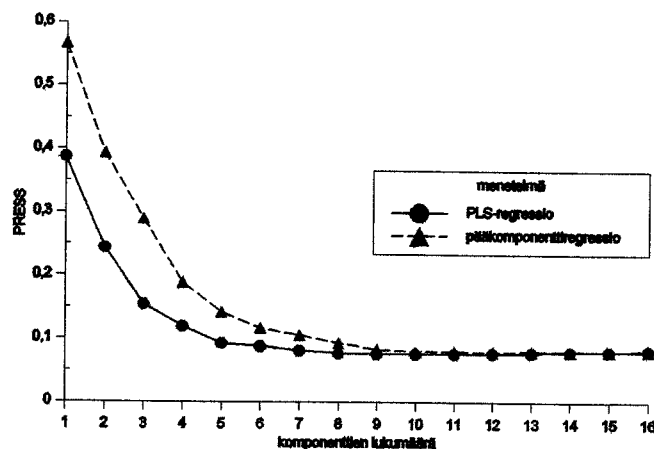
^a Asetettu PLS-komponenttien lukumääräksi 16

PLS-regressiolla saatujen estimointiaineiston sovitteiden ja testiaineiston ennusteiden selityskertoimien arvot ovat taulukossa 6.9.

Kun regressiomalli (6.3) PLS-regressiolla ja tarvittavien PLS-komponenttien lukumäärä valitaan *cross-validation* menetelmällä, niin käytetään Wangin ja Chowin (1994) lauseketta (6.19) ennusteen virheneliösummalle, jolloin

$$\text{PRESS}(A) = \sum_{j=1}^3 \sum_{i=1}^{176} \frac{y_{A+1j} - \hat{y}_{A+1j}}{(1 - h_{ii})^2}, \quad (6.27)$$

missä y_{A+1j} on havainnon y_j jäännös ja H on projektiomatriisi (5.73) PLS-algoritmin askeleessa A ja $\text{diag}(H) = \text{diag}(h_{ii})$, $i = 1, \dots, 176$. PLS-regressiolla ja pääkomponenttiregressiolla saatujen ennusteiden virheneliösummat, jotka on jaettu selitettävien muuttujien lukumäärällä kolme, ovat valittujen komponenttien lukumäärän A funktiona kuvassa 6.2.



Kuva 6.2: PLS-regressiolla ja pääkomponenttiregressiolla saatujen ennusteiden virheneliösummat

Taulukko 6.9: PLS-regressiolla saatujen soviteiden ja ennusteiden selityskertoimien arvot

A	$R^{2(1.5)}$	$R^{2(13)}$	$R^{2(25)}$	R^2	$Q^{2(1.5)}$	$Q^{2(13)}$	$Q^{2(25)}$	Q^2
A_{CROSS}	0.972	0.917	0.911	0.933	0.912	0.917	0.841	0.902
$A_{ Y }$	0.974	0.918	0.917	0.937	0.970	0.924	0.913	0.934
S_p^a	0.974	0.919	0.918	0.937	0.970	0.921	0.910	0.931

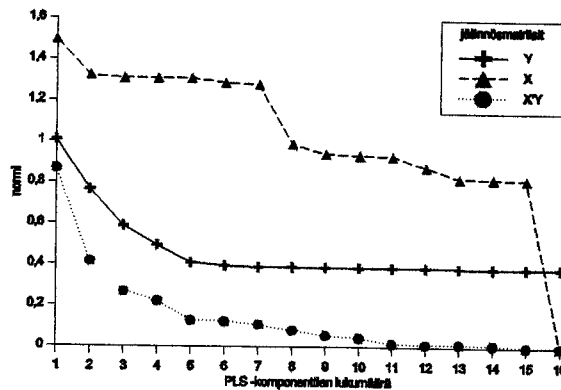
^a Asetettu PLS-komponenttien lukumääräksi 16

Kun verrataan PLS-regressiolla saatua ennusteen virheneliösummaa (6.27) monimuuttujaisella pääkomponenttiregressiolla saatuun ennusteen virheneliösummaan (6.20), niin edellinen pienenee nopeammin kuin jälkimmäinen komponenttien lukumäärän A funktiona. PLS-regressios-

sa saavutetaan ennusteen virheneliösumman minimi yhdeksällä komponentilla, kun vastaavasti pääkomponenttiregressiossa saavutetaan ennusteen virheneliösumman minimi 11 komponentilla. Täten PLS-regressiomallin ulottuvuus on pienempi kuin pääkomponenttiregressiomallin ulottuvuus, kun komponentit on valittu *cross-validation* -menetelmällä.

