

Notater

Sverre Langaas

28. april 2021

Innhold

1	Sannsynlighet	1
1.1	Aksiom og teknisk rammeverk	1
1.2	Sannsynlighetsregning	3
1.2.1	Betinget sannsynlighet og uavhengighet	4
1.3	Tilfeldige variabler	5
1.4	Univariat fordeling	6
1.4.1	Momenter	7
1.4.2	Kvantiler	8
1.5	Multivariat fordeling	9
1.5.1	Simultanfordeling	9
1.5.2	Betinget fordeling	11
1.5.3	Uavhengighet	12
1.5.4	Generalisering til N dimensjoner	12
1.5.5	Tilfeldig utvalg	14
1.5.6	Hierarki	15
1.6	Transformasjon av tilfeldig variabel	16
1.6.1	Fordeling til transformert variabel	16
1.6.2	Kontinuerlig	17
1.6.3	Diskret	18
1.6.4	Bivariat transformasjon	18
1.7	Momentgenererende funksjoner	19
1.8	Projeksjon av tilfeldige variabler	21
1.8.1	Projektering i L_2	22
1.8.2	Payoff	23
2	Stokastiske prosesser	25
2.1	Asymptotisk teori	26
2.1.1	Store talls lov	27
2.1.2	Sentralgrenseteoremet	28
2.1.3	Delta-metoden	28
2.1.4	LNN og CLT med flere variabler	29

2.1.5	LLN og CLT med avhengighet mellom observasjoner	29
2.2	Markov-kjeder	29
2.2.1	Overgangssannsynlighet	30
2.3	Annet	30
3	Noen kjente fordelinger	31
3.1	Normalfordeling	31
3.1.1	Normalfordelt utvalg	31
3.1.2	Truncated normalfordeling	32
3.1.3	Multivariat normalfordeling	33
3.2	Fordelinger assosiert med normalfordeling	33
3.2.1	χ^2 -fordeling	33
3.2.2	t-fordeling	34
3.2.3	F-fordeling	34
3.3	Fordelinger fra bernoulli-prosess	34
3.3.1	Binomialfordeling	34
3.3.2	Negativ binomialfordeling	34
3.3.3	Geometrisk fordeling	34
3.3.4	Multinomialfordeling	35
3.4	Fordelinger fra poisson-prosess	35
3.4.1	Poissonfordeling	35
3.4.2	Eksponentialfordeling	35
3.5	Andre fordelinger	35
3.5.1	Uniformfordeling	35
3.5.2	Gammafordeling	36
3.5.3	Betafordeling	37
3.5.4	Hypergeometrisk	37
4	Inferens	38
4.1	Motivasjon	38
4.1.1	Generalisering	39
4.2	Formelt rammeverk	40
4.2.1	Modell	41
4.2.2	Lære egenskaper til fordeling	41
4.2.3	Utvalgsfordelingen til estimatorer	42
4.2.4	Estimere utvalgsfordelingen	43
4.3	Egenskaper til estimatorer	44
4.4	Estimering	45
4.4.1	Punktestimat	45
4.4.2	Konfidensmengder	46

4.4.3	Hypotesetester	47
4.4.4	Modellering	50
4.5	David og Mac	51
5	Momentestimatorer	53
5.1	Utvalgsanalogprinsippet	53
5.1.1	Motivere OLS som utvalgsanalog	53
5.2	Momentestimator	54
5.2.1	Egenskaper	55
5.3	GMM	55
5.3.1	2SLS	57
6	Maximum likelihood	59
6.1	Begreper	60
6.1.1	Score	61
6.1.2	Informasjon	62
6.1.3	Alternativ utledning	63
6.2	Eksempler	64
6.2.1	Bernoulli	64
6.2.2	Normalfordeling med kjent varians	66
6.2.3	Uniform	66
6.2.4	Andre hendelser	66
6.3	Oppsummere informasjon fra likelihoodfunksjonen	67
6.3.1	Kvadratisk tilnærming	67
6.3.2	Konfidensintervall	68
6.4	Likelihood i flere dimensjoner	69
6.4.1	Betinget likelihood	69
6.4.2	Generell fremgangsmåte til å finne likelihood til betinget fordeling	69
6.4.3	Betinget normal	70
6.4.4	Betinget bernoulli	70
6.5	Prinsipp for å utlede tester	71
6.5.1	Greier fra DM	71
6.5.2	Wald-test	72
6.5.3	Likelihood ratio	72
6.5.4	Lagrange multipliar	73
6.6	Egenskaper ved feilspesifikasjon	73
6.6.1	Total variation distance og KL-divergence	73
6.6.2	MLE fra empirisk risikominimering	74
6.6.3	Kvasi-MLE	75
6.6.4	Extremum estimators	76

