

Itseopiskelumateriaalia: Kausaalimallintamisen perusteet tilastotieteessä

Juho Kopra `juho.j.kopra@jyu.fi`
Santtu Tikka `santtu.tikka@jyu.fi`

Jyväskylän yliopisto,
Matematiikan ja tilastotieteen laitos

19. toukokuuta 2016

Tämä moniste on tarkoitettu itseopiskelumateriaaliksi tilastotieteen maisterivaiheen opiskelijoille (tai vastaavat tiedot omaaville). Erityisesti todennäköisyyslaskennan ja yleistettyjen lineaaristen mallien tuntemus on tarpeen. Materiaalin tarkoituksena on selvittää lukijalle perusteet Judea Pearl'n kehittämästä kausaalimallintamisesta ja -laskennasta. Materiaali perustuu Judea Pearl'n kirjaan *Causality* [Pearl, 2009]. Lauseiden ja määritelmien kohdalla annetaan aina kirjan osio, josta nämä löytyvät.

1 Perusteet

1.1 Graafiteorian perusteita

Graafiteoriaa voidaan hyödyntää osana kausaalipäätelyä, koska graafien avulla on mahdollista havainnollistaa muuttujien välisiä riippuvuuksia tehokkaasti. Graafien avulla voidaan esittää monimuotoisia tilastollisia malleja sekä johtaa muuttujien välisiä riippumattomuusominaisuuksia tehdyistä oletuksista. Graafit tekevät helpoksi esimerkiksi tutkimusten koeasetelmien esittämisen ja visualisoinnin täsmällisesti. Graafin avulla voidaan esittää tutkimuksen kaikki vaiheet ja siihen liittyvien muuttujien väliset yhteydet. Esimerkiksi tutkimuksen otantamenetelmä ja mahdolliset puuttuvan tiedon mekanismit voidaan tuoda helposti esille.

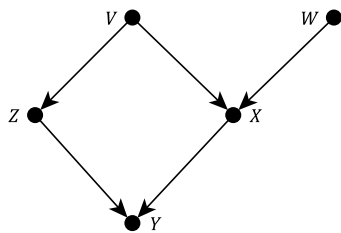
Tässä monisteessa kausaalisuus liittyy läheisesti interventioon, eli ulkopuolisen toiminnon käsitteeseen. Graafien näkökulmasta ulkopuolisella toiminnolla tarkoitetaan tutkijan aiheuttamaa muutosta graafissa, minkä jälkeen on mahdollista tarkastella, kuinka muuttujat reagoivat kyseiseen interventioon. Mikäli graafi kuvaa mallinnettavaa tilannetta oikein, voidaan tällaisten interventioiden vaikutus usein estimoida aineiston perusteella, tai saada tietoa siitä mitä lisäoletuksia tarvittaisiin, jotta päätelmien pätevyys ei olisi kyseenalainen.

Graafi G on pari $\langle V, E \rangle$, missä joukko V sisältää graafin solmut ja joukko E solmuja yhdistävät särmät. Solmujen ajatellaan vastaavan joitakin kiinnostuksen kohteena olevia muuttujia ja särmien näiden välisiä yhteyksiä. Solmut piirretään joko suljettuina ympyröinä (\bullet) tai avoimina ympyröinä (\circ) riippuen siitä onko muuttuja havaittu vai ei. Särmät, eli kahdesta solmusta koostuvat joukot $\{X, Y\}$ puolestaan piirretään solmujen X ja Y välisenä viivana. Kun muuttujien väliset yhteydet mielletään kausaalisiksi, käsitellään graafin särkeä suunnattuina. Suunnattu särmä solmusta X solmuun Y on pari (X, Y) ja se piirretään näiden välisenä nuolena siten, että nuolen kärki osoittaa solmuun Y ($X \bullet \longrightarrow \bullet Y$). Suunnattuja särkeä voidaan merkitä tekstissä nuoli-operaattorilla (\rightarrow tai \leftarrow), esim. $X \rightarrow Y$. Graafi on joko suunnattu tai suuntaamaton eli joko sen kaikki särmät ovat suunnattuja tai suuntaamattomia. Jatkossa käsiteltävät graafit eivät sisällä särkeä (X, X) eli särkeä joiden lähtösolmu ja päätesolmu ovat samat. Lisäksi solmujoukko V ja särkeajoukko E oletetaan äärellisiksi.

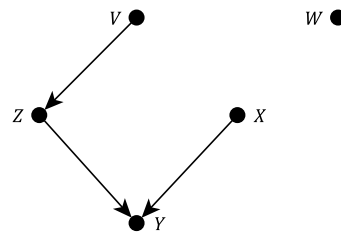
Polulla tarkoitetaan jonoa särkeä $(\{X_1, X_2\}, \{X_2, X_3\}, \dots, \{X_{n-1}, X_n\})$, jossa jokainen särmä on yhteydessä seuraavaan toisen päätesolmunsu kautta, ja jokainen solmu toistuu vain kerran. Poikkeuksena sallitaan polku, joka päättyy siihen solmuun mistä se alkoi eli $X_1 = X_n$. Tällaista polkua sanotaan silmukaksi. Polut ovat usein kiinnostuksen kohteena myös suunnatuissa graafeissa. Esimerkiksi jono suunnattuja särkeä $((X, Z), (Y, Z), (Y, W)) = X \rightarrow Z \leftarrow Y \rightarrow W$ voidaan mieltää polkuna kun särmien suunta jätetään huomiotta. Tätä vastaava polku olisi siis $(\{X, Z\}, \{Y, Z\}, \{Y, W\})$, jossa suunnatut parit on nyt korvattu kahden solmun joukoilla. Suunnattujen graafien yhteydessä voidaan puhua myös suunnatuista poluista, jolloin jonon jokaisen särkeän päätesolmun on oltava sama kuin seuraavan särkeän lähtösolmu. Edellisen esimerkin jono ei siis ole suunnattu polku, vaikka se voidaankin mieltää poluksi. Vastaavasti voidaan myös määrittellä suunnattu silmukka. Jatkossa käsitellään suunnattuja silmukattomia graafeja (engl. directed acyclic graph, DAG) eli suunnattuja graafeja, jotka eivät sisällä lainkaan suunnattuja silmukoita. Tällaiset graafit ovat tärkeässä erityisasemassa kausaalipäätelyssä.

Suunnatut särmät ja polut määrittelevät graafin solmuille lukuisia suhteita

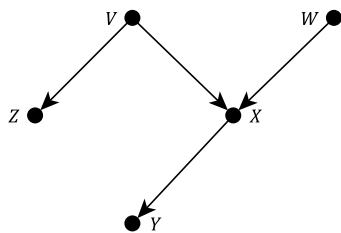
ta, joista esitellään tässä tärkeimmät. Solmu Y on solmun X jälkeläinen jos graafi sisältää suunnatun polun solmusta X solmuun Y (esim. $X \rightarrow Z \rightarrow Y$). Tällöin sanotaan myös, että X on solmun Y esivanhempi. Erityistapaus edellisistä on tilanne, jossa polku solmujen välillä koostuu vain lähtö- ja päätesolmusta: solmu Y on solmun X lapsi jos graafi sisältää suunnatun särmän solmusta X solmuun Y . Tällöin sanotaan myös, että X on solmun Y vanhempi.



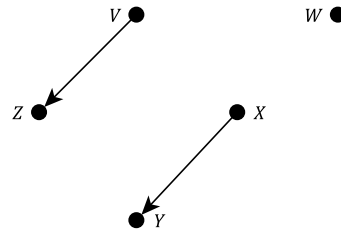
(a) Graafi G .



(b) Graafi $G_{\bar{X}}$.



(c) Graafi $G_{\underline{Z}}$.

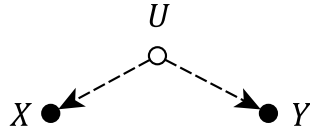


(d) Graafi $G_{\bar{X}\underline{Z}}$.

Kuva 1: Havainnollistetaan särmäpoistoja graafissa G .

Graafiin G voidaan myös kohdistaa erilaisia operaatioita. Merkinnällä $G_{\underline{X}}$ tarkoitetaan graafia, joka saadaan graafista G poistamalla kaikki solmusta X lähtevät särmät. Merkinnällä $G_{\bar{X}}$ tarkoitetaan graafia, joka saadaan graafista G poistamalla kaikki solmuun X saapuvat särmät. Näitä merkintöjä voidaan luonnollisesti yhdistellä, esimerkiksi $G_{\bar{X}\underline{Z}}$ tarkoittaa graafia, joka saadaan graafista G poistamalla solmusta X lähtevät särmät ja solmuun Z saapuvat särmät. Merkinnät voidaan myös yleistää tilanteeseen, jossa X tai Z on solmujoukko. Kaikki edellä mainitut operaatiot tuottavat graafin G aligraafin eli graafin $G' = (V', E')$, missä $V' \subseteq V$ ja $E' \subseteq E$.

Kausaalisuuden alalla on myös otettu käyttöön muutamia graafeihin liittyviä lyhennysmerkintöjä ja konventioita. Tarkastellaan tilannetta, jossa yksi havaitsematon muuttuja on kahden havaitun muuttujan vanhempi kuten kuvassa 2a.



(a) Graafi, joka sisältää yhden havaitsemattoman muuttujan U , sekä havaitut muuttujat X ja Y .



(b) Esimerkki lyhennysmerkinnästä, jossa havaitsemattoman muuttujan vaikutus esitetään kaksisuuntaisella särmällä.

Kuva 2: Muuttuja U voidaan merkitä graafiin kahdella tavalla.

Tällaisessa tilanteessa graafiin yleensä piirretään muuttujia X ja Y yhdistävä kaksisuuntainen särmä kuten kuvassa 2b. On syytä huomata, että kyseessä on vain yleistynyt merkintätapa, jossa tällaiset havaitsemattomat muuttujat jätetään piirtämättä graafiin, minkä tarkoituksena on selkeyttää graafin rakennetta. Tästä on hyötyä erityisesti silloin, kun havaitsemattomia tekijöitä on useita. Kaksisuuntainen särmä ei siis tarkoita kahta yksisuuntaista särmää, vaan merkinnän takana on aina yksi havaitsematon sekoittava tekijä U .

Tiettyjen oletusten vallitessa voidaan solmuja vastaavien muuttujien ehdolliset riippumattomuusominaisuudet todeta suoraan graafista. Tämä tapahtuu d -separoituvuuden avulla.