Kun PLS-komponenttien lukumäärä valitaan *cross-validation* -menetelmällä, niin $R^{2(25)} = 0.911$, $Q^{2(1.5)} = 0.912$ ja $Q^{2(25)} = 0.841$ ovat pieniä verrattuna muilla menetelmillä saatuihin vastaaviin selityskertoimien arvoihin. Täten PLS-regressiomalli, jossa PLS-komponenttien lukumäärä valitaan *cross-validation* -menetelmällä, ei anna parasta testiaineiston ennustetta.

Jäännösmatriisien Y_{A+1} , X_{A+1} ja $X_{A+1}'Y_{A+1}$ normit PLS-komponenttien lukumäärän A funktiona ovat kuvassa 6.3. Koska PLS-algoritmi perustuu epäyhtälöön (5.65), niin $\|X_A'Y_A\| \leq \|X_{A+1}'Y_{A+1}\|$ (ks. kuva 6.3).



Kuva 6.3: PLS-regressiossa jäännösmatriisien Y_{A+1} , X_{A+1} ja $X_{A+1}'Y_{A+1}$ normit PLS-komponenttien lukumäärän A funktiona.

Kuten kuvasta 6.3 nähdään, niin normi $\|X_{A+1}'Y_{A+1}\|$ vakiintuu PLS-komponenttien lukumäärällä 6. Normi $\|X_{A+1}\|$ vakiintuu PLS-komponenttien lukumäärällä 9, paitsi $\|X_{A+1}\| = 0$, kun $A = 16$. Valitaan kuvan 6.3 perusteella PLS-komponenttien lukumääräksi $A = 12$, jolloin normi $\|Y_{A+1}\|$ on vakiintunut ja tällöin myös normi $\|X_{A+1}'Y_{A+1}\|$ on vakiintunut.

Kun PLS-regressiossa valitaan PLS-komponenttien lukumääräksi A selittävien muuttujien lukumäärä p , niin $X_{p+1} = \mathbf{0}$. Kun mallissa (6.24) PLS-komponenttien lukumääräksi on valittu $A = 12$, niin saadut sovitteiden selityskertoimien arvot ovat lähes yhtä suuret kuin PLS-komponenttien lukumäärällä $A = 16$ saadut sovitteiden selityskertoimien arvot. Koska PLS-komponenttien lukumäärällä $A = 12$ saadaan hiukan suuremmat testiaineiston ennusteiden selityskertoimien arvot

kuin valinnalla $A = 16$, niin PLS-komponenttien lukumäärällä $A = 16$ ylisovitetään mallia (6.24).

Kun verrataan taulukossa 6.9 kaikkia PLS-regressiolla saatuja sovitteiden ja ennusteiden selityskertoimien arvoja, niin valitaan PLS-komponenttien lukumäärän A valintamenetelmäksi jäännösmatriisin Y_{A+1} normin graafinen tarkastelu, koska tällä valintamenetelmällä saadaan parhaimmat testiaineiston ennusteiden selityskertoimien arvot. Täten valitaan mallin (6.24) PLS-komponenttien lukumääräksi $A = 12$.

7 Yhteenveto

Tässä tutkielmassa tehtävänä oli muodostaa monimuuttujainen regressiomalli teräksen karkenevuudelle. Muodostettava regressiomalli oli monimuuttujainen, koska Jominy-kokeessa karkaistaessa terästä sen kovuuteen vaikuttaa etäisyys terässauvan alapäästä ja tavoitteena oli mallittaa teräksen karkenevuutta eri etäisyyksiltä koesauvan alapäästä. Muodostettavassa regressiomallissa selittävinä muuttujina olivat teräksen seosaineiden pitoisuudet, jotka olivat multikollineaarisia, sillä selittävä muuttuja x_{Bjab} oli seosaineiden B, Ti ja N pitoisuuksien linearikombinaatio. Tällöin seosaineiden pitoisuuksien rakennematriisin neliömatriisin käänteismatriisia ei ole olemassa, mutta laskusuorituksissa olevasta pyöristysvirheistä johtuen, tämä käänteismatriisi saatiin. Koska laskettaessa regressiomallin regressiokertoimien pienimmän neliösumman estimaatteja käytettiin tätä teräksen seosaineiden pitoisuuksien neliömatriisin käänteismatriisin approksimaatiota, niin regressiokertoimien pienimmän neliösumman estimaatit eivät ole luotettavia.