7	Lineær regresjon	77
7.1	Egenskap ved simultanfordeling	77
7.1.1	Projeksjon	78
7.1.2	Dekomponering av varians	78
7.1.3	Tolkning av feilledd	78
7.2	Numeriske egenskaper	78
7.2.1	Ortogonal projeksjon	78
7.2.2	Frisch-Waugh-Lovell	79
7.2.3	In-sample fit	80
7.3	Statistiske egenskaper	81
7.3.1	Små utvalg	81
7.3.2	Store utvalg	82
7.3.3	Presisjon til koeffisient	82
7.3.4	Presisjon til prediksjon	82
7.3.5	Residual	83
7.4	Hypotesetester	83
7.4.1	Lineære restriksjoner av koeffisient	83
7.4.2	Diagnose	83
7.5	Funksjonell form	84
7.6	Fordeling til feilledd	85
7.6.1	Generalisert minste kvadrat	85
7.6.2	Vektet minste kvadrat	86
7.7	Robust estimering (sandwich)	87
7.8	Annet	87
8	Tidsserier	89
8.1	Deskriptivt	90
8.1.1	Stasjonaritet	91
8.1.2	Empirisk	91
8.2	Modellering	92
8.2.1	AR(p)	93
8.2.2	Autoregressiv, AR(k)	93
8.2.3	MA(q)	93
8.2.4	ARMA(p,q)	93
8.3	Tidsserier	93
8.3.1	Lineær trend	94
8.4	Annet	94
8.4.1	Mer annet	95

9	Statistisk læring	96
9.1	Bakgrunn og oversikt	96
9.1.1	Generativ modell	97
9.1.2	Utfordringer	97
9.2	Empirisk risikominimering	99
9.2.1	Kryssvalidering	100
9.2.2	Dekomponering av risiko med kvadratisk tap	101
9.3	Lineær regresjon	102
9.3.1	Feature space	102
9.3.2	Regularisering	104
9.4	Andre regresjonsmetoder	105
9.4.1	Splines	105
9.4.2	Ikke-paramerisk regresjon	105
9.4.3	Kvantilregresjon	106
9.5	Klassifikasjon	106
9.5.1	Logistisk regresjon	107
9.5.2	Bayesianske metoder	109
9.5.3	KNN	110
9.5.4	Support vector machines	110
9.5.5	Beslutningstrær	110
9.6	Ensemble	111
9.6.1	Bagging og tilfeldig skog	111
9.6.2	Boosting	112
9.6.3	Stacking	112
9.7	Vurderingskriterier	112
9.7.1	Confusion matrix	112
9.7.2	Presisjon vs Recall trade-off	113
10	Læring uten tilsyn	115
10.1	Dimensjonalitetsreduksjon	115
10.1.1	Principal component analysis	115
10.1.2	Andre metoder	116
10.2	Clustering	116
10.2.1	K-means	116
10.3	Tetthetsestimering	117
10.3.1	Mål på avstand mellom tetthetsfunksjoner	117
10.3.2	Histogram	118
10.3.3	Kernel density estimation	119
10.3.4	Mixture models	119

11 Økonometri	120
11.1 Programevaluering	120
11.1.1 Potensielle utfall	122
11.1.2 Matching	123
11.1.3 Dårlig kontroll	124
11.1.4 Knytte potensielle utfall til regresjonsligning	126
11.1.5 Målefeil	127
11.1.6 Utelatte variabler	127
11.2 Instrumentelle variabler	129
11.2.1 Estimering	129
11.2.2 Heterogen behandlingseffekt	130
11.2.3 Eksperiment med delvis compliance	133
11.2.4 Simultane ligningssystem	133
11.2.5 Generalisering av wald	135
11.3 Regresjonsdiskontinuitet	135
11.3.1 Fuzzy rdd	137
11.4 Paneldata	137
11.4.1 Dekomponering av feilledd	138
11.4.2 Identifikasjon	139
11.4.3 Estimering	139
11.4.4 Dynamisk panel	141
11.5 Dynamiske modeller	141
11.5.1 Lagged uavhengig variabel	141
11.5.2 Lagged avhengig variabel	142
11.5.3 Instrument	143
11.6 Forskjeller i forskjeller	143
11.6.1 Identifikasjon	144
11.6.2 Flere grupper og flere tidsperioder	144
11.7 Limited Dependent Variable	145
11.7.1 Stokastisk nytte	145
11.7.2 Sensurert regresjon (tobit)	147
11.7.3 Heltallsverdier (Poisson-regresjon)	151
11.8 Modellere seleksjon	152
11.8.1 Avkortet (truncated) regression	153
11.8.2 Heckit	153
12 Kalkulus	156
12.1 Litt bakgrunn	156
12.1.1 Algebra	156

12.2	Litt analyse	157
12.2.1	Følger	157
12.2.2	Rekker	157
12.2.3	Grenser og kontinuitet	157
12.2.4	Topologi	158
12.3	Litt om funksjoner...	158
12.3.1	Real valued functions	158
12.3.2	Inverse funksjoner	158
12.4	Lineær tilnærming av funksjoner	159
12.4.1	Derivasjon	159
12.4.2	Taylor-tilnærming	160
12.5	Noen vanlige funksjoner	161
12.5.1	Polynomial	161
12.5.2	Eksponential og logaritmer	161
12.5.3	Trigonometriske funksjoner	162
12.5.4	Sammensatte funksjoner	162
12.6	Kort om integral	162
12.7	Flervariable funksjoner	163
12.8	Ubetinget optimering	164
12.8.1	Gradient descent	164
12.9	Betinget optimering	164
13	Lineær algebra	165
13.1	Vektorer	165
13.2	Matriser	166
13.2.1	Derivasjon	167
13.3	Lineære transformasjoner	168
13.4	Ortogonale projeksjoner	169
13.5	Kvadratisk form	170
14	Matematisk tenkemåte	171
14.1	Logikk	171
14.1.1	Utsagnslogikk	172
14.1.2	Første ordens logikk	174
14.2	Mengdelære	175
14.3	Relasjoner	178
14.3.1	Tillukning	179
14.3.2	Funksjoner	180
14.4	Grafer	181
14.4.1	Trær	182