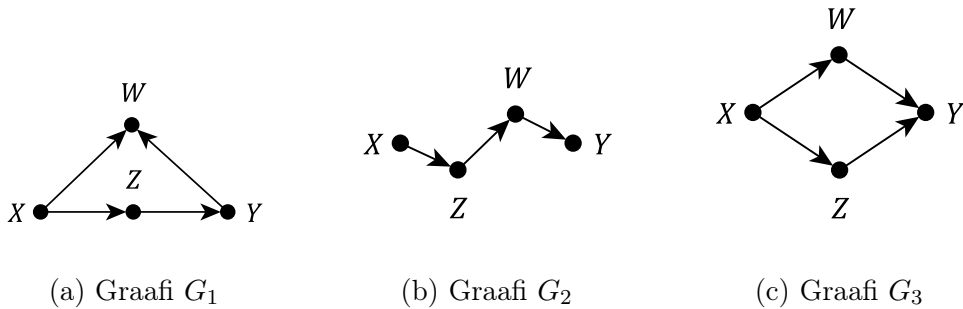
Määritelmä 1 (d -separoituvuus, Causality Def. 1.2.3). *Suuntaamaton polku p suunnatussa graafissa on solmujoukon Z d -separoima täsmälleen silloin, kun molemmat seuraavista ehdoista toteutuvat*

1. Jos p sisältää ketjun $i \rightarrow m \rightarrow j$ tai haarukan $i \leftarrow m \rightarrow j$ niin m kuuluu solmujoukkoon Z .
2. Jos p sisältää yhden tai useamman käänteisen haarukan (collider) $i \rightarrow m \leftarrow j$, niin ainakin yhden haarukan solmu m ei kuulu joukkoon Z , eikä mikään solmun m jälkeläinen kuulu joukkoon Z .

Erilliset muuttujajoukot X ja Y ovat solmujoukon Z d -separoimia, jos ja vain jos Z d -separoi kaikki suuntaamattomat polut muuttujajoukkojen X ja Y välillä.

Edellisessä määritelmässä oleellista on se, että kahden muuttujan välillä voi kulkea useita eri polkuja. Tällöin, jotta muuttujat olisivat d -separoituneita, on kaikkien polkujen d -separoituneisuus tarkistettava. Polkuja käsitellään suuntaamattomina, mutta särmien suunnalla on merkitystä sen kannalta toteutuuko d -separoituvuus vai ei. Tarkastellaan tätä esimerkin avulla.

Esimerkki 1. Tarkastellaan graafeja (a) - (c) kuvassa 3. Halutaan selvittää, ovatko solmut X ja Y solmun Z d -separoimia kyseisissä graafeissa. Kuvan 3a graafissa G_1 muuttujat X ja Y ovat solmun Z d -separoimia, sillä Z on keskosolmu ketjussa $X \rightarrow Z \rightarrow Y$, ja W toimii käänteisenä haarukkana $X \rightarrow W \leftarrow Y$. Jos taas haluttaisiinkin tutkia, ovatko X ja Y solmujoukon $\{W, Z\}$ d -separoimia olisi vastaus kielteinen, sillä nyt W kuuluukin ehtojoukkoon $\{W, Z\}$, joten polku $(\{X, W\}, \{Y, W\})$ ei d -separoidu.



Kuva 3: Esimerkkigraafeja d -separoituvuudesta

Kuvan 3b graafissa G_2 on solmujen X ja Y välillä vain yksi polku, jonka solmu Z d -separoi. Tässä tilanteessa ehtojoukko voisi myös koostua solmuista W ja Z tai pelkästään solmusta W . Kuvan 3c graafissa G_3 on solmujen X ja Y välillä puolestaan kaksi polkua, joista vain toinen sisältää solmun Z . Nyt Z ei siis d -separoi solmuja X ja Y , vaan tähän tarvitaan myös solmu W .

Pelkkä d -separoituvuus ei vielä riitä siihen, että graafin avulla voitaisiin tehdä päätelmiä muuttujien välisistä ehdollisista riippumattomuuksista. Lisäksi muuttujien yhteisjakauman on oltava graafin rakenteen mukainen.

Määritelmä 2 (yhteensopivuus, Causality Def. 1.2.2). Olkoot suunnattu silmukatonta graafi $G = \langle V, E \rangle$ ja muuttujajoukon V yhteisjakauma $P(V)$. Graafi G ja jakauma P ovat yhteensopivia, mikäli

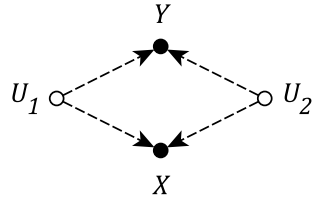
$$P(V_1, \dots, V_n) = \prod_{i=1}^n P(V_i | PA_{V_i}),$$

missä PA_{V_i} on joukko, joka sisältää solmun V_i vanhemmat graafissa G .

Seuraava lause määrittelee ehdollisen riippumattomuuden ja d -separoituvuuden välisen yhteyden

Lause 1 (Causality Thm. 1.2.4). *Jos muuttujajoukko Z d-separoi muuttujat X ja Y graafissa $G = \langle V, E \rangle$, niin muuttujat X ja Y ovat riippumattomia ehdolla Z kaikkien graafin G kanssa yhteensopivien jakaumien $P(V)$ suhteen. Jos X ja Y eivät ole muuttujajoukon Z d-separoimia, niin X ja Z ovat riippuvia ehdolla Z ainakin yhden graafin G kanssa yhteensopivan jakauman P suhteen.*

Graafin kanssa yhteensopivan jakauman ehdolliset riippumattomuudet voidaan siis aina todeta d-separoituvuuden avulla. Käänteiseen suuntaan tämä väite ei kuitenkaan päde. Jos X ja Y eivät ole solmujoukon Z d-separoimia graafissa G , niin X ja Y eivät välttämättä ole ehdollisesti riippuvia kaikkien graafin G kanssa yhteensopivien jakaumien P suhteen. Tällaista tilannetta havainnollistaa kuvan 4 graafi G .



Kuva 4: Graafi G , jossa muuttujat X ja Y eivät ole d-separoituneita, mutta ovat silti riippumattomia.

Oletetaan, että muuttujat U_1 ja U_2 noudattavat kaksiulotteista normaali-jakaumaa $N(\mathbf{0}, I_2)$, missä $\mathbf{0} = (0, 0)$ ja I_2 on identtinen 2×2 matriisi. Olkoot nyt muuttujat X ja Y sellaisia, että niiden arvot määräytyvät seuraavista yhtälöistä

$$X = \beta U_1 + \gamma U_2 \qquad Y = \gamma U_1 - \beta U_2,$$

jolloin saadaan

$$\begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(\beta U_1 + \gamma U_2, \gamma U_1 - \beta U_2) \\ &= \text{Cov}(\beta U_1, \gamma U_1) - \text{Cov}(\beta U_1, \beta U_2) \\ &\quad + \text{Cov}(\gamma U_2, \gamma U_1) - \text{Cov}(\gamma U_2, \beta U_2) \\ &= \beta \gamma \text{Var}(U_1) - 0 + 0 - \gamma \beta \text{Var}(U_2) \\ &= \beta \gamma - \gamma \beta \\ &= 0. \end{aligned}$$

Koska myös X ja Y noudattavat yhdessä kaksiulotteista normaalijakaumaa, seuraa tästä kyseisten muuttujien välinen riippumattomuus. Kyse on siitä,

että sopivasti valituilla mallin parametreilla saadaan aikaan muuttujien välinen riippumattomuus, jota ei voida määrittää mallia vastaavasta graafista käyttäen d-separoituvuutta (tälle on englanninkielinen termi *incidental cancellation*). Jätettäkään harjoitustehtäväksi varmistaa, että muuttujien X ja Y yhteisjakauma todella on normaali ja että $P(X, Y, U_1, U_2)$ on graafin G kanssa yhteensopiva.

Tehtäviä 1.

1. Olkoon graafi G kuten kuvassa 2b. Piirrä graafi $G_{\overline{X}}$.
2. Olkoon graafi $G = \langle V, E \rangle$, missä $V = \{X, Y, Z\}$ ja $E = \{(X, Y), (Z, X)\}$. Piirrä kaikki graafin G aligraafit.
3. Kuinka monta erilaista suunnattua silmukatonta graafia voidaan muodostaa tilanteessa, jossa graafi sisältää kolme solmua?
4. a) Osoita Kuvaan 4 liittyen, että jakauma $P(X, Y)$ on 2-ulotteinen normaalijakauma.
b) Osoita, että $P(X, Y, U_1, U_2)$ on graafin G kanssa yhteensopiva.

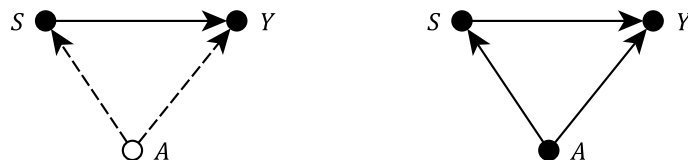
1.2 Kausaalisuudesta

Kokeellisia tutkimuksia lukuun ottamatta on tilastotieteessä perinteisesti välitetty päätelemästä että X on Y :n syy, vaikkakin monessa tilanteessa soveltajat haluaisivatkin tällaisen päätelmän tehdä. Kausaalipäätelmien tekeminen edellyttääkin tiettyjä oletuksia. Tässä monisteessa rajoitumme tilanteisiin, joissa muuttujien väliset suhteet eli kausaalimallin rakenne (graafi) tunnetaan. Todellisuudessa näin ei aina ole, vaan rakenteeksi voidaan valita useita vaihtoehtoisia rakenteita ja tutkia niitä yksi kerrallaan.

Aiemmin on ajateltu, että kausaalisuus voidaan todentaa vain tekemällä kokeita, joissa muuttujan X (joka potentiaalisesti vaikuttaa Y :hyn) arvo on satunnaistettu. Kokeissa saadaan selville X :n kausaalivaikutus Y :hyn, sillä satunnaistamisen takia X :n arvoon on vaikuttanut vain kokeen tekijä.

Tyypillisesti kokeen toteuttaminen on kallista, ja on myös paljon tilanteita joissa kokeen toteuttaminen ei ole joko eettisesti sallittua, tai lainkaan mahdollista. Näistä rajoituksista huolimatta kausaalipäätelmien tekeminen on joissain tilanteissa mahdollista myös silloin, kun X :n arvoon ei olla tutkimuksen aikana vaikutettu, vaan se on ainoastaan mitattu. Tällaista tutkimusta sanotaan havainnoivaksi tutkimukseksi (engl. observational study).

Kausaalisuudella ajatellaan tässä monisteessa ensisijaisesti tilannetta:



Kuva 5: Vasemmalla ennustemallin graafi G_1 , jossa muuttuja A ei ole havaittu. Oikealla graafi G_2 , jossa muuttuja A on havaittu ja jolloin kausaalipäätely onnistuu.

Jos muuttujalle X asetetaan jokin arvo, niin mikä on muuttujan Y jakauma?

Vertaa kysymystä ennustamiseen todennäköisyyslaskennassa:

Jos havaitsemme muuttujan X arvon, niin mikä on muuttujan Y jakauma?