Harjaregressiolla, pääkomponenttiregressiolla ja osittaisella pienimmän neliösumman regressiolla approksimoitiin paremmin kääntyviä teräksen seosaineiden pitoisuuksien rakennematriisin neliömatriiseja. Näillä menetelmillä muodostettuja teräksen karkenevuuden regressiomallien riittävyttä arvioitiin soviteen ja ennusteen selityskertoimien avulla. Kun estimoitu teräksen karkenevuuden regressiomalli sovitettiin estimointiaineistoon eli havaintoihin, joista regressiokertoimen estimaatit oli laskettu, niin regressiomallin soviteen selityskertoimen arvoksi saatiin pienimmän neliösumman menetelmällä 0.937. Tämä soviteen selityskerroin on suurin mahdollinen saaduista soviteen selityskertoimista, koska pienimmän neliösumman menetelmä minimoi soviteen ja havaintojen välisen jäännöseliösumman. Kun estimoitu teräksen karkenevuuden regressiomalli oli sovitettu testiaineistoon, jonka havainnot eivät olleet mukana estimoitaessa regressiokertoimien estimaatteja, niin ennusteen suurimmat selityskertoimien Q^2 arvot olivat pienimmän neliösumman menetelmällä $Q^2 = 0.931$, harjaregressiolla $Q^2 = 0.932$, pääkomponenttiregressiolla $Q^2 = 0.932$ ja PLS-regressiolla $Q^2 = 0.934$. Vaikka säännöllistämismenetelmät regressioanalyysissä perustuvat regressiokertoimien estimaattoreiden keskineliövirheen minimointiin ja teräksen karkenevuuden regressiomallissa regressiokertoimien arvot olivat tuntemattomia, niin tästä huolimatta säännöllistämismenetelmillä pystyttiin kutistamaan regressiokertoimien pienimmän neliösumman estimaattoria kohti origoa ja vähentämään regressiomallin ulottuvuutta siten, että säännöllistämismenetelmillä saatiin

keskimäärin hiukan tarkempia selitettävien muuttujien ennusteita havaitsemattomissa otospisteissä kuin pienimmän neliösumman menetelmällä. Koska PLS-regressiolla saatiin suurin ennusteen selityskertoimen Q^2 arvo, niin teräksen karkenevuuden regressiomalliksi valitaan malli (6.24), missä PLS-komponenttien lukumäärä on 12, joka on pienempi kuin selittävien muuttujien lukumäärä 16. Koska tällä mallilla on suurin ennusteen selityskertoimen Q^2 arvo, niin sillä on Copasin (1983) mukaan paras validiteetti muodostetuista teräksen karkenevuuden regressiomalleista.

Lähteet

- Brown, P. J. & Payne C. (1975): Election Night Forecasting. *Journal of the Royal Statistical Society A* **4**, 463 - 482.
- Brown, P. J. & Zidek J.V. (1980): Adaptive Multivariate Ridge Regression. *The Annals of Statistics* **8**, 64 - 74.
- Brown, P. J. & Zidek J. V. (1982): Multivariate Regression Shrinkage Estimators with Unknown Covariance Matrix. *Scandinavian Journal of Statistical* **9**, 209 - 215.
- Brown, P. J. (1993): *Measurement, Regression, and Calibration*. Oxford: Clarendoon Press.
- Copas, J. B. (1983): Regression, Prediction and Shrinkage. *Journal of the Royal Statistical Society B* **3**, 311 - 354.
- Draper, N.R. & Smith, H. (1981): *Applied Regression Analysis*. New York: Wiley.
- Garthwaite, P. H. (1994): An Interpretation of Partial Least Squares. *Journal of the American Statistical Association* **89**, 122 - 127
- Gunst, R. F. & Mason, R. L. (1980): Regression Analysis and Its Applications: A Data-Oriented Approach. *Statistics: textbook and monographs* **141**.
- Frank, I. E. & Friedman, J. H. (1993): Statistical view of chemometrics regression tools. *Technometrics* **35**, 127 - 131.
- Füle, Erika (1995): On egological regression and ridge estimation. *Communications in Statistics.Simula*. **24**, 385-398.
- Haaland, D. M. & Thomas, E. V. (1988): Partial Least-Squares Methods for Spectral Analyses. 1. Relation to other Quantitative Calibration Methods and the Extraction of Quantitative Information. *Analytical Chemistry* **60**, 1193 - 1201.
- Heikka, R., Minkkinen, P. & Taavatsainen, V.-M. (1994): Comparision of variable selection and regression methods in multivariate calibration of a process analyzer. *Process Control and Quality* **6**, 47 - 54.
- Helland, I. S. (1988): On the structure of partial least squares regression. *Communications in Statistics* **17**, 581 - 607.
- Hoerl, A. E. & Kennard, R. W. (1970a): Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics* **12**, 55 - 67.