14.4.2	Vektete grafer	182
14.5	Induksjon	182
14.6	Kombinatorikk	182
14.6.1	Multiplikasjonsprinsippet	182
14.6.2	Permutasjoner	182
14.6.3	Kombinasjoner	183
14.7	Informasjonsteori	183
14.7.1	Entropi	183

Kapittel 1

Sannsynlighet

Vi kan bruke sannsynlighetsteori til å modellere situasjoner med usikkerhet eller uforutsigbarhet. Det innebærer at vi ikke kan beregne eksakte egenskaper til en tilstand med utgangspunkt i den informasjonen vi ha tilgjengelig. I statistikk observerer vi gjerne et utvalg eller delmengde av en populasjonen.¹ På bakgrunn av denne ufullstendige informasjonen vil vi si noe om egenskaper til hele populasjonen. Dette vil det ikke være mulig å gjøre eksakt. Sannsynlighetsteori gir oss likevel et rammeverk for å håndtere og kvantifisere usikkerheten. Mer generelt gir det oss et rammeverk for å håndtere vår ignoranse og kan brukes i veldig mange sammenhenger.

1.1 Aksiom og teknisk rammeverk

Sannsynlighetsteorien tar utgangspunkt i et stokastisk forsøk. Dette består for det første av et utfallsrom Ω som er mengden av alle de ulike utfallene ω som kan bli realisert i eksperimentet. Utfallsrommet er fullstendig og gjensidig utelukkene slik at ett, og bare ett, utfall blir realisert. Hvis vi tenker at forsøket blir gjentatt flere ganger vil vi gjerne observere ulike utfall. En mulig tolkning av sannsynligheten til utfall er den relative frekvensen i uendelig gjentatte forsøk.

Det er en utfordring at vi kan ha utfallsrom som er ikke-tellbare. Det vil si at det ikke er mulig å liste opp de ulike utfallene som $\Omega = \{\omega_1, \omega_2, \dots\}$ ². Et eksempel på et slikt utfallsrom er enhetsdisken

$$\Omega = \{(i, j) \in \mathbb{R}^2 : |i + j| \leq 1\} \quad (1.1)$$

¹Terminologien stammer fra da statistikk var omtrent ensbetydende med biostatistikk. Mer generelt kan vi betrakte populasjonen som en datageneringsprosess og utvalget er en følge med realiseringer fra denne prosessen.

²I praksis er alle målbare utfall tellbare siden vi kun kan måle utfallene med begrenset presisjon. Vi har altså tall med begrenset antall desimaler slik at det i prinsippet ville være mulig å liste opp utfallene. Det vil likevel være praktisk å behandle utfallsrommet som om det var kontinuerlig siden vi da kan bruke kalkulus i stedet for diskret matematikk.

Det er da ikke mulig å gi sannsynlighet til enkeltutfall (i, j) siden man kan vise at sannsynligheten nødvendigvis er null. Vi tallfester derfor kun sannsynlighet til delmengder og ikke enkelt-utfall.³ Merk at hvis utfallsrommet er tellbart kan vi også angi sannsynlighet til delmengder som kun inneholder enkeltutfall, eks: $A = \{\omega_k\}$. Vi kaller delmengder for hendelser. En hendelse A inntreffer hvis et utfall $\omega \in A$ blir realisert. Hvis sannsynligheten er uniform så kan vi i tellbare utfallsrom finne sannsynlighet for hendelser som antall gunstige utfall delt på antall mulige,

$$\mathbb{P}(A) = \frac{|A|}{|\Omega|} \quad (1.2)$$

Den naturlige generalisering til ikke-tellbare mengder er det relative arealet til mengden A .

$$\mathbb{P}(A) = \frac{\lambda(A)}{\lambda(\Omega)} \quad (1.3)$$

der λ er en funksjon som gir arealet. Det er en utfordring at ikke alle delmengder har et veldefinert mål på areal. Vi avgrenser oss derfor til å se på delmengdene av Ω som oppfører seg bra. Såkalte σ -algebraer har egenskapene vi ønsker. \mathcal{F} er en σ -algebra på Ω hvis

1. $A \in \mathcal{F} \implies A^C \in \mathcal{F}$
2. $A_1, A_2, \dots \in \mathcal{F} \implies \cup_n A_n \in \mathcal{F}$
3. $\Omega \in \mathcal{F}$

Dette sikrer at delmengdene oppfører seg bra, men sikrer ikke at de inneholder alle hendelser vi er interessert i. Et trivielt eksempel er $\mathcal{F} = \{\emptyset, \Omega\}$. I praksis er utfallsmengden enten tellbar slik at vi kan betrakte alle delmengder, eller den består av (en delmengde av) tallinjen \mathbb{R}^N og vi betrakter borel-mengdene $\mathcal{B}(\mathbb{R}^N) = \mathcal{F}$.⁴ Vi kan definere en *probability measure* $\mathbb{P} : \mathcal{F} \rightarrow \mathbb{R}$ som angir sannsynlighet for hendelser. Den må tilfredstille aksiomene:

1. $\mathbb{P}(A) \geq 0, \quad \forall A \in \mathcal{F}$
2. $\mathbb{P}(\Omega) = 1$
3. $\mathbb{P}(A_1 \cup A_2 \cup \dots) = \sum_{n=1}^{\infty} \mathbb{P}(A_n)$ for disjunkte delmengder

Disse egenskapene korresponderer med tolkningen av sannsynlighet som relativt frekvens av hendelser i uendelig gjennsatte forsøk, men aksiomene er agnostisk for tolkning av

³Sannsynlighetsteori bygger i stor grad på mengdelære. Se appendiks for definisjoner og regneregler for mengder.