Lisäksi kausaalimallintaminen antaa meille mahdollisuuden vastata kysymyksen:

Jos olisimme toimineet toisin ja asettaneet muuttujan X arvoksi jonkin sen havaitusta arvosta poikkeavan arvon, niin mikä olisi ollut muuttujan Y jakauma?

Ensinmainittuun (kausaali)kysymykseen liittyen käytämme merkintää $P(Y|do(X = x))$, kun taas toiseen, ennustamiseen liittyvään kysymykseen liittyy todennäköisyyslaskennasta tuttu merkintä $P(Y|X = x)$, joka on ehdollinen todennäköisyys. Kolmas kysymys on niin kutsuttu kontrafaktuaaliksi kysymykseksi, jota käsittelemme myöhemmin. Kysymysten 1 ja 2 välinen ero ei ole välttämättä ilmeinen, joten tarkastellaanpa esimerkkiä.

Esimerkki 2. Tarkastellaan aineistoa, jossa on mitattu äitien raskausajan tupakointia (S), äitien ikää (A), sekä vasteena (Y) onko syntyneellä lapsella Downin oireyhtymä. Mikäli vain tupakointi ja vaste ovat mitattuja, niin ennustaminen on mahdollista:

$$P(Y|S).$$

Sen sijaan huomioimalla kausaalirakenne, eli tekemämme oletus siitä, että ikä A vaikuttaa sekä tupakointiin että vasteeseen Y , voidaan selvittää onko tehdyillä mittauksilla mahdollista tehdä kausaalipäätelyä. Toisin sanoen,

voimmeko vastata kysymyksiin: "Mikä on todennäköisyys Downin oireyhtymälle jos kaikki äidit tupakoisivat raskauden aikana?" ja "Mikä on todennäköisyys Downin oireyhtymälle jos yksikään äiti ei tupakoisi raskauden aikana?". Edellisiä kysymyksiä voidaan tarkastella myös ikäryhmäkohtaisesti. Ilmenee, että mikäli ikää ei ole mitattu, niin em. kysymyksiin ei ole graafin G_1 (Kuva 5) tapauksessa mahdollista vastata.

Taulukko 1: Esimerkin 2 aineisto ilman äitien ikää (A).

	Y=0	Y=1
S=0	424	285
S=1	115	176

Ennustemallia käyttämällä saadaan $P(Y = 1|S = 0) = 0.402$ ja $P(Y = 1|S = 1) = 0.605$. Ennustemallin perusteella raskausajan tupakointi näyttää olevan yhteydessä lapsen Downin oireyhtymään. Mikäli äidin ikä A on tiedossa, voimme käyttää hyväksi graafia. Näin ehdollistamalla myös muuttujalla A saadaan:

$$\begin{aligned}
 P(Y = 1|S = 0, A < 35) &= 0.198 \\
 P(Y = 1|S = 1, A < 35) &= 0.223 \\
 P(Y = 1|S = 0, A \geq 35) &= 0.712 \\
 P(Y = 1|S = 1, A \geq 35) &= 0.742
 \end{aligned}$$

Huomataan, että raskausajan tupakoinnilla näyttäisi olevan pienempi rooli Downin oireyhtymän synnyssä, kuin äidin iällä.

Taulukko 2: Esimerkin 2 aineisto kokonaisuudessaan.

	A < 35		A ≥ 35	
	Y=0	Y=1	Y=0	Y=1
S=0	344	84	80	201
S=1	59	18	56	158

Esimerkissä 2 ikä (A) on sekoittavan tekijän (engl. confounder) roolissa. Mikäli sekoittavaa tekijää ei ole havaittu tai mitattu, se voi aiheuttaa näennäisen assosiaation kahden muun muuttujan välille. Todellisuudessa Downin

oireyhtymän ja tupakoinnin välillä voi olla muitakin sekoittavia tekijöitä, jotka on jätetty huomiotta tässä mallissa.

Ero ennustamisen ja kausaalipäätelyn välillä voidaan edellisessä esimerkissä käsittää tutkijan roolin kautta. Mikäli halutaan vain ennustaa Downin oireyhtymän ilmenemistä, niin riittää että tutkija toimii vain ulkopuolisena tarkkailijana tutkimuksen ajan vaikuttamatta itse äitien tupakointiin. On selvää, että intervention tekeminen esimerkin tilanteessa ei ole mahdollista. Tutkija ei voi pakottaa tutkimukseen osallistuneita äitejä aloittamaan tai lopettamaan tupakointia. Esimerkin tilanteessa kausaalipäätely on tästä huolimatta mahdollista, mutta tarvitsemme seuraavassa luvussa esiteltävää kausaalimallia, jotta voimme vastata kysymykseen tupakoinnin kausaalisesta yhteydestä Downin oireyhtymään. Huomautettakoon että esimerkin aineisto on simuloitua, eivätkä siitä lasketut todennäköisyydet kuvasta todellista riskiä saada Downin oireyhtymä.

Tehtäviä 2.

- Toista Esimerkki 2:n analyysi esim. R-ohjelmistolla käyttäen logistista regressiota. Onko raskausajan tupakoinnin regressiokerroin tilastollisesti merkitsevä?*
 - kun ikä ei ole mallissa?*
 - kun ikä on mallissa?*
- Esimerkin aineisto on simuloitua, mutta sovellus todellinen. Lue artikkeli Chen et al. [1999] ja laske todellisella aineistolla ehdolliset jakaumat $P(Y|S)$ ja $P(Y|S, A)$ sekä toista Tehtävä 1 tällä aineistolla. Artikkelissa käytetään tapaus-verrokki -asetelmaa, mutta sitä ei tässä tehtävässä tarvitse ottaa huomioon.*

1.3 Graafi ja kausaalisuus yhdessä

Tässä kappaleessa määrittelemme kausaalimallin ja käymme lävitse kausaalilaskentaa. Kausaalimallinnuksen pyrkimyksenä on laskea vasteen Y jakauma, kun pakotamme muuttujan X arvoon x , vaikka kaikkia mallirakenteen muuttujia ei olisi mitattu.

Määritelmä 3 (Kausaalimalli, Causality Def. 7.1.1 & Def. 7.1.6). *Kausaalimalli on nelikko $M = \langle U, V, F, P \rangle$, jossa*

- $U := \{U_1, \dots, U_m\}$ on taustamuuttujien joukko, jonka arvo määräytyy mallin ulkopuolelta. Muuttujia U sanotaan eksogeenisiksi muuttujiksi.*

(ii) $V := \{V_1, V_2, \dots, V_n\}$ on joukko muuttujia, jotka määräytyvät mallissa olevien muuttujien $U \cup V$ perusteella. Muuttujia V sanotaan endogeenisiksi muuttujiksi.

(iii) F on joukko f_1, f_2, \dots, f_n funktioita, siten että kukin f_{V_i} on kuvaus $f_{V_i} : U_{V_i} \cup PA_{V_i} \rightarrow \{V_i\}$, missä $U_{V_i} \subseteq U$ on muuttujaan V_i liittyvät endogeeniset muuttujat (joita voi olla 0-n kpl) ja $PA_{V_i} \subseteq V \setminus \{V_i\}$ ovat ne havaitut muuttujat, jotka ovat V_i :n vanhempia. Kukin f_{V_i} määrittää muuttujan V_i :n arvon

$$v_i = f_{V_i}(pa_{V_i}, u_{V_i}), \quad i = 1, \dots, n, \quad (1)$$

joka riippuu ko. muuttujan vanhempien PA_{V_i} arvosta pa_{V_i} sekä siihen liittyvien eksogeenisten muuttujien U_{V_i} arvosta u_{V_i} . Lisäksi oletetaan, että mikäli funktiot F ovat tunnettuja, niin on mahdollista ratkaista yksikäsitteisesti muuttujien V arvot. Toisin sanoen, on olemassa yksikäsitteinen kokoelma kuvauksia $H(U) := \{h_{V_1}, \dots, h_{V_n}\}$ siten, että mille tahansa $u \in U$ on voimassa

$$v_i = h_{V_i}(u).$$

(iv) $P(u)$ on todennäköisyysjakauma, joka on määritelty satunnaismuuttujien U otosavaruudessa.

Tärkeää kausaalimallissa on tässä kohtaa ymmärtää, että muuttujan X_i arvon $x_i = f_{X_i}(pa_{X_i}, u_{X_i})$ määrittävä funktio f ajatellaan deterministiseksi, vaikka itse funktio onkin tuntematon. Kaikki muuttujan X_i havaittava vaihtelu johtuu siis sekä sen havaittavien vanhempien pa_{X_i} vaihtelusta että havaitsemattomien tekijöiden u_{X_i} :n satunnaisvaihtelusta. Voidaankin sanoa, että kausaalimalli käsittää muuttujien väliset suhteet deterministisinä, mutta jonka havaitsemattomien tekijöiden U ajatellaan olevan satunnaisia ja toisistaan riippumattomia.

Kausaalimallin määritelmässä ei mainita graafia G , joka liittyy kuitenkin oleellisesti malliin. Graafi on keino kuvata mallia siten, että graafi ja mallin jakauma ovat yhtäsoivia (Määritelmä 2), eikä se siksi ole mukana määritelmässä. Graafin suunnatut särmät muodostuvat yhtälön (1) perusteella siten, että kustakin funktion f_{V_i} argumentista (eli kaikista muuttujista PA_{V_i} ja U_{V_i} , jotka voivat olla vektorimuotoisia) lähtee suunnattu särmä muuttujaan V_i . Huomautettakoon myös, että vaatimus kuvauksien $H(U)$ kokoelmalle ehdossa (iii) liittyy alimallin kausaalimalliuteen ja kontrafaktuaalisuuteen, joita käsittelemme myöhemmin. Funktioiden kokoelman $H(U)$ yksikäsitteisyys ei siis ole tässä kohtaa vielä keskeistä. Kausaalimallin rakennetta muokkaamalla syntyy alimalli.

Määritelmä 4 (Alimalli, Causality Def. 7.1.2). *Olkoon M kausaalimalli, X joukko muuttujia V :ssä, ja x on reaalisatio X :stä. Kausaalimallin M alimalli M_x on kausaalimalli*

$$M_x = \langle U, V, F_x, P \rangle,$$

jossa $F_x = \{f_{V_i} : V_i \notin X\} \cup \{X = x\}$.

Alimallin funktioiden joukko F_x muodostetaan siten, että muuttujajoukon X muuttujien V_i arvon määräävät funktiot f_{V_i} korvataan vakiofunktioilla siten että yhtälö $X = x$ toteutuu. Funktiot, jotka eivät liity muuttujajoukkoon X pysyvät samoina. Tällaista funktioiden korvaamista kutsutaan myös interventioksi. Esimerkiksi alimallissa $M_{z'}$, jossa joukkoon X kuuluu vain yksi muuttuja Z , korvataan yhtälön (1) alkuperäinen funktio $z = f_Z(p_Z, u_Z)$ vakiofunktioilla $z = z'$. Koska interventio voidaan kohdistaa vain havaittuihin muuttujiin, niin havaitsemattomien muuttujien todennäköisyysjakauma P on sama, kuin aiemminkin. Alimalli todellakin on kausaalimalli, vaikka osa funktioista F korvataan vakiofunktioilla. Määritelmän mukaisella joukolle F_x on nyt yksikäsitteinen ratkaisu $H_x(U)$, kuten kausaalimallin Määritelmä 3 edellyttää. Todistus sivuutetaan.