- Hoerl, A. E. & Kennard, R. W. (1970b): Ridge Regression: Applications to Nonorthogonal Problems. *Technometrics* **12**, 69 - 82.
- Huhtala, K. (1991): Teräksen karkenevuuden mallittamisesta. *Suomen tilastoseuran vuosikirja 1991*. Helsinki: Haakapaino.
- Höskulsson, A. (1988): PLS Regression Methods. *Journal of Chemometrics* **2**, 211 - 228.
- Jackson, J. E. (1991): *A User's Guide to Principal Components*. New York: John Wiley & Sons, Inc.
- Liski, E. & Puntanen, S. (1976): *Regressioanalyysin perusteita*. Tampere: Tampereen Pikakopio Oy.
- Leskinen, E. (1977): Multikollineaarisuuden ongelmista regressioanalyysissä. *Jyväskylän yliopiston tilastotieteen laitoksen julkaisuja* **7**.
- Leskinen, E. (1981): On Principal Component Analysis and Estimation in Regression. *Department of Statistics University of Jyväskylä*, **4**.
- Naes, T., Irgens, C. & Martens, H. (1986): Comparison of Linear Statistical Method for Calibration of NIF Instruments. *Applied Statistics* **35**, 195 - 206.
- Otto, M. & Wegscheider, W. (1985): Spectrophotometric Multicomponent Analysis Applied to Trace Metal Determinations. *Analytical Chemistry* **57**, 63 - 69.
- Stone, M. & Brooks, R. J. (1990): Continuum Regression: Cross-validated Sequentially Constructed Prediction Embracing Ordinary Least Squares, Partial Least Squares and Principal Components Regression. *Journal of the Royal Statistical Society* **B 52**, 237 - 269.
- Ranta, E. (1991): *Biometria: Tilastotiedettä ekologeille*. Helsinki: Yliopistopaino.
- Wang, S.-G. & Chow, S.-C. (1994): Advanced linear models - Theory and Applications. *Statistics: textbook and monographs* **141**.

Liite A

Taulukko A.1: Estimointi- ja testiaineiston teräksen seosaineiden pitoisuuksien keskiarvot ja -hajonnat

Seos	Estimointiaineisto		Testiaineisto	
	Keskiarvo	Keskihajonta	Keskiarvo	Keskihajonta
C	0.2547	0.0937	0.2615	0.0989
Si	0.2614	0.0369	0.2593	0.0389
Mn	1.1992	0.1580	1.1977	0.1575
P	0.0145	0.0036	0.0135	0.0029
S	0.0217	0.0157	0.0235	0.0183
Cr	0.4922	0.3225	0.4941	0.3206
Ni	0.1689	0.0795	0.1643	0.0725
Mo	0.0360	0.0270	0.0421	0.0500
V	0.0056	0.0020	0.0054	0.0020
Ti	0.0333	0.0064	0.0327	0.0072
Cu	0.1784	0.0460	0.1723	0.0449
Als	0.0272	0.0074	0.0280	0.0069
B	0.0031	0.0005	0.0031	0.0006
Mb	0.0041	0.0009	0.0039	0.0009
N	0.0100	0.0017	0.0100	0.0015
Bjab	0.0029	0.0018	0.0027	0.0022