⁴Inneholder intervall i \mathbb{R} og rektangler i \mathbb{R}^N .

sannsynligheten. Det vil også være mulig å bruke en mer subjektiv oppfatning av sannsynlighet der sannsynlighetsfunksjon tilfredstiller aksiom. Tilsammen utgjør $(\Omega, \mathcal{F}, \mathbb{P})$ et *probability space*.

1.2 Sannsynlighetsregning

Fra aksiomene kan vi utlede diverse regneregler som kan brukes til å finne sannsynlighet for ulike hendelser. Et veldig enkelt og nyttig resultat er at

$$\mathbb{P}(\Omega) = \mathbb{P}(A \cup A^c) = \mathbb{P}(A) + \mathbb{P}(A^c) = 1 \quad (1.4)$$

$$\implies \mathbb{P}(A) = 1 - \mathbb{P}(A^c). \quad (1.5)$$

Samlingen av mengder $\{A, A^c\}$ er et eksempel på en partisjonering av Ω siden de er parvis disjunkt, $A \cap A^c = \emptyset$, og inneholder alle elementene i mengden, $A \cup A^c = \Omega$. Slike partisjoner er praktiske å jobbe med siden vi enkelt kan plusse sammen sannsynligheter.⁵ Dette medfører også at

$$\mathbb{P}(\Omega \cup \Omega^c) = \mathbb{P}(\Omega) + \mathbb{P}(\Omega^c) = 1 \quad (1.6)$$

$$\implies \mathbb{P}(\emptyset) = 0 \quad (1.7)$$

og

$$\mathbb{P}(A) \leq \mathbb{P}(A) + \mathbb{P}(A^c) = 1 \quad (1.8)$$

$$\implies \mathbb{P}(A) \leq 1 \quad (1.9)$$

Vi kan være interessert i om én eller flere hendelser inntreffer. Dette er ikke helt trivielt siden hendelsene kan overlappe og vi vil ikke telle hvert utfall mer enn én gang.⁶

$$\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B) \quad (1.10)$$

Kan vise at

$$A \subset B \implies \mathbb{P}(B) = \mathbb{P}(B \cap A) + \mathbb{P}(B \cap A^c) = \mathbb{P}(A) + \mathbb{P}(B \cap A^c) \quad (1.11)$$

$$\implies \mathbb{P}(A) \leq \mathbb{P}(B) \quad (1.12)$$

⁵Dette kommer vi tilbake til under lov om total sannsynlighet.

⁶Vet ikke helt hvordan jeg kan utlede det resultatet formelt...

1.2.1 Betinget sannsynlighet og uavhengighet

Vi vil oppdatere våre sannsynlighetsberegninger når vi får mer informasjon om utfallet. Den nye informasjonen gjør det mulig å utelukke noen utfall og avgrense oss til å betrakte en delmengde $B \subset \Omega$ som det nye utfallsrommet. Vi kan likevel finne betingede sannsynligheter med utgangspunkt i sannsynlighetsfunksjonen $\mathbb{P}(\cdot)$ assosiert med det opprinnelige utfallsrommet,

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}, \quad (1.13)$$

der vi kan merke oss at $\mathbb{P}(\cdot|B)$ er et fullverdig probability measure på en sub- σ -algebra av \mathcal{F} siden det tilfredstiller aksiomene. Vi skal senere se at definisjonen over gir sammenhengen mellom simultanfordeling og betinget fordeling til to tilfeldige variabler ettersom realiserte verdier av variablene implisitt avgrenser delmengder av utfallsrommet Ω . Vi kan også merke at det den betingede fordelingen er en skalering av simultanfordeling.

Vi sier at to hendelser er uavhengige dersom

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B) \quad (1.14)$$

som impliserer $\mathbb{P}(A|B) = \mathbb{P}(A)$ og $\mathbb{P}(B|A) = \mathbb{P}(B)$. Det medfører at informasjon om hvorvidt den ene hendelsen har inntruffet ikke gir oss noe informasjon om sannsynligheten for den andre hendelsen. Vi kan generalisere uavhengighet til en samling av N mengder, A_1, A_2, \dots, A_N , ved å kreve at $\mathbb{P}(\cap_{j \in R} A_j) = \prod_{j \in R} \mathbb{P}(A_j)$ for alle $R \subset I = \{1, 2, \dots, N\}$. Det er altså ikke tilstrekkelig at de er parvis uavhengige.