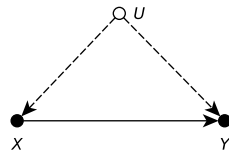
Määritelmä 5 (Kausaalivaikutus, Causality Def. 7.1.3). *Olkoon M kausaalimalli, X muuttujajoukko havaittuja muuttujia V :ssä ja olkoon x jokin X :n arvo. Muuttujan X kausaalivaikutus Y :hyn, merkitään $P(Y|do(X = x))$, mallissa M saadaan alimallin M_x avulla.*

Merkintä $P(Y|do(X = x))$ tarkoittaa Y -muuttujan jakaumaa alimallissa M_x , kun $X = x$, eikä sitä tule sekoittaa ehdolliseen jakaumaan $P(Y|X = x)$ alkuperäisessä mallissa M .

Määritelmä 6 (Kausaalivaikutuksen identifioitavuus, Causality: Definition 3.2.4). *Muuttujan X kausaalivaikutus Y :hyn on identifioituva graafista G , jos mille tahansa kahdelle mallille M_1 ja M_2 , joilla on keskenään samat graafit $G(M_1) = G(M_2)$ (mutta eri funktiot) pätee $P_{M_1}(y|do(X = x)) = P_{M_2}(y|do(X = x))$, kun $P_{M_1}(v) = P_{M_2}(v) > 0$.*

Tämä määritelmä tarkoittaa siis sitä, että identifioituva kausaalivaikutus $P(y|do(X = x))$ on yksikäsitteisesti laskettavissa graafin G kanssa yhteensopivasta havaittujen muuttujien yhteisjakaumasta olivatpa kausaalimallin tuntemattomat funktiot mitä muotoa tahansa. Määritelmää voidaan käyttää todistamaan, että jokin algoritmi, tai laskusäännöt, antavat yksikäsitteisen tuloksen. Sovellettaessa kausaalilaskentaa, ei tätä määritelmää yleensä tarvitse hyödyntää.

Shpitser and Pearl [2006] huomasivat että Määritelmässä 6 riittää olettaa $P_{M_1}(X = x|PA_X) > 0$, josta seuraa että $P_{M_2}(X = x|PA_X) > 0$. Vaadittu ehto $P_{M_1}(X = x|PA_X) > 0$, tarkoittaa, että voimme laskea kausaalivaikutuksen jakauman $P(y|do(X = x))$ vain sellaisille x , joilla on positiivinen (ei nolla) todennäköisyys. Todellisen aineiston kanssa sovellettaessa tämä tarkoittaa sitä, että meillä on oltava havaintoja $X = x$, sillä muutoin frekventistinen estimaattori antaisi $P(X = x) = 0$. Emme siis voi päätellä intervention $X = x$ vaikutusta Y :hyn, ellei ole havaittu $X = x$.



Kuva 6: Graafi esimerkkiin 3.

Esimerkki 3. Osoitetaan Määritelmän 6 avulla, että kausaalivaikutus $P(Y|do(X))$ ei ole identifioituvaa graafissa Kuvassa 6. Olkoot mallit M_1 ja M_2 , joille havaitut jakaumat ovat samat: $P_{M_1}(X, Y) = P_{M_2}(X, Y)$. Samoin molemmille malleille muuttujat X, Y ja U ovat binäärisiä, muuttuja U saa arvon 1 todennäköisyydellä 0.5 (samoin arvon 0 todennäköisyydellä 0.5), ja muuttujan X arvo määräytyy funktiosta $f_x(u) = u$. Olkoon \oplus ”joko tai” -operaattori, joka palauttaa luvun 1 kun täsmälleen yksi argumenteista saa arvon 1 ja luvun 0 muutoin. Muuttujan Y arvo määräytyy funktion $f_y(x, u) = x \oplus u$ avulla mallissa M_1 ja $f_y(x, u) = 0$ mallissa M_2 . Tällöin $P_{M_1}(X = 0) = P_{M_2}(X = 0) = 0.5 > 0$ ja $P_{M_1}(X = 1) = P_{M_2}(X = 1) = 0.5 > 0$, joten voidaan laskea kausaalivaikutuksia sekä interventiolle $do(X = 1)$ että $do(X = 0)$. Määritellään jakaumien laskemista varten indikaattorifunktio $\mathbb{1}(t)$, joka palauttaa arvon 1, kun argumentti t on tosi, ja 0 kun argumentti on epätosi. Nyt havaitut jakaumat $P_{M_1}(X, Y) = P_{M_2}(X, Y)$ ovat samat, sillä mallille M_1 pätee

$$\begin{aligned}
 P_{M_1}(X = 0, Y = 0) &= P(f_x(u) = 0, f_y(x, u) = 0) \\
 &= \sum_u \mathbb{1}(u = 0) \mathbb{1}(x \oplus u = 0) P(U = u) \\
 &= 1 \cdot 1 \cdot P(U = 0) + 0 \cdot 1 \cdot P(U = 1) = P(U = 0) = 0.5 \\
 P_{M_1}(X = 1, Y = 0) &= P(f_x(u) = 1, f_y(x, u) = 0) \\
 &= \sum_u \mathbb{1}(f_x(u) = 1) \mathbb{1}(x \oplus u = 0) P(U = u) \\
 &= 0 \cdot 1 \cdot P(U = 0) + 1 \cdot 1 \cdot P(U = 1) = P(U = 1) = 0.5
 \end{aligned}$$

$$\begin{aligned}
P_{M_1}(X = 0, Y = 1) &= P(f_x(u) = 0, f_y(x, u) = 1) \\
&= \sum_u \mathbb{1}(f_x(u) = 0) \mathbb{1}(x \oplus u = 1) P(U = u) \\
&= 1 \cdot 0 \cdot P(U = 0) + 0 \cdot 0 \cdot P(U = 1) = 0 \\
P_{M_1}(X = 1, Y = 1) &= P(f_x(u) = 1, f_y(x, u) = 1) \\
&= \sum_u \mathbb{1}(f_x(u) = 1) \mathbb{1}(x \oplus u = 1) P(U = u) \\
&= 0 \cdot 0 \cdot P(U = 0) + 1 \cdot 0 \cdot P(U = 1) = 0,
\end{aligned}$$

ja mallille M_2 pätee

$$\begin{aligned}
P_{M_2}(X = 0, Y = 0) &= P(f_x(u) = 0, f_y(x, u) = 0) \\
&= \sum_u \mathbb{1}(u = 0) \mathbb{1}(0 = 0) P(U = u) \\
&= 1 \cdot 1 \cdot P(U = 0) + 0 \cdot 1 \cdot P(U = 1) = P(U = 0) = 0.5 \\
P_{M_2}(X = 1, Y = 0) &= P(f_x(u) = 1, f_y(x, u) = 0) \\
&= \sum_u \mathbb{1}(f_x(u) = 1) \mathbb{1}(0 = 0) P(U = u) \\
&= 0 \cdot 1 \cdot P(U = 0) + 1 \cdot 1 \cdot P(U = 1) = P(U = 1) = 0.5 \\
P_{M_2}(X = 0, Y = 1) &= P(f_x(u) = 0, f_y(x, u) = 1) \\
&= \sum_u \mathbb{1}(f_x(u) = 0) \mathbb{1}(0 = 1) P(U = u) \\
&= 1 \cdot 0 \cdot P(U = 0) + 0 \cdot 0 \cdot P(U = 1) = 0 \\
P_{M_2}(X = 1, Y = 1) &= P(f_x(u) = 1, f_y(x, u) = 1) \\
&= \sum_u \mathbb{1}(f_x(u) = 1) \mathbb{1}(0 = 1) P(U = u) \\
&= 0 \cdot 0 \cdot P(U = 0) + 1 \cdot 0 \cdot P(U = 1) = 0.
\end{aligned}$$

Sen sijaan kausaalivaikutukset eivät ole samat, sillä

$$\begin{aligned}
P_{M_1}(Y = 0 | do(X = 0)) &= \sum_u P(u) \mathbb{1}(f_y(0, u) = 0) = P(u = 1) = 0.5 \\
P_{M_1}(Y = 0 | do(X = 1)) &= P(u = 0) = 0.5 \\
P_{M_2}(Y = 0 | do(X = 0)) &= \sum_u P(u) \mathbb{1}(f_y(0, u) = 0) \\
&= P(u = 0) + P(u = 1) = 1 \\
P_{M_2}(Y = 0 | do(X = 1)) &= P(u = 0) + P(u = 1) = 1,
\end{aligned}$$

Siis kausaalivaikutus ei ole identifioituva, koska $P_{M_1}(Y = 0|do(X = 0)) = 0.5 \neq 1 = P_{M_2}(Y = 0|do(X = 0))$ vaikka havaitut jakaumat olivat samat.

1.3.1 Markov-mallit ja semi-Markov-mallit

Kausaalimalleista puhuttaessa tarkoitetaan yleensä joko Markov-malleja tai semi-Markov-malleja. Markov-malleissa kaikki muuttujat ovat havaittuja (graafissa merkitään \bullet), eli potentiaalisia (havaitsemattomia) sekoittavia tekijöitä ei ole, kun taas semi-Markov-malleissa osa muuttujista on havaitsemattomia (graafissa merkitään \circ). Markov-malleissa kausaalivaikutuksen jakauksen laskeminen on yksinkertaisempaa kuin semi-Markov-malleissa. Käytännön sovellutuksissa, joihin kausaalilaskentaa sovelletaan, on kuitenkin tyyppillisesti kyse semi-Markov-malleista. Semi-Markov malleissa havaitsemattomia sekoittavia tekijöitä ei välttämättä voida ottaa huomioon, jolloin on selvitettävä kausaalivaikutuksen identifioituvuus. Tietyissä tilanteissa kausaalivaikutuksen identifioituvuus voidaan todeta yksinkertaisten kriteerien avulla (Määritelmä 7 ja Lause 2). Lauseessa 3 esitellään kausaalilaskennan säännöt (engl. do-calculus), joita käyttäen voidaan todeta minkä tahansa kausaalivaikutuksen identifioituvuus.