Taulukko A.2: Teräksen seosaineiden pitoisuuksien korrelaatiomatriisi

	C	Si	Mn	P	S	Cr
C	1.0000	.0501	.7996	.0971	-.2345	-.7379
Si	.0501	1.0000	.1049	-.0685	.1419	.1634
Mn	.7996	.1049	1.0000	.0558	-.1972	-.6779
P	.0971	-.0685	.0558	1.0000	.1400	-.0700
S	-.2345	.1419	-.1972	.1400	1.0000	.4523
Cr	-.7379	.1634	-.6779	-.0700	.4523	1.0000
Ni	-.2537	-.2359	-.4550	-.1676	-.1083	.1046
Mo	-.4137	.4589	-.2799	-.1671	.4453	.6232
V	-.1465	.2389	-.0160	-.0473	.2255	.3481
Ti	.3383	-.1032	.3214	-.0334	-.3383	-.3689
Cu	-.0969	-.0524	-.1825	.2552	.2028	.1062
Als	.0365	.0616	-.0865	-.1144	.2895	.2562
B	-.3293	.1335	-.2565	.0635	.3615	.3447
Mb	.0622	.1318	.1149	-.1092	.1096	.0739
N	-.4683	.1086	-.3721	-.2238	.0413	.4627
Bjab	.4752	-.1133	.4195	.1587	-.1729	-.4859

	Ni	Mo	V	Ti	Cu	Als
C	-.2537	-.4137	-.1465	.3383	-.0969	.0365
Si	-.2359	.4589	.2389	-.1032	-.0524	.0616
Mn	-.4550	-.2799	-.0160	.3214	-.1825	-.0865
P	-.1676	-.1671	-.0473	-.0334	.2552	-.1144
S	-.1083	.4453	.2255	-.3383	.2028	.2895
Cr	.1046	.6232	.3481	-.3689	.1062	.2562
Ni	1.0000	-.0038	-.1069	.1155	-.0953	.1476
Mo	-.0038	1.0000	.4561	-.4172	.0022	.1443
V	-.1069	.4561	1.0000	-.0564	.0465	.1827
Ti	.1155	-.4172	-.0564	1.0000	-.1301	.1146
Cu	-.0953	.0022	.0465	-.1301	1.0000	-.0231
Als	.1476	.1443	.1827	.1146	-.0231	1.0000
B	.0108	.4566	.1843	-.2837	.2044	.0642
Mb	.0176	.0891	.2098	.1264	.0531	.1815
N	.2095	.4197	.2322	-.0833	-.0957	.0116
Bjab	-.0549	-.4751	-.1450	.7346	.0271	.1139

	B	Mb	N	Bjab
C	-.3293	.0622	-.4683	.4590
Si	.1335	.1318	.1086	-.0965
Mn	-.2565	.1149	-.3721	-.3990
P	.0635	-.1092	-.2238	.1607
S	.3615	.1096	.0413	-.1483
Cr	.3447	.0739	.4627	-.4611
Ni	.0108	.0176	.2095	-.0513
Mo	.4566	.0891	.4197	-.4506
V	.1843	.2098	.2322	-.1435
Ti	-.2837	.1264	-.0833	.7153
Cu	.2044	.0531	-.0957	.0343
Als	.0642	.1815	.0116	.1494
B	1.0000	.0358	.1460	.0074
Mb	.0358	1.0000	-.0109	.1031
N	.1460	-.0109	1.0000	-.6741
Bjab	.0271	.1012	-.6724	1.0000

Taulukko A.3: Teräksen seosaineiden pitoisuuksien ja teräksen kovuuksien $y^{(1.5)}$, $y^{(13)}$ ja $y^{(25)}$ korrelaatiokertoimet

Seos	Teräksen kovuudet		
	$y^{(1.5)}$	$y^{(13)}$	$y^{(25)}$
C	0.985	0.934	0.847
Si	0.38	0.51	180
Mn	0.796	0.781	0.737
P	00.096	1.000	0.057
S	-0.262	-0.305	-0.097
Cr	-0.736	-0.671	-0.370
Ni	-0.250	-0.244	-0.229
Mo	-0.427	-0.360	-0.101
V	-0.178	-0.144	0.054
Ti	0.347	0.378	0.299
Cu	-0.097	-0.041	-0.047
Als	0.003	0.040	0.208
B	-0.340	-0.276	-0.188
Mb	0.051	0.062	0.174
N	-0.444	-0.453	-0.336
Bjab	0.463	0.489	0.356