Det kan også være praktisk å betinge av informasjon selv om vi er interessert i en ubetinget hendelse. Dette gjør at vi kan dele problem inn i enklere delproblem og finne løsningen som en vektet sum. Begynner med å observere at

$$\mathbb{P}(A) = \mathbb{P}(A \cap \Omega) = \mathbb{P}(A \cap (B \cup B^c)) = \mathbb{P}((A \cap B) \cup (A \cap B^c)) \quad (1.15)$$

$$= \mathbb{P}(A \cap B) + \mathbb{P}(A \cap B^c) \quad (1.16)$$

$$= \mathbb{P}(A|B)\mathbb{P}(B) + \mathbb{P}(A|B^c)(1 - \mathbb{P}(B)) \quad (1.17)$$

Mer generelt vil en mengde av delmengder $\{B_m\}_{m \geq 1}$ utgjøre en partisjonering av Ω hvis

1. $\cup_m B_m = \Omega$

2. $B_j \cap B_k = \emptyset$ hvis $j \neq k$

Vi kan da finne $\mathbb{P}(A) = \sum_m \mathbb{P}(A|B_m)\mathbb{P}(B_m)$. Dette resultatet er kjent som loven om total sannsynlighet. Vi kan representere fremgangsmåten grafisk med et tre.

Gitt at hendelsen A inntreffer kan vi også være interessert å finne sannsynlighet for

utfallet er i de ulike mengdene av partisjoneringen, for eksempel

$$\mathbb{P}(B_j|A) = \frac{\mathbb{P}(A|B_j)\mathbb{P}(B_j)}{\sum_m \mathbb{P}(A|B_m)\mathbb{P}(B_m)} \quad (1.18)$$

som er kjent som bayes regel.

I praksis er det litt tungvint å jobbe direkte med delmengder av utfallsrommet. Vi liker bedre å jobbe med tall. Vi vil derfor finne en representasjon som gjør det enklere å få svar på de spørsmål vi er interessert i. I praksis er det enklere å jobbe med tall siden de har naturlig rangering, mål på avstand og vi kan blant annet bruke resultat fra kalkulus. Dette motiverer tilfeldige variabler.

1.3 Tilfeldige variabler

På tross av navnet er en tilfeldig variabel verken tilfeldig eller variabel. Det er en deterministisk funksjon x som mapper fra utfallsrommet til talllinjen, $x : \Omega \rightarrow \mathbb{R}$. Vi kan illustrere bruken av tilfeldige variabler med et eksempel. Anta at vi ser på en uendelig coin-flips og la mynt være 1 og krone være 0. Da har vi uendelig antall utfall der hvert utfall er en uendelig følge. Vi er interessert i hvor mange kast det tar før det blir en mynt.

$$\Omega = \{(a_1, a_2, \dots) | a_n \in \{0, 1\}, \forall n \in \mathbb{N}\} \quad (1.19)$$

$$x(\omega) = \min\{n \in \mathbb{N} | a_n = 1\} \quad (1.20)$$

Merk at vi da - for hver realisering i utfallsrommet - bare får det tallet som sier antall kast før første mynt i stedet for hele den uendelige følgen. Denne transformasjonen medfører et tap av informasjon, men vi får den informasjonen vi trenger. Det er et generelt poeng at for å løse problem må vi finne en egnet representasjon av informasjon.

Det faktum at tilfeldige variabler bare er en deterministisk funksjon og all action skjer i probability space er skjult av notasjonelle konvensjoner. Når vi skriver $\{x = 1\}$ så refererer vi implisitt til delmengden av utfallsrommet $\{\omega \in \Omega | x(\omega) = 1\}$. Det betyr at når vi snakker om sannsynligheten til en tilfeldig variabel P_x så er det egentlig \mathbb{P} som jobber under the hood.

$$P_x(1) = P(\{x = 1\}) = \mathbb{P}(\{\omega \in \Omega | x(\omega) = 1\}) = \mathbb{P}(A) \quad (1.21)$$

En annen konvensjon er at sammenligninger mellom tilfeldige variabler blir gjort punktvis i utfallsrommet

$$x = y \iff x(\omega) = y(\omega), \forall \omega \in \Omega \quad (1.22)$$

Det er altså ikke tilstrekkelig at de har samme fordeling. Eksempel: la $X \sim U(0, 1)$ slik

at $F(x) = x$ når $x \in (0, 1)$, og la $Y = g(X) = 1 - X$. Har da at $F(y) = P(1 - X \leq y) = P(X \geq 1 - y) = 1 - P(X < 1 - y) = 1 - (1 - y) = y$. Dette medfører at $X \stackrel{d}{=} Y$, men $X \neq Y$.

Kan også nevne at binære tilfeldige variabler er mye brukt siden vi ofte er interessert i om et eller annet inntreffer eller ikke

$$\mathbb{I}_A(\omega) = \mathbb{I}\{\omega \in A\} \quad (1.23)$$

1.4 Univariat fordeling

Vi har sett at vi kan definere en tilfeldig variabel x på et sannsynlighetsrom $(\Omega, \mathcal{F}, \mathbb{P})$ som gjør at vi for hver $B \in \mathcal{B}(\mathbb{R})$ kan tallfeste $P(x \in B) = \mathbb{P}\{\omega \in \Omega | x(\omega) \in B\}$. Dette er litt omstendelig. Vi kan også bare omdefinere utfallsrommet slik at $\Omega = \mathbb{R}$. Den tilfeldige variabelen gjør denne transformasjonen eksplisitt. Alternativt kan vi abstrahere vekk fra transformasjonen og jobbe direkte med sannsynlighetsrommet $(\mathbb{R}, \mathcal{B}(\mathbb{R}), P)$. Fordelingen P er *probability measure* der $\Omega = \mathbb{R}$ og $\mathcal{F} = \mathcal{B}(\mathbb{R})$. Vi sier at P er *supported* av S hvis $P(S) = 1$.