Määritelmä 7 (Ei-sekoittuneisuus, Causality Def. 6.2.1). *Olkoon M kausaalimalli, josta aineisto on generoitunut. Siis formaali kuvaus miten kunkin muuttujan arvo määräytyy. Merkitään $P(y|do(x))$ vasteen $Y = y$ todennäköisyyttä laskien mallista M hypoteettisessa tilanteessa, jossa on pakotettu $X = x$. Tällöin sanomme, että X ja Y eivät ole sekoittuneita M :ssä, täsmälleen silloin kun*

$$P(y|do(x)) = P(y|x) \quad (2)$$

kaikille x ja y , joille $P(y|x)$ on ehdollinen todennäköisyys. Jos yhtälö (2) pätee, niin sanomme $P(y|x)$:n olevan harhaton.

Jos Määritelmän 7 nojalla harhattomuus on voimassa, niin kausaalivaikutus on identifioituva eikä Lausetta 2 tarvita. Lauseen 2 avulla voidaan helposti määrittää identifioituvuus tilanteessa, jossa kaikki X :n vanhemmat ovat havaittuja.

Lause 2 (Takaovikriteeri, Causality Thm. 3.2.5). *Olkoon semi-Markov-kausaalimalli graafista G , jonka muuttujien osajoukko V on mitattu. Tällöin kausaalivaikutus $P(y|do(X = x))$ on identifioituva aina kun $\{X \cup Y \cup PA_X\} \subseteq V$, eli kun X , Y ja kaikki X :n vanhemmat ovat mitattuja. Tällöin lauseke $P(y|do(X = x))$ saadaan ottamalla X :n vanhemmat PA_X huomioon, eli laskemalla*

$$P(y|do(X = x)) = \sum_{pa_x} P(y|x, pa_x)P(pa_x).$$

Todistus. Harjoitustehtävä. □

Esimerkki 4. Lauseen 2 avulla saamme suoraan lausekkeen tupakoinnin kausaalivaikutukselle Downin oireyhtymän syntyyn esimerkin 2 graafissa G_2 ja se on

$$P(Y|do(S = s)) = \sum_a P(y|s, a)P(a).$$

Voimme nyt laskea Downin oireyhtymän todennäköisyyden tilanteessa, jossa kaikki äidit olisivat aloittaneet tupakoinnin. Ikäryhmien osuudet ovat lähes yhtäsuuret aineistossa, joten käytetään approksimaatiota $P(A < 35) \approx 0.5$. Todennäköisyys että lapsella on Downin oireyhtymä, kun äiti tupakoi raskauden aikana on

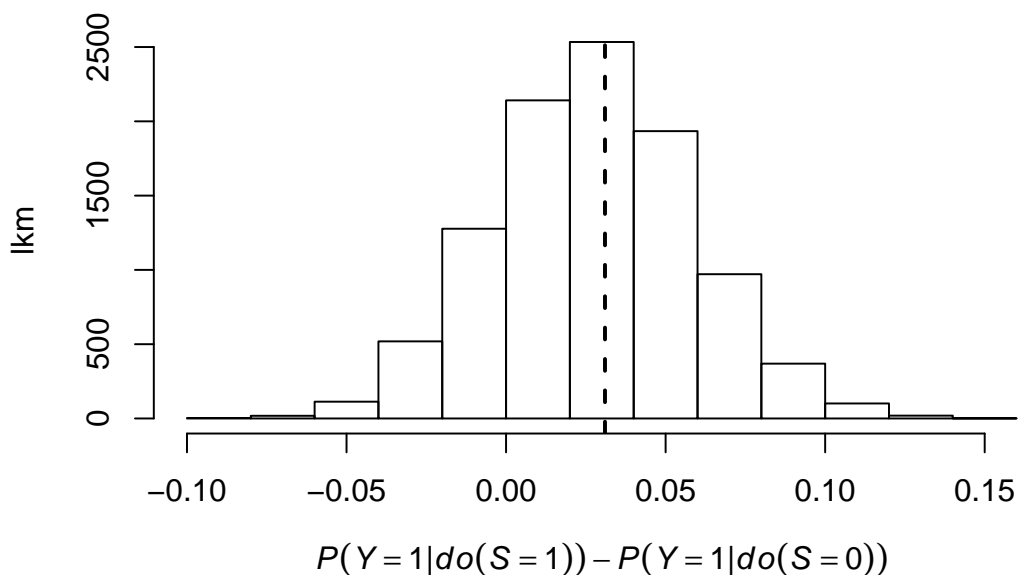
$$\begin{aligned} P(Y = 1|do(S = 1)) &= \\ P(Y = 1|S = 1, A < 35)P(A < 35) + P(Y = 1|S = 1, A \geq 35)P(A \geq 35) &\approx \\ 0.223 \cdot 0.5 + 0.742 \cdot 0.5 &= 0.483, \end{aligned}$$

ja kun äiti ei tupakoi raskauden aikana

$$\begin{aligned} P(Y = 1|do(S = 0)) &= \\ P(Y = 1|S = 0, A < 35)P(A < 35) + P(Y = 1|S = 0, A \geq 35)P(A \geq 35) &\approx \\ 0.198 \cdot 0.5 + 0.712 \cdot 0.5 &= 0.455. \end{aligned}$$

Aineistosta laskettu todennäköisyyksien erotus on $0.483 - 0.455 = 0.028$. Tutkitaan vielä bootstrap-menetelmällä eroavatko nämä todennäköisyydet toisistaan tilastollisesti merkitsevästi. Aineistosta bootstrap-otoksia generoiden lasimme kausaalivaikutusten erotuksen $P(Y = 1|do(S = 1)) - P(Y = 1|do(S = 0))$ bootstrap-jakauman, ks. Kuva 7. Bootstrap-otosten avulla saadaan p -arvo ≈ 0.481 ja 95 % luottamusväli $(-0.033, 0.090)$. Siis kausaalivaikutukset eivät eroa toisistaan tilastollisesti merkitsevästi, joten päätellään että ei ole näyttöä että tupakointi vaikuttaisi Downin oireyhtymän syntyyn.

Kausaalivaikutusten erotus



Kuva 7: Kausaalivaikutusten erotuksen $P(Y = 1|do(S = 1)) - P(Y = 1|do(S = 0))$ bootstrap-jakauma. Havaittu erotus on merkitty katkoviivalla.

Lause 3 (Kausaalilaskenta, Causality Thm. 3.4.1). *Olkoon G suunnattu silmukatonta graafi (DAG) liittyen kausaalimalliin (kuten määritelmässä 3), ja olkoon $P(\cdot)$ mallin todennäköisyysjakauma. Mille tahansa erilliselle osajoukolle X, Y, Z ja W pätevät seuraavat säännöt.*

1. *Havaintojen poistaminen ja lisääminen*

$$P(Y|do(X), Z, W) = P(Y|do(X), W) \text{ jos } Y \perp\!\!\!\perp Z|\{X, W\} \text{ graafissa } G_{\overline{X}}$$

2. *Toiminnan ja havainnon vaihtaminen*

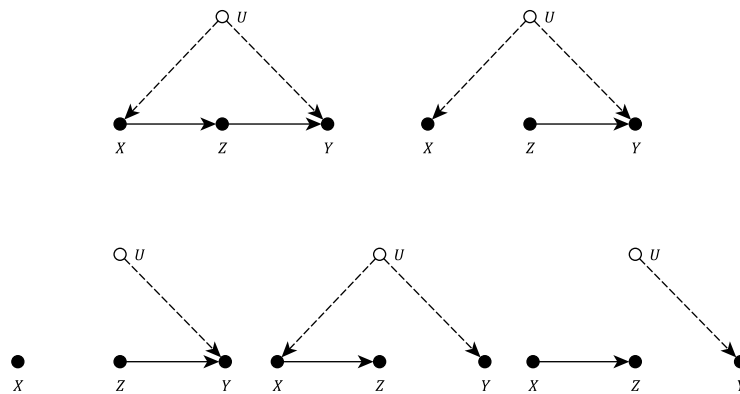
$$P(Y|do(X), do(Z), W) = P(Y|do(X), Z, W) \text{ jos } Y \perp\!\!\!\perp Z|\{X, W\} \text{ graafissa } G_{\overline{XZ}}$$

3. *Toiminnan poistaminen ja lisääminen*

$$P(Y|do(X), do(Z), W) = P(Y|do(X), W) \text{ jos } Y \perp\!\!\!\perp Z|\{X, W\} \text{ graafissa } G_{\overline{XZ(W)}}$$

missä joukko $Z(W)$ muodostuu niistä joukon Z solmuista, jotka eivät ole minkään joukon W solmun esivanhempia graafissa $G_{\underline{X}}$.

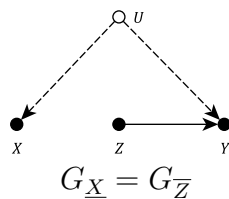
Seuraus 1 (Causality Cor. 3.4.2). *Kausaalivaikutus $q = P(y_1, \dots, y_k | do(X_1 = x_1), \dots, do(X_n = x_n))$ on identifioituva mallissa, jota kuvaa graafi G , jos äärellisellä määrällä iteraatioita kausaalilaskennan sääntöjä käyttäen voidaan kirjoittaa q muodossa, jossa todennäköisyysjakauman ehdossa ei ole yhtään do-operaattoria.*



Kuva 8: Graafi G , jota käytetään laskuissa 1-3.

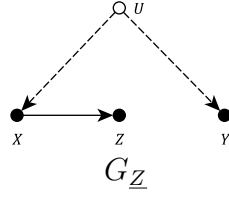
Seuraavaksi havainnollistetaan kausaalilaskennan sääntöjä esimerkkien avulla. Esimerkit liittyvät graafiin G Kuvassa 8.