Taulukko A.4: Estimointiaineiston rakennematriisin X neliömatriisin XX ominaisarvot sekä vastaavat ominaisarvot laskettuna testiaineistolle

Indeksi	Estimointiaineisto	Testiaineisto
1	4.484	4.798
2	2.051	2.051
3	1.706	1.706
4	1.580	1.580
5	0.942	0.942
6	0.882	0.882
7	0.825	0.825
8	0.689	0.689
9	0.889	0.669
10	0.618	0.618
11	0.563	0.563
12	0.449	0.449
13	0.260	0.260
14	0.152	0.152
15	0.125	0.125
16	0.005	0.005

Taulukko A.5: Regressiomallin (6.3) parametrin $B = (\beta^{(1.5)}, \beta^{(13)}, \beta^{(25)})$ pienimmän neliösumman estimaattori B_{PNS} ja suluissa sen varianssit.

Seosaine	B_{PNS}		
	$b^{(1.5)}_{PNS}$	$b^{(13)}_{PNS}$	$b^{(25)}_{PNS}$
C	0.9800 (0.0007)	0.9611 (0.0022)	0.9389 (0.0022)
Si	-0.0047 (0.0002)	-0.0026 (0.0008)	-0.0482 (0.0008)
Mn	0.0426 (0.0008)	0.2810 (0.0025)	0.5832 (0.0025)
P	0.0009 (0.0002)	0.0011 (0.0006)	0.0215 (0.0006)
S	-0.0221 (0.0003)	-0.1836 (0.0009)	-0.1444 (0.0009)
Cr	0.0437 (0.0008)	0.4110 (0.0024)	0.7521 (0.0023)
Ni	0.0052 (0.0003)	0.0962 (0.0009)	0.1900 (0.0009)
Mo	-0.0076 (0.0005)	0.1166 (0.0017)	0.2137 (0.0017)
V	-0.0382 (0.0002)	-0.0664 (0.0007)	-0.0612 (0.0007)
Ti	-0.0537 (0.0096)	-0.1749 (0.0302)	-0.1725 (0.0301)
Cu	0.0106 (0.0002)	0.1170 (0.0006)	0.0982 (0.0006)
Als	-0.0277 (0.0002)	0.0149 (0.0008)	0.0109 (0.0008)
B	-0.0271 (0.0016)	-0.0637 (0.0051)	-0.0739 (0.0050)
Nb	-0.0015 (0.0002)	-0.0078 (0.0006)	0.0041 (0.0006)
N	0.0796 (0.0069)	0.1440 (0.0216)	0.1265 (0.0215)
Bjab	0.0829 (0.0166)	0.3522 (0.0520)	0.3048 (0.0519)

Taulukko A.6: Regressiomallin (6.3) ja (6.4) parametrin A pienimmän neliösumman estimaattori A_{PNS} ja harjaestimaattori $A_H(K_{ad})$, missä K_{ad} on valittu (3.38) mukaan, ja näiden estimaattoreiden varianssit suluissa