Fordelingen P er i likhet med andre probability measures \mathbb{P} en funksjon som mapper mengder til $[0, 1]$. Det er veldig fleksibelt og generelt, men det er litt vanskelig å karakterisere funksjonen P . Det er enklere å jobbe med funksjoner som som mapper tall. Dermed er det veldig greit at er en-til-en korrespondanse mellom P og en kumulativ fordelingsfunksjon F der

$$F(s) = P(x \leq s) = P((-\infty, s]), s \in \mathbb{R} \quad (1.24)$$

Denne funksjonen oppfyller en del egenskaper

1. $\lim_{s \rightarrow \infty} F(s) = 1$
2. $\lim_{s \rightarrow -\infty} F(s) = 0$
3. $b > a \implies F(a) \leq F(b)$
4. Den er høyre-kontinuerlig, $\lim_{s \rightarrow s^+} F(s) := F(s^+) = F(s)$

Merk at definisjonsmengden er hele tallinjen uavhengig av om fordelingen er supported av mindre delmengde. De fleste fordelinger er enten diskret eller absolutt kontinuerlige. Fordelingen er diskret hvis den er støttet av en tellbar mengde, altså at det eksisterer en mengde $\{s_j\}_{j \geq 1}$ der $P(\{s_j\}_{j \geq 1}) = 1$. Sannsynlighetsmengden på et gitt element s_j i mengden er $p_j := P(\{s_j\})$ og følgen $\{p_j\}_{j \geq 1}$ utgjør en *pmf*. Hvis fordelingen derimot er absolutt kontinuerlig kan den representeres med en tetthet f som er en ikke-negativ

funksjon på \mathbb{R} som integrerer til 1, der

$$P(B) = \int_B f(s)ds, \quad \forall B \in \mathcal{B}(\mathbb{R}) \quad (1.25)$$

Det er poeng at med lebesgue integral trenger vi ikke alltid skille mellom diskret og absolutt kontinuerlig fordeling siden vi kan integrere begge. Dette er en fordel når vi utvikler teori. På en annen side er distinksjonen vesentlig når vi anvender teori.

$$f_x(s) = F'_x(s) \quad (1.26)$$

$$p_x(s) = F(s) - F(s^-) \quad (1.27)$$

Tilfeldige variabler er funksjoner som transformerer vilkårlige utfallsrom til tallinjen som er enklere å jobbe med. Vi har sett at vi kan velge hvorvidt vi vil være eksplisitt om denne transformasjonen. Hvis vi velger å være eksplisitte kan vi betegne fordelingen P som fordelingen til x , der

$$P(B) = \mathbb{P}(\{x \in B\}) \quad (1.28)$$

For et gitt probability space så vil hver tilfeldig variabel definere en fordeling. Tilsvarende vil det for hver fordeling være mulig å finne en tilfeldig variabel som har denne fordelingen. Vi kan bruke notasjonen $\mathcal{L}(x)$ for å betegne fordelingen til x . Merk at når vi snakker om fordelingen til en tilfeldig variabel så eksisterer det alltid et underliggende sannsynlighetsrom.

1.4.1 Momenter

Vi har lyst på sammendragsmål som beskriver egenskap til funksjon. Forventningsverdi er første moment

$$\mathbb{E}[X] = \int x d(f(x)) = \begin{cases} \int x f(x) dx, & \text{hvis kontinuerlig} \\ \sum x f(x), & \text{hvis diskret} \end{cases} \quad (1.29)$$

Det første integralet er noe lebesgue integral som i prinsippet kan evalueres, men jeg bruker det bare for enhetlig notasjon. Forventningsverdi gir et vektet gjennomsnitt av utfallene til tilfeldig variabel og er et mål på sentraltendensen. Det er forholdsvis enkelt å finne forventningsverdi til transformasjoner av X dersom vi kjenner fordelingen til denne, fordi

$$\mathbb{E}[Y] = \mathbb{E}[g(X)] = \sum_x g(x)p_X(x) \quad (1.30)$$

For å få et mål på spredningen kan vi definere en ny variabel som angir avvik fra forventningsverdi og se på hvor stor størrelsen på dette avviket er i gjennomsnitt.

$$\mathbb{E}[Y^2] = \mathbb{E}[(X - \mu)^2] = \mathbb{E}[X^2] - \mu^2 \quad (1.31)$$

Merk generelt at forventningsverdien til en transformasjon av X ikke tilsvarer transformasjonen evaluert i forventningsverdien, altså

$$\mathbb{E}f(X) \neq f(\mathbb{E}X) \quad (1.32)$$

Forsikringsselskap tjener penger fordi $\mathbb{E}f(X) < f(\mathbb{E}X)$ er lavere når f er konkav. Folk er derfor villig til å betale for å redusere variasjon i X . Kan knytte dette til *Jensens ulikhet*.