Lasku 1: $P(z | do(X = x))$:



Ehto säännölle 2 täyttyy, sillä $Z \perp\!\!\!\perp X | \emptyset$ graafissa $G_{\underline{X}}$. Polku $X \leftarrow U \rightarrow Y \leftarrow Z$ on d--separoitu, sillä se sisältää käänteisen haarukan $U \rightarrow Y \leftarrow Z$ ja Y ei ole ehdossa. Saadaan

$$P(z | do(X = x)) = P(z | x). \quad (3)$$



Lasku 2: $P(y|do(Z = z))$:

Nyt emme voi käyttää sääntöä 2 muuttaaksemme $do(Z = z)$:n z :ksi, koska $G_{\underline{Z}}$ sisältää "takaovi"-polun $Z \leftarrow X \leftarrow U \rightarrow Y$. Tietysti haluaisimme katkaista tämän polun mittaamalla muuttujia (kuten X), jotka ovat kyseisellä polulla. Tämä johtaa siihen että on ehdollistettava X :llä ja summattava yli sen mahdollisten arvojen:

$$P(y|do(Z = z)) = \sum_x P(y|x, do(Z = z))P(x|do(Z = z)). \quad (4)$$

Nyt on käsiteltävänä kaksi termiä, joissa on do-operaattori: $P(y|x, do(Z = z))$ ja $P(x|do(Z = z))$. Jälkimmäinen voidaan laskea käyttämällä sääntöä 3:

$$P(x|do(Z = z)) = P(x) \text{ jos } Z \perp\!\!\!\perp X \text{ graafissa } G_{\overline{Z}}, \quad (5)$$

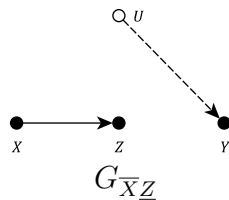
koska X ja Z ovat d-separoituja $G_{\overline{Z}}$:ssa. Tämä on luonnollistakin, sillä Z :n arvon asettaminen ei vaikuta X :ään, koska Z on X :n jälkeläinen G :ssä. Ensimmäinen termi voidaan laskea käyttäen sääntöä 2:

$$P(y|x, do(Z = z)) = P(y|x, z) \text{ jos } Z \perp\!\!\!\perp Y|X \text{ graafissa } G_{\underline{Z}}, \quad (6)$$

koska X d-separoi Z :n ja Y :n graafissa $G_{\underline{Z}}$. Nyt sijoittamalla yhtälöt (5) ja (6) yhtälöön (4) saadaan

$$P(y|do(Z = z)) = \sum_x P(y|x, z)P(x). \quad (7)$$

Lasku 3: $P(y|do(X = x))$:



Kirjoittamalla

$$P(y|do(X = x)) = \sum_z P(y|z, do(X = x))P(z|do(X = x)) \quad (8)$$

huomataan, että termi $P(z|do(X = x))$ laskettiin jo aiemmin (yhtälö (3)), mutta mikään laskusäännöistä ei salli poistaa do-operaattoria termistä $P(y|z, do(X = x))$. Voidaan kuitenkin käyttää sääntöä 2 lisäämään do-operaattori

$$P(y|z, do(X = x)) = P(y|do(Z = z), do(X = x)),$$

koska ehto $Y \perp\!\!\!\perp Z|X$ pätee graafissa $G_{\overline{XZ}}$. Nyt voidaan poistaa $do(X = x)$ käyttäen sääntöä 3 (koska $Y \perp\!\!\!\perp X|Z$ graafissa $G_{\overline{XZ}}$). Siten saadaan

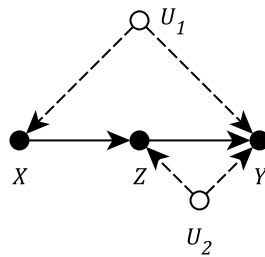
$$P(y|z, do(X = x)) = P(y|do(Z = z)),$$

minkä oikea puoli on jo laskettu yhtälössä (7). Nyt sijoittamalla yhtälöön (8) yhtälöt (3) ja (7) saadaan

$$P(y|do(X = x)) = \sum_z \left[\sum_{x'} P(y|x', z)P(x') \right] P(z|x), \quad (9)$$

missä x' on vain merkintä, jotta summattava x' erottuu havaitusta arvosta x . Yhtälö (9) on erityistapaus etuovikriteeristä [Pearl, 2009, s. 81].

On osoitettu, että Seuraus 1 ei ole pelkästään riittävä, vaan myös välttämätön ehto identifioituvuudelle. Kausaalilaskennan säännöillä pystytään johtamaan kaikkien identifioituvien kausaalivaikutusten lausekkeet [Shpitser and Pearl, 2006]. Annetusta lausekkeesta ja graafista on kuitenkin hyvin vaikea todeta suoraan vaikutuksen identifioituvuus tai identifioitumattomuus. Tarkastellaan esimerkkinä kuvan 9 graafia H , joka muistuttaa laskujen 1–3 graafia G .



Kuva 9: Graafi H , jossa kausaalivaikutus $P(y|do(x))$ ei identifioitu.

Muuttuja U_2 aiheuttaa nyt sen, että yhtälöä (4) ei voidakaan kirjoittaa muodossa (7), sillä $Z \not\perp\!\!\!\perp Y|X$ graafissa $H_{\underline{Z}}$. Tämänkaltainen nopea heuristiikka osoittaa, että edellistä päättelyä ei nyt voida soveltaa uudelleen graafissa H . Tämä ei kuitenkaan vielä riitä, sillä voihan olla jokin toinen äärellinen jono kausaalilaskennan sääntöjä, jolla vaikutus saadaan ilmaistua

ilman do-operaattoria. Mikäli eri sääntöjä yrittää soveltaa graafissa H , huomaa kuitenkin pian ettei yksikään säännöistä tai todennäköisyyslaskennan perustyökaluista auta ongelman ratkaisussa.

Kyseinen vaikutus on todellisuudessa identifioitumaton, mutta tätä ei voi todeta pelkästään siten, ettei laskusääntöjen avulla löytynyt haluttua tulosta. Kausaalivaikutuksen identifioitumattomuus voidaan todeta joko suoraan määritelmällä esittämällä kaksi kausaalimallia M_1 ja M_2 , joissa havaittujen muuttujien jakaumat yhtenevät, mutta kausaalivaikutuksen jakaumat eivät. Tämä on usein hyvin työlästä, ja ongelmaan onkin kehitetty uusia lähestymistapoja, joista eräs on identifioituvuusalgoritmi, jonka kehittivät Shpitser and Pearl [2006]. Kyseisen algoritmin avulla voidaan selvittää minkä tahansa kausaalivaikutuksen identifioituvuus tai identifioitumattomuus. Mikäli kausaalivaikutus on identifioituva, tuottaa algoritmi myös tämän lausekkeen havaittujen jakaumien avulla.

Tehtäviä 3.

1. *Todista Lause 2 (Takaovikriteeri).*

2 Kontrafaktuaaliset tilanteet

Kontrafaktuaaleilla tarkoitetaan todellisuuden vastaisia tilanteita, kuten seuraavassa.

Jos olisimme toimineet toisin ja asettaneet muuttujan X arvoksi jonkin sen havaitusta arvosta poikkeavan arvon, niin mikä olisi ollut muuttujan Y jakauma?

Nyt muuttujalle X asetettu arvo on eri kuin havaittu eli todellisuuden vastainen, siitä siis nimi kontrafaktuaali. Voimme siis käsitellä kontrafaktuaalisessa tilanteessa Y :n jakaumaa do-operaattorin avulla.

2.1 Kontrafaktuaalilaskentaa

Tässä kappaleessa käsittelemme kontrafaktuaalikysymyksiä ja niiden laskentaa. Aluksi käydään läpi kontrafaktuaalilaskennan perusteita (kausaalimalliin nojautuen) ja sen jälkeen käsittelemme kaksosgraafia ja lyhyesti d-separoituvuutta kaksosgraafissa. Kontrafaktuaalitodennäköisyyksien laskemiseksi tarvitaan kausaalimallin alimalli ja määritellään potentiaalivasteen käsite.

Määritelmä 8 (Potentiaalivaste, Causality Def. 7.1.4). *Olkoot X ja Y kaksi osajoukkoa havaituista muuttujista V . Muuttujajoukon Y potentiaalivaste toiminnolle $do(X = x)$, merkitään $Y_x(u)$, on ratkaisu Y :lle funktioiden joukossa F_x , eli $Y_x(u) = Y_{M_x}(u)$.*

On syytä huomata, että $Y_x(u)$ on sama kuin $h_Y(u)$ funktioiden joukossa $H_x(U)$. Jos X ja Y ovat kuten Määritelmässä 8, niin kontrafaktuaalilause "Y saisi arvon y (tilanteessa u), jos X olisi ollut x " kirjoitetaan yhtäsuuruutena $Y_x(u) = y$, missä $Y_x(u)$ on Y :n potentiaalivaste toiminnolle $do(X = x)$. Voidaan myös puhua kontrafaktuaalilauseesta $Y_x = y$, joka ei rajoitu tilanteeseen u . Esimerkiksi todennäköisyys $P(Y_x = y)$ voidaan laskea, ja se on

$$P(Y_x = y) = \sum_{\{u|Y_x(u)=y\}} P(u). \quad (10)$$

Kaavassa (10) summaus käy läpi kaikki taustamuuttujien u realisaatiot, jotka toteuttavat kontrafaktuaalilauseen $Y_x(u) = y$.

Monimutkaisempiakin kontrafaktuaalitodennäköisyyksiä voidaan laskea. Todennäköisyys $P(Y_x = y, X = x')$ tarkoittaa yhteisjakaumaa "Y saisi arvon y , jos olisi ollut $X = x$ vaikka todellisuudessa havaittiin $X = x'$ ". Sen sijaan todennäköisyys $P(Y_x = y, Y_{x'} = y')$ tarkoittaa "Y saisi arvon y , jos olisi ollut $X = x$ ja Y saisi arvon y' , jos olisi ollut $X = x'$ ". Tällaisen todennäköisyyden

laskeminen on mahdollista, koska kontrafaktuaalilauseet eivät (välttämättä) ole todellisia tapahtumia. Vaikka kontrafaktuaalitalanteet, kuten edellä, voivat olla hyvinkin monimutkaisia, niin kaavasta (10) huomataan, että laskettiinpa mitä kontrafaktuaalia tahansa niin kyse on oikeastaan vain sellaisten taustamuuttujien realisaatioiden todennäköisyyksien yhteenlaskusta, joissa kontrafaktuaalilause on totta. Edellä mainittujen kontrafaktuaalilauseiden todennäköisyydet lasketaan seuraavasti:

$$P(Y_x = y, X = x') = \sum_{\{u|Y_x(u)=y \& X(u)=x'\}} P(u)$$

ja

$$P(Y_x = y, Y_{x'} = y') = \sum_{\{u|Y_x(u)=y \& Y_{x'}(u)=y'\}} P(u).$$

Käyttäen Lausetta 4 voidaan laskea kontrafaktuaalitodennäköisyyksiä ehdolla havaittu aineisto.

Pohditaan lyhyesti kontrafaktuaalilauseiden yhteisjakaumaa yleisemmin. Yleisesti voidaan laskea

$$P(Y_x = y, Z_w = z), \tag{11}$$

mikä tarkoittaa "Y saisi arvon y , jos olisi ollut $X = x$, ja Z saisi arvon z , jos olisi ollut $W = w$ ". Se ei ole sama kuin $P(Y = y, Z = z | do(X = x), do(W = w))$. Todennäköisyys (11) on määritelty mille tahansa muuttujajoukoille Y, X, Z ja W , ja joukkojen ei tarvitse olla erilliset. Huomataan, että voidaan laskea mielivaltaisen monen kontrafaktuaalilauseen yhteisjakauma, esim. $P(Y_{x_1} = y_1, \dots, Y_{x_n} = y_n)$.

Lause 4. (Kontrafaktuaalitodennäköisyyksien laskeminen, *Causality Theorem 7.1.7*) Kontrafaktuaalilauseen "jos A niin silloin B " todennäköisyys ehdolla aineisto e , merkitään $P(B_a|e)$, voidaan laskea seuraavasti.