Indeksi	A_{PNS}			$A_H(K_{ad})$		
	$a_{PNS}^{(1.5)}$	$a_{PNS}^{(13)}$	$a_{PNS}^{(25)}$	$a_{H(K_{ad})}^{(1.5)}$	$a_{H(K_{ad})}^{(13)}$	$a_{H(K_{ad})}^{(25)}$
1	0.3731 (0.4×10^{-4})	0.3230 (1.1×10^{-4})	0.2224 (1.1×10^{-4})	0.3730 (0.3×10^{-4})	0.3226 (1.0×10^{-4})	0.2219 (1.0×10^{-4})
2	-0.2208 (0.8×10^{-4})	-0.2828 (2.5×10^{-4})	-0.3631 (2.5×10^{-4})	-0.2207 (0.7×10^{-4})	-0.2822 (2.3×10^{-4})	-0.3614 (2.2×10^{-4})
3	-0.0690 (1.0×10^{-4})	-0.1108 (3.0×10^{-4})	-0.1644 (3.0×10^{-4})	-0.0690 (0.9×10^{-4})	-0.1105 (2.7×10^{-4})	-0.1635 (2.6×10^{-4})
4	0.1319 (1.0×10^{-4})	0.0520 (3.2×10^{-4})	0.0526 (3.2×10^{-4})	0.1318 (0.9×10^{-4})	0.0519 (2.9×10^{-4})	0.0523 (2.9×10^{-4})
5	-0.1303 (1.7×10^{-4})	-0.0803 (5.4×10^{-4})	-0.1192 (5.4×10^{-4})	-0.1302 (1.6×10^{-4})	-0.0800 (4.9×10^{-4})	-0.1118 (4.9×10^{-4})
6	0.0863 (1.8×10^{-4})	0.0376 (5.7×10^{-4})	0.0679 (5.7×10^{-4})	0.0862 (1.7×10^{-4})	0.0375 (5.2×10^{-4})	0.0672 (5.21×10^{-4})
7	-0.0062 (2.0×10^{-4})	-0.0022 (6.1×10^{-4})	-0.0553 (6.1×10^{-4})	0.0062 (1.8×10^{-4})	-0.0022 (5.6×10^{-4})	-0.0546 (5.5×10^{-4})
8	-0.0259 (2.4×10^{-4})	-0.0727 (7.4×10^{-4})	-0.0144 (7.3×10^{-4})	-0.0259 (2.1×10^{-4})	-0.0722 (6.7×10^{-4})	-0.0142 (6.6×10^{-4})
9	0.0193 (2.4×10^{-4})	0.0844 (7.6×10^{-4})	0.0153 (7.6×10^{-4})	0.0193 (2.2×10^{-4})	0.0839 (6.9×10^{-4})	0.0150 (6.8×10^{-4})
10	0.4164 (2.6×10^{-4})	0.3893 (8.2×10^{-4})	0.3781 (8.2×10^{-4})	0.4157 (2.4×10^{-4})	0.3865 (7.4×10^{-4})	0.3724 (7.4×10^{-4})
11	-0.0158 (2.9×10^{-4})	-0.0231 (9.0×10^{-4})	0.0049 (9.0×10^{-4})	-0.0158 (2.6×10^{-4})	0.0230 (8.2×10^{-4})	0.0048 (8.1×10^{-4})
12	0.0235 (3.6×10^{-4})	-0.0669 (11.2×10^{-4})	-0.0181 (11.2×10^{-4})	0.0235 (3.3×10^{-4})	-0.0662 (10.2×10^{-4})	-0.0177 (10.1×10^{-4})
13	-0.0968 (6.2×10^{-4})	-0.3921 (19.4×10^{-4})	-0.6105 (19.4×10^{-4})	-0.0964 (5.7×10^{-4})	-0.3855 (17.5×10^{-4})	-0.5893 (17.1×10^{-4})
14	-0.06856 (10.7×10^{-4})	-0.8475 (33.5×10^{-4})	-0.9924 (33.4×10^{-4})	-0.6808 (9.7×10^{-4})	-0.8234 (29.7×10^{-4})	-0.9347 (28.7×10^{-4})
15	-0.2784 (13.0×10^{-4})	0.0637 (40.6×10^{-4})	0.4214 (40.1×10^{-4})	-0.2760 (11.8×10^{-4})	0.0615 (35.8×10^{-4})	0.3921 (34.5×10^{-4})
16	-0.1307 (342.5×10^{-4})	-0.4167 (1071.1×10^{-4})	-0.3744 (1068.3×10^{-4})	-0.1067 ($255.6.1 \times 10^{-4}$)	-0.2149 ($505.1.1 \times 10^{-4}$)	-0.1258 (328.13×10^{-4})

Taulukko A.7: Regressiomallin (6.3) ja (6.4) parametrin A harjaestimaattorit $A_H(K_{FF})$, missä K_{FF} on valittu (3.59) mukaan, ja $A_H(K_{CROSS})$, missä K_{CROSS} on valittu (3.61) mukaan, sekä suluissa näiden estimaattoreiden varianssit