1.4.2 Kvantiler

Vi har sett at vi kan karakterisere fordelinger med den kumulative fordelingen

$$F_x(s) = \begin{cases} \int_{-\infty}^s f_x(t)dt, & \text{hvis absolutt kontinuert} \\ \sum_{j:x_j \leq s} p_j, & \text{hvis diskret} \end{cases} \quad (1.33)$$

Litt usikker på hvilken notasjon jeg vil bruke. Tenker at det er greit å spesifisere hvilken variabel vi betrakter fordelingen til slik at vi ikke må bruke så mange ulike bokstaver til å betegne funksjonene. Tror også jeg foretrekker å betegne tilfeldige variabler med stor bokstav. Må avklare dette senere. Uansett, vi kan ofte være interessert i å finne ζ der $F_X(\zeta) = \tau$. Det finner vi ved å evaluere den inverse av cdf i τ . For å håndtere tilfellet der cdf ikke er strengt voksende kan vi definere

$$F_X^{-1}(\tau) = \inf\{s : F_X(s) \leq \tau\} \quad (1.34)$$

Vi får blant annet bruk for kvantilfunksjonen når vi vil finne kritisk verdi i hypotesetester. Anta at vi har en standardnormalfordelt testobservator Z . Vi vil finne et (sentrert) intervall (l, u) der $P(l \leq Z \leq u) = 1 - \alpha$. Vi utnytter at fordeling er symmetrisk slik at $F_Z(-z) = 1 - F_Z(z)$. Dette medfører at

$$F_{|Z|}(s) = P(-s \leq Z \leq s) = F_Z(s) - F_Z(-s) = F_Z(s) - (1 - F_Z(s)) = 2F_Z(s) - 1 \quad (1.35)$$

La $F_{|Z|} := F$. Vi vil finne kritisk verdi c der

$$c = F^{-1}(1 - \alpha/2) \quad (1.36)$$

$$P(|Z| \leq c) = 2F(c) - 1 = 2F[F^{-1}(1 - \alpha/2)] - 1 = 1 - \alpha \quad (1.37)$$

Vi betegner ofte $c := z_{\alpha/2} := \Psi^{-1}(1 - \alpha/2)$

1.5 Multivariat fordeling

Vi kan definere flere funksjoner $X_1(\cdot), \dots, X_N(\cdot)$ på det samme sannsynlighetsrommet. Hvert utfall ω kan være en ganske omfattende representasjon av den tilstanden som blir realisert slik at det ikke trenger å være noe deterministisk sammenheng mellom de ulike transformerte utfallene. Til sammen utgjør de en tilfeldig vektor,

$$X : \Omega \rightarrow \mathbb{R}^N \quad (1.38)$$

$$: \omega \mapsto \begin{bmatrix} X_1(\omega) \\ \vdots \\ X_N(\omega) \end{bmatrix} \quad (1.39)$$

Vi kan forholdsvis enkelt generalisere konseptene om fordeling og support fra én dimensjon til flere dimensjoner. Samtidig er det nye konsepter knyttet til fordeling betinget av informasjon... Jeg begynner med å utlede for bivariate transformasjoner og utvider deretter til N dimensjoner.

1.5.1 Simultanfordeling

Vi kan definere flere tilfeldige variabler på samme utfallsrom Ω . Anta for eksempel at vi kaster et kronestykket to ganger slik at $\Omega = \{HH, HT, TH, TT\}$ og la X være antall heads på første kast og Y være antall heads totalt. Vi kan ordne disse tilfeldige variablene i en tuple slik hver realisering utgjør et punkt i \mathbb{R}^2 , f.eks: $(X(HT), Y(HT)) = (1, 1)$. På samme måte som med én variabel har denne vektoren en fordeling $\mathcal{L}(X, Y) = P$ som angir sannsynlighet til delmengder av \mathbb{R}^2 . Vi kan finne

$$P((X, Y) \in A) = \begin{cases} \sum_{(x,y) \in A} p(x, y) \\ \int \int_A f(x, y) dx dy \end{cases} \quad (1.40)$$

der $p(\cdot)$ og $f(\cdot)$ er henholdsvis den simultane punktsannsynligheten og tettheten. Tetthetsfunksjonen må tilfredstille $f(x, y) \geq 0 \forall (x, y) \in S(X, Y)$ og at $\int \int_{S(X, Y)} f(x, y) dx dy = 1$ (og analogt for diskret). Det er en rett fram generalisering. Utfordringen er å bestemme grenseverdiene i integrasjonen for at den skal samsvare med $A \subset S(X, Y)$ dersom dette ikke bare er en enkel rektangel.

Vi kan også ha fordelinger der vektingen av punktene (tetthet) bare avhenger av én

variabel, men der supporten avhenger av begge variablene. Eksempel:

$$g(y) = e^{-y}, \quad S(x, y) = \{(x, y) : 0 < x < y < \infty\} \quad (1.41)$$

For å gjøre det mer eksplisitt at dette er en simultanfordeling kan vi skrive det på formen

$$f(x, y) = e^{-y} I\{(x, y) \in S(x, y)\} \quad (1.42)$$

som opplagt også avhenger av x siden funksjonen tar verdi 0 når $(x, y) \notin S(x, y)$. Vi kan vise at dette er en gyldig pdf ved å integrere over mengden $S(x, y)$.

$$\int_0^\infty \int_0^y e^{-y} dx dy = \int_0^\infty y e^{-y} dy = 1 \quad (1.43)$$