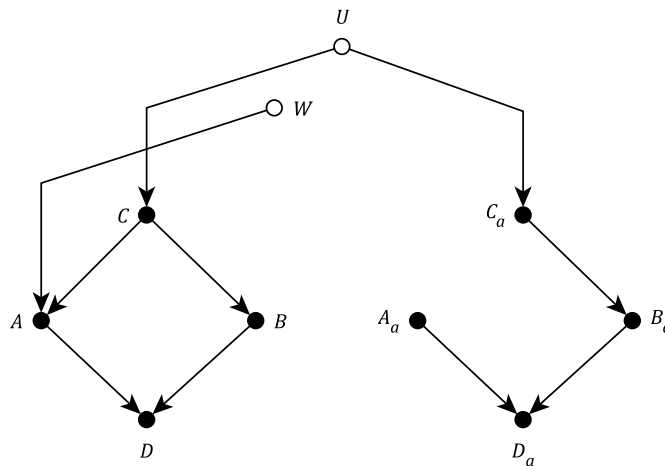
1. Päivittäminen – Päivitä jakauma $P(u)$ aineiston e avulla jakaumaksi $P(u|e)$.
2. Toiminnolla muokkaaminen – Muokkaa kausaalimallia M toiminnolla $do(A = a)$, jossa A on vasteen vanhempi tai esivanhempi, näin saat alimallin M_a .
3. Ennustaminen – käytä toiminnolla muokattua alimallia $M_a = \langle U, V, F_a, P(u|e) \rangle$ laskeaksesi todennäköisyysjakauman $P(B_a|e)$.

Huomautettakoon tässä kohtaa, että mikäli funktio u :lta havaituille soluille on tuntematon, on tyypillisesti mahdotonta laskea Päivittämis-askel. Samoin liittyen kaavaan (10) on vaikeaa tietää mitkä u :t ovat sellaisia, joille pätee $Y_x(u) = y$, mikäli em. funktio on tuntematon. Lukijan on siis syytä huomata, että Esimerkissä 5 tarvittavat funktiot tunnetaan (ne määritellään yhtäsuuruuksina). Käytännön sovellutuksissa tällaisia funktioita ei välttämättä tunneta.

2.2 Kaksosgraafi

Kaksosgraafi on graafi, jossa kuvataan sekä kausaalimalli että sen intervention seurauksena saatu alimalli rinnakkain. Kaksosgraafin tarkoituksena on helpottaa kontrafaktuaalipäätelyä, samaan tapaan kuin kausaalimallin graafin tarkoituksena on helpottaa kausaalipäätelyä. Myös kaksosgraafissa voidaan tarkastella muuttujien välisiä ehdollisia riippumattomuuksia d-separoituvuuden avulla.

Kaksosgraafi liittyy kontrafaktuaaliseen tilanteeseen, jossa ajatellaan että oltaisiin tehty interventio johonkin muuttujaan. Tällöin kaksosgraafi muodostuu kausaalimallista ja sen intervention jälkeisestä alimallista. Alimallissa kausaalimallin havaitut muuttujat on kopioitu, joten ne piirretään omina soluinaan. Havaitsemattomat muuttujat ovat samat sekä kausaalimallille että alimallille.



Kuva 10: Kaksosgraafi esimerkkiin 5. Graafin vasen puolisko (ml. U) kuvaa kausaalimallin M ja oikea puolisko (ml. U) kuvaa alimallin M_a .

Eräs kaksosgraafi on Kuvassa 10, jota käsittelemme myös myöhemmin

esimerkissä. Varsinainen kausaalimalli, kuvaajan vasen puolisko, koostuu havaituista muuttujista A, B, C ja D sekä havaitsemattomiasta tekijöistä (ns. taustamuuttujat), U ja W .¹ Kontrafaktuaalista tilannetta kuvaava alimalli M_a on kuvaajan oikealla puoliskolla ja se yhtyy kausaalimalliin havaitsemattomien muuttujien osalta. Muuttujasta W ei kuitenkaan mene suunnattua särmää A_a :han, koska alimallin graafissa A_a :han saapuvat särmät on poistettu.

Esimerkki 5 (Kaksosgraafi). *Kuva 10 esittää teloitusryhmän toimintaa. Teloitusryhmä muodostuu kahdesta ampujasta A ja B sekä komentajasta C . Ampuja A on kokematon, ja saattaa siksi ampuu hermostuneisuudesta johtuen. Ampuja B toimii aina käskyjen mukaan. Mikäli komentaja komentaa "Ampukaa!", niin C saa arvon 1, muutoin $C = 0$. Jos ampuja A ampuu, niin merkitään $A = 1$. Muuttuja W määrittää ampuuko ampuja A hermostuneisuudesta johtuen, joten $A = 1$ jos $C = 1$ tai $W = 1$, muutoin $A = 0$. Ampujalle B asetetaan, jos $C = 1$, niin $B = 1$, muutoin $B = 0$. Muuttuja U on teloituspäätös, joka ei ole tiedossamme. Samoin emme tiedä ampuuko ampuja A hermostuneisuuden takia, eli W on tuntematon. Halutaan laskea kontrafaktuaalinen todennäköisyys menehtyykö henkilö D , jolloin $D = 1$, mikäli D ei menehdy, niin $D = 0$.*

Koska muuttujat U ja W ovat havaitsemattomia, niille määräytyy todennäköisyysjakauma $(u, w) \sim P(U, W)$, kuten kausaalimallissa aina. Oletetaan $U \perp W$. Tätä jakaumaa voidaan mallintaa seuraavasti:

$$P(u, w) = \begin{cases} pq & \text{jos } u = 1, w = 1 \\ p(1 - q) & \text{jos } u = 1, w = 0 \\ (1 - p)q & \text{jos } u = 0, w = 1 \\ (1 - p)(1 - q) & \text{jos } u = 0, w = 0. \end{cases} \quad (12)$$

Jakaumaa (12) voidaan sanoa ns. priorijakaumaksi.

Päivittämällä jakaumaa aineistolla $D = 1$ saadaan Bayesin kaavalla posteriorijakauma

$$\begin{aligned} P(u, w|D = 1) &= \frac{P(u, w)P(D = 1|u, w)}{P(D = 1)} \\ &= \begin{cases} \frac{p(u, w)}{1 - (1 - p)(1 - q)} & \text{jos } u = 1 \text{ tai } w = 1, \\ 0 & \text{jos } u = 0 \text{ ja } w = 0. \end{cases} \end{aligned} \quad (13)$$

¹Sivuhuomautuksena mainittakoon että Pearl merkitsee kirjassaan tässä kohtaa muuttujia U ja W graafissa tummennetulla pallolla.

Edellä huomataan, että henkilö menehtyy kaikissa muissa tilanteissa paitsi, kun $u = 0$ & $w = 0$. Siksi $P(D = 1) = 1 - (1 - p)(1 - q)$. Aineistolla $D = 1$ päivittäminen sai aikaan sen, että päivittämisen jälkeen on mahdotonta (todennäköisyys on 0), että olisi $u = 0$ ja $w = 0$. Näin siis saimme aineiston avulla lisätietoa havaitsemattomien tekijöiden jakaumasta.

Posteriorijakauman $P(u, w | D = 1)$ avulla voidaan laskea **toiminnolla muokatun** alimallin M_a (kaksosgraafin oikea puolisko) todennäköisyyksiä.

Niitä molempia käyttäen voimme **ennustaa** eli laskea kontrafaktuaalisia todennäköisyyksiä, esimerkiksi voidaan laskea todennäköisyys että "henkilö D jää eloon, kun A ei ammu (toiminto $do(A = 0)$), kun todellisuudessa D menehtyi", eli

$$P(D_a = 0 | D = 1) = P(U = 0 | D = 1) = \frac{q(1 - p)}{1 - (1 - p)(1 - q)}, \quad (14)$$

mikä on sama kuin todennäköisyys että teloitusmääräystä ei annettu, mutta henkilön D havaittiin menehtyneen.

Esimerkki 6 (Kaksosgraafi ja d -separoituvuus). d -separoituvuutta muuttujajoukkojen välillä on mahdollista tarkastella myös kaksosgraafissa, sillä se on edelleen nimensä mukaisesti graafi. Kaksosgraafin avulla voidaan siis tutkia muuttujien ja kontrafaktuaalimuuttujien välillä vallitsevia ehdollisia riippumattomuuksia samaan tapaan kuin graafeissa muutenkin.

Esimerkiksi kontrafaktuaalinen muuttuja C_a on ehdollisesti riippumaton muuttujajoukosta $\{A, B\}$ ehdolla muuttuja C kuvan 8 kaksosgraafissa, sillä C d -separoi polut $A \leftarrow C \leftarrow U \rightarrow C_a$ ja $B \leftarrow C \leftarrow U \rightarrow C_a$. Vastaavasti voidaan myös tarkastella kontrafaktuaalien välisiä ehdollisia riippumattomuuksia. Kontrafaktuaalinen muuttuja A_a on ehdollisesti riippumaton kontrafaktuaalisesta muuttujasta B_a ehdolla \emptyset , sillä D_a toimii käänteisenä haarukkana polulla $A_a \rightarrow D_a \leftarrow B_a$. Muuttujalla D_a ehdollistaminen sen sijaan avaisi kyseisen polun, jolloin A_a ja B_a eivät enää olisi ehdollisesti riippumattomia (ehdolla D_a).

Tehtäviä 4.

1. Laske todennäköisyys lausekkeessa (14) tilanteessa, jossa $P(U = 1) = 0.5$ ja $P(W = 1) = 0.1$.
2. Piirrä kaksosgraafi Kuvan 8 graafille tilanteessa, jossa on tehty interventio $X = 1$. Laske kontrafaktuaalinen todennäköisyys $P(Y_x = 1 | Y =$

1) kun kausaalimallin muuttujat määräytyvät seuraavasti

$$\begin{aligned}f_X(u) &= 1 - u \\f_Z(x) &= x \\f_Y(z, u) &= \max(2z, u),\end{aligned}$$

ja kun $P(U = 1) = p$ ja $P(U = 0) = 1 - p$.

Lähteet

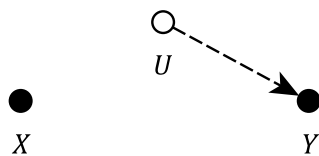
Chi-Ling Chen, Tim J Gilbert, and Janet R Daling. Maternal smoking and Down syndrome: the confounding effect of maternal age. *American Journal of Epidemiology*, 149(5):442–446, 1999.

Judea Pearl. *Causality: models, reasoning and inference*. Cambridge University Press, 2009.

Ilya Shpitser and Judea Pearl. Identification of joint interventional distributions in recursive semi-markovian causal models. In *Proceedings of the 21st National Conference on Artificial Intelligence - Volume 2*, pages 1219–1226. AAAI Press, 2006.