Indeksi	$A_H(K_{FF})$			$A_H(K_{CROSS})$		
	$a^{(1.5)}(k_{FF}^{(1.5)})$	$a_H^{(13)}(k_{FF}^{(13)})$	$a_H^{(25)}(k_{FF}^{(25)})$	$a_H^{(1.5)}(0.0072)$	$a_H^{(13)}(0.0072)$	$a_H^{(25)}(0.0072)$
1	0.3730 (0.4×10^{-4})	0.3230 (1.1×10^{-4})	0.2224 (1.1×10^{-4})	0.3725 (0.4×10^{-4})	0.3225 (1.1×10^{-4})	0.2220 (1.1×10^{-4})
2	-0.2207 (0.8×10^{-4})	-0.2828 (2.5×10^{-4})	-0.3631 (2.3×10^{-4})	-0.2200 (0.8×10^{-4})	-0.2818 (2.5×10^{-4})	-0.3618 (2.4×10^{-4})
3	-0.0690 (0.9×10^{-4})	-0.1108 (3.0×10^{-4})	-0.1644 (3.0×10^{-4})	-0.0687 (1.0×10^{-4})	-0.1103 (3.0×10^{-4})	-0.1637 (3.0×10^{-4})
4	0.1319 (1.0×10^{-4})	0.0520 (3.2×10^{-4})	0.0526 (3.2×10^{-4})	0.1313 (1.0×10^{-4})	0.0518 (3.2×10^{-4})	0.0523 (3.2×10^{-4})
5	-0.1302 (1.7×10^{-4})	-0.0803 (5.3×10^{-4})	-0.1129 (5.4×10^{-4})	-0.1293 (1.7×10^{-4})	-0.0797 (5.4×10^{-4})	-0.1120 (5.4×10^{-4})
6	0.0862 (1.8×10^{-4})	0.0376 (5.7×10^{-4})	0.06979 (5.7×10^{-4})	0.0856 (1.8×10^{-4})	0.0373 (5.7×10^{-4})	0.06974 (5.7×10^{-4})
7	0.0062 (2.0×10^{-4})	-0.0022 (6.1×10^{-4})	-0.0553 (6.1×10^{-4})	0.0061 (2.0×10^{-4})	-0.0022 (6.1×10^{-4})	-0.0547 (6.1×10^{-4})
8	-0.0259 (2.4×10^{-4})	-0.0726 (7.4×10^{-4})	-0.0144 (7.3×10^{-4})	-0.0257 (2.3×10^{-4})	-0.0719 (7.3×10^{-4})	-0.0143 (7.3×10^{-4})
9	0.0193 (2.4×10^{-4})	0.0844 (7.6×10^{-4})	0.0153 (7.6×10^{-4})	0.0191 (2.4×10^{-4})	0.0835 (7.5×10^{-4})	0.0151 (7.5×10^{-4})
10	0.4158 (2.6×10^{-4})	0.3892 (8.2×10^{-4})	0.3781 (8.2×10^{-4})	0.4116 (2.6×10^{-4})	0.3848 (8.2×10^{-4})	0.3737 (8.2×10^{-4})
11	-0.0158 (2.9×10^{-4})	0.0231 (9.0×10^{-4})	0.0049 (9.0×10^{-4})	-0.0156 (2.9×10^{-4})	0.0229 (9.0×10^{-4})	0.0048 (8.9×10^{-4})
12	0.0235 (3.6×10^{-4})	-0.0669 (11.3×10^{-4})	-0.0181 (11.3×10^{-4})	0.0232 (3.6×10^{-4})	-0.0658 (11.2×10^{-4})	-0.0180 (11.2×10^{-4})
13	-0.0964 (6.2×10^{-4})	-0.3918 (19.5×10^{-4})	-0.6105 (19.4×10^{-4})	-0.0942 (6.1×10^{-4})	-0.3816 (19.1×10^{-4})	-0.5941 (13.1×10^{-4})
14	-0.6814 (10.5×10^{-4})	-0.8456 (33.4×10^{-4})	-0.9923 (33.4×10^{-4})	-0.6545 (10.3×10^{-4})	-0.8091 (32.2×10^{-4})	-0.9473 (32.1×10^{-4})
15	-0.2763 (12.7×10^{-4})	0.0636 (40.1×10^{-4})	0.4214 (40.4×10^{-4})	-0.2632 (12.4×10^{-4})	0.0602 (38.7×10^{-4})	0.3984 (38.6×10^{-4})
16	-0.1091 (285.9×10^{-4})	-0.4009 (1030.0×10^{-4})	-0.3733 (1058.0×10^{-4})	-0.0519 (137.0×10^{-4})	-0.1653 (427.8×10^{-4})	-0.1485 (427.4×10^{-4})