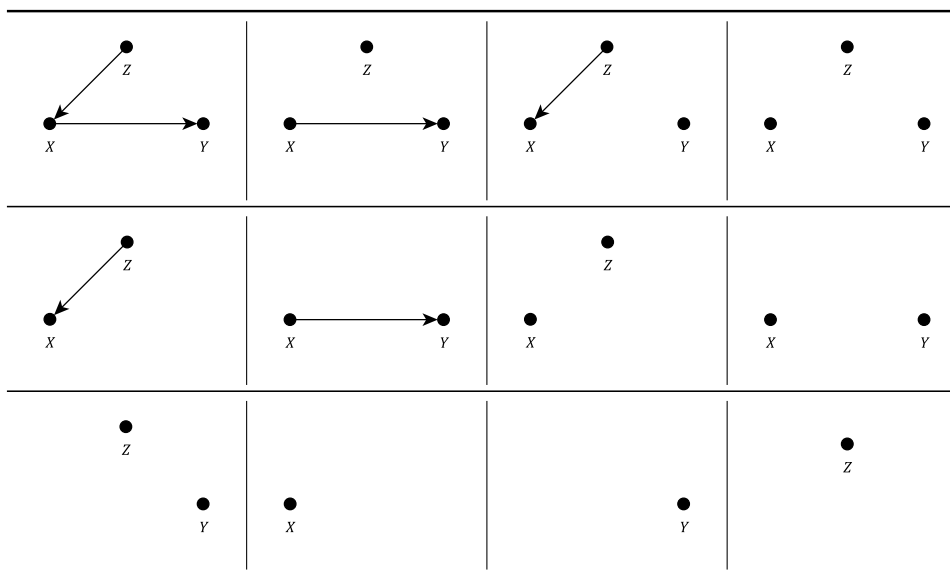
3 Harjoitustehtävien ratkaisut

1. 1.1 Poistamalla solmuun X saapuvat särmät saadaan seuraava graafi:



$G_{\overline{X}}$ kuvan 2b graafille G .

1.2 Graafin $G = \langle \{X, Y, Z\}, \{(X, Y), (Z, X)\} \rangle$ itsensä lisäksi sen ali-graafija ovat



1.3 Kolme solmua sisältäviä suuntaamattomia graafeja voidaan muodostaa

$$\binom{3}{0} + \binom{3}{1} + \binom{3}{2} + \binom{3}{3} = 2^3 = 8$$

kappaletta, eli nolla, yksi, kaksi ja kolme särmää sisältävien graafien lukumäärien summan verran. Termit $\binom{3}{k}$, $k = 0, 1, 2, 3$ kertovat siis kuinka monta kolme solmua ja k särmää sisältävää suuntaamatonta graafia on mahdollista muodostaa. Summan termejä

voidaan nyt painottaa särmien suunnilla ja määrillä. Mahdollisista tavoista valita särmät täytyy kuitenkin poistaa vaihtoehdot, jotka synnyttävät silmukoita. Jos kahden solmun X ja Y välillä on särmä $X \rightarrow Y$, sulkee tämä vaihtoehdon $Y \rightarrow X$ pois. Kahden solmun välillä voi siis kulkea särmä ainoastaan yhteen suuntaan, jolloin jokaiselle suuntaamattomalle särmälle on tasan 2 suuntavaihtoehtoa. Tällä rajoituksella ei yhden ja kahden solmuparin tapauksessa voi vielä muodostua silmukkaa, mutta myös kolmen solmun väliset silmukat on estettävä. Tällaisia silmukoita on 2 kappaletta, $X \rightarrow Y \rightarrow Z$ ja $X \leftarrow Y \leftarrow Z$. 2^3 tavasta valita kolmelle särmälle suunta on siis poistettava nämä kaksi silmukan tuottavat vaihtoehdot. Laskemalla eri vaihtoehdot yhteen saadaan nyt

$$\binom{3}{0} + 2 \cdot \binom{3}{1} + 2^2 \cdot \binom{3}{2} + (2^3 - 2) \cdot \binom{3}{3} = 25.$$

Suunnattuja silmukattomia graafeja voidaan siis muodostaa 25 kappaletta tilanteessa, jossa graafi sisältää 3 solmua.

- 1.4 a) Osoitamme, että jokainen X : ja Y :n lineaarikombinaatio on normaalijakautunut. Yksinkertaisella laskulla saadaan (ol. $\beta \neq 0$ tai $\gamma \neq 0$):

$$U_1 = \frac{\beta X}{\beta^2 + \gamma^2} + \frac{\gamma Y}{\beta^2 + \gamma^2}, \quad U_2 = \frac{\gamma X}{\beta^2 + \gamma^2} - \frac{\beta Y}{\beta^2 + \gamma^2}.$$

Siis jokainen X :n ja Y :n lineaarikombinaatio on myös U_1 :n ja U_2 :n lineaarikombinaatio, ja siten normaalijakautunut.

- b) Todennäköisyyslaskennan laskusääntöjen mukaan

$$P(X, Y, U_1, U_2) = P(X, Y|U_1, U_2)P(U_1, U_2).$$

Edelleen U_1 ja U_2 ovat riippumattomia joten

$$P(X, Y, U_1, U_2) = P(X, Y|U_1, U_2)P(U_1)P(U_2).$$

Voidaan myös kirjoittaa

$$P(X, Y|U_1, U_2) = P(X|Y, U_1, U_2)P(Y|U_1, U_2).$$

Koska U_1 ja U_2 määräävät Y :n arvon, niin $X \perp\!\!\!\perp Y|U_1, U_2$ ja $P(X|Y, U_1, U_2) = P(X|U_1, U_2)$. Nyt

$$\begin{aligned} P(X, Y, U_1, U_2) &= P(X|U_1, U_2)P(Y|U_1, U_2)P(U_1)P(U_2) \\ &= P(X|PA_X)P(Y|PA_Y)P(U_1|PA_{U_1})P(U_2|PA_{U_2}), \end{aligned}$$

joten $P(X, Y, U_1, U_2)$ on kuvan 4 graafin G kanssa yhteensopiva.

2. 2.1 a) Analyysi voidaan tehdä R:ssä seuraavasti

```
> dat <- expand.grid(Y = c(0, 1), S = c(0, 1), A = c(0, 1))
> dat$n <- c(344, 84, 59, 18, 80, 201, 56, 158)
> fit.a <- glm(Y ~ S, family = binomial,
  weights = n, data = dat)
```

Muuttujien Y , S ja A eri kombinaatioista on eri määrä havaintoja, mikä tulee ottaa huomioon `glm`-funktion `weights`-parametrilla. Tässä mallissa tupakoinnin regressiokerroin on merkitsevä.

b) Lisätään ikä malliin

```
> fit.b <- glm(Y ~ S + A, family = binomial,
  weights = n, data = dat)
```

Tässä mallissa tupakoinnin regressiokerroin ei ole merkitsevä, ikäryhmän sen sijaan on.

2.2 Artikkelin aineisto voidaan kirjoittaa R:ssä seuraavasti (tässä ikäryhmien 'Unknown'-luokan havainnot on poistettu yksinkertaisuuden vuoksi)

```
> dat <- expand.grid(Y = c(0, 1), S = c(0,1), A = c(0, 1))
> dat$n <- c(5214, 421, 1411, 112, 611, 186, 108, 15)
```

Mallit voidaan nyt sovittaa aivan kuten edellisessä tehtävässä. Tulokset ovat vastaavat: kun ikä ei ole mallissa saadaan tupakoinnille tilastollisesti merkitsevä regressiokerroin, kun ikä on mallissa ei kerroin enää ole merkitsevä. Ehdolliset jakaumat $P(Y|S)$ ja $P(Y|S, A)$ ovat nyt:

$$P(Y = 1|S = 0) = 0.094$$
$$P(Y = 1|S = 1) = 0.077,$$

ja

$$P(Y = 1|S = 0, A < 35) = 0.076$$
$$P(Y = 1|S = 1, A < 35) = 0.067$$
$$P(Y = 1|S = 0, A \geq 35) = 0.221$$
$$P(Y = 1|S = 1, A \geq 35) = 0.199.$$

3. 3.1 Todennäköisyyslaskennan laskusääntöjen avulla voidaan kirjoittaa

$$P(y|do(X = x)) = \sum_{pa_X} P(y|do(X = x), pa_X)P(pa_X|do(X = x)).$$

Oletuksen mukaan kaikki solmun X vanhemmat ovat havaittuja, jolloin PA_X d-separoi kaikki solmuun X saapuvat polut. Tällöin graafissa $G_{\underline{X}}$ ei solmusta X lähde polkuja ja kaikki siihen saapuvat polut on d-separoitu. Voimme hyödyntää kausaalilaskennan sääntöä 2, sillä $Y \perp\!\!\!\perp X|PA_X$ graafissa $G_{\underline{X}}$, ja saamme

$$P(y|do(X = x), pa_X) = P(y|x, pa_X).$$

Edelleen $PA_X \perp\!\!\!\perp X$ graafissa $G_{\overline{X}}$, sillä millä tahansa polulla solmujoukosta PA_X solmuun X on oltava käännteinen haarukka (collider). Kausaalilaskennan sääntöä 3 käyttäen saamme nyt

$$P(pa_X|do(X = x)) = P(pa_X).$$

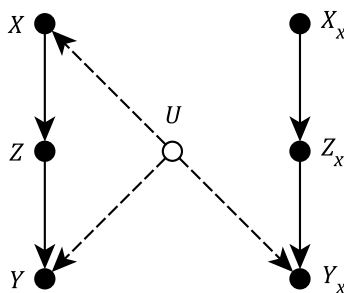
4. 4.1 Sijoittamalla lukuarvot saadaan

$$\frac{q(1-p)}{1-(1-p)(1-q)} = \frac{0.1 \cdot (1-0.5)}{1-(1-0.5)(1-0.1)} = \frac{0.05}{0.55} = \frac{1}{11} \approx 0.091$$

- 4.2 Käytämme Lausetta 4 kontrafaktuaalitodennäköisyyden laskemiseksi. Ensimmäiseksi on päivitettävä jakauma $P(U)$ aineistolla $Y = 1$ jakaumaksi $P(U|Y = 1)$. Funktion f_Y avulla saamme

$$1 = f_Y(z, u) = \max(2(1-u), u).$$

On siis oltava $U = 1$, eli $P(U = 1|Y = 1) = 1$. Seuraavaksi muokkaamme kausaalimallia toiminnolla $X = 1$, jolloin saamme kaksosgraafin kuvan 8 graafille.



Lopuksi voimme laskea todennäköisyyden $P(Y_x = 1|Y = 1)$ (enustaminen) määrittämällä kontrafaktuaalisen muuttujan Y_x arvon, joka on

$$f_{Y_x}(z_x, u) = \max(2z_x, u) = \max(2x_x, u) = \max(2, 1) = 2.$$

Siis $P(Y_x = 1|Y = 1) = 0$.