Kumulativ fordeling

Denne kan beskrives med en kumulativ fordeling F som tar vektor som argument og der

$$F(x, y) = P(\{X \leq x\} \cap \{Y \leq y\}) \quad (1.44)$$

$$= \mathbb{P}\{\omega \in \Omega : X(\omega) < x \wedge Y(\omega) < y\} \quad (1.45)$$

Denne delmengden tilsvarer et rektangel med øvre høyre hjørne i (x, y) . Hvis fordelingen er kontinuerlig er det en simultanfordeling $f(x, y)$ der

$$F(x, y) = \int_{-\infty}^x \int_{-\infty}^y f(w_1, w_2) dw_1 dw_2 \quad (1.46)$$

Hvis fordelingen derimot er diskret er sammenheng mellom simultan og kumulativ gitt ved

$$F(x, y) = \sum_{w_1 < x} \sum_{w_2 < y} p(w_1, w_2) \quad (1.47)$$

Marginal fordeling

Vi kan utlede de marginale fordelingene ved å observere at

$$P(X \in A) = \int_A \int_{\mathbb{R}} f(x, y) dy dx \quad (1.48)$$

eller..

$$P(X = x) = P(\{X = x\} \cap \{Y < \infty\}) = \mathbb{P}(\{X = x\} \cap \Omega) \quad (1.49)$$

$$= \sum_y p(x, y) \quad (1.50)$$

Hm. Integrerer ut de andre variablene. Det er poeng at vi kan få ut all mulig informasjon fra simultanfordelingen. Det er også poeng at vi ikke nødvendigvis kan rekonstruere simultanfordelingen fra informasjon om de marginale fordelingene fordi disse ikke inneholder informasjonen om hvordan de ulike variablene er relatert til hverandre.

1.5.2 Betinget fordeling

Vi har nå definert simultanfordeling og sett hvordan vi kan få ut igjen marginal fordeling fra dette. Det er veldig interessant å betrakte fordelingen til én av variablene gitt verdien av de andre... Det kan også

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \implies p_{X|Y}(x, y) = \frac{p_{X,Y}(x, y)}{p_Y(y)} \quad (1.51)$$

Den betingede fordelingen er en fullverdig fordeling. Vi kan selvfølgelig finne dens momenter.

$$E[Y|X = x] = \int y f(y|x) dy \quad (1.52)$$

$$V[Y|X = x] = \int y^2 f(y|x) dy - \left(\int y f(y|x) dy \right)^2 \quad (1.53)$$

For gitt $X = x$ er dette et uttrykk som i prinsippet kan evalueres og gir oss et tall. Før X er realisert er den betingede forventningen en tilfeldig variabel der $\mathbb{E}[Y|X = x] = g(X)$. I tillegg til betinget forventning har vi også betinget varians kan dekomponere den samlede variansen i ..

$$\mathbb{V}(Y) = \mathbb{E}[\mathbb{V}(Y|X)] + \mathbb{V}[\mathbb{E}(Y|X)] \quad (1.54)$$

Bruk av betinget fordeling i statistisk modellering

Veldig viktig i statistisk modellering. Har én output og mange input. Har i teorien en egen betinget fordeling for hver verdi av input, men vanskelig å jobbe med. Kan bruke sentraltendens av betingede fordelinger som sammendragsmål og se på hvordan sentraltendens endrer når vi endrer input.

Det er også praktisk dersom vi modellerer en univariat fordeling.

1. Partisjonere utfallsrom ut fra verdi til andre variablene
2. Modellere sannsynlighet for hver av delmengdene
3. Modellere den betingede sannsynligheten for utfallet vi er interessert i på hver av delmengdene

4. Aggregere betingede sannsynligheter med vekting ut fra sannsynlighet for delmengdene

Dette er eksempel på anvendelse av lov om total sannsynlighet der vi bruker tilfeldige variabler. For å se koblingen kan vi betrakte

$$\mathbb{P}(A) = \mathbb{P}(A|B_1)\mathbb{P}(B_1) + \mathbb{P}(A|B_2)\mathbb{P}(B_2) \quad (1.55)$$

og

$$p(y) = P(y|x_1)p(x_1) + P(y|x_2)p(x_2) \quad (1.56)$$

$$= \sum_x P(y|X=x)p(x) \quad (1.57)$$

1.5.3 Uavhengighet

Vi har sett på uavhengighet av hendelser. Vi kan utvide til uavhengighet mellom tilfeldige variabler. Disse er uavhengige hvis simultanfordelingen tilsvarer produktet av de marginale fordelingene,

$$f(x, y) = f(x)f(y) \implies f(x|y) = \frac{f(x, y)}{f(y)} = f(x) \quad (1.58)$$

Uavhengige variabler har egenskaper som gjør de enkle å jobbe med. Blant annet er forventningsverdien av produktet lik produktet av forventningsverdier,

$$E[XY] = \int \int xyf(x, y)dxdy \quad (1.59)$$

$$= \int yf(y) \left[\int xf(x)dx \right] dy \quad (1.60)$$

$$= \int yf(y)E[X]dy \quad (1.61)$$

$$= E[X]E[Y] \quad (1.62)$$

1.5.4 Generalisering til N dimensjoner

Vi kan utvide til flerdimensjonale fordelinger der $P : \mathcal{B}(\mathbb{R}^N) \rightarrow [0, 1]$, $\mathbf{x} : \Omega \rightarrow \mathbb{R}^N$ og

$$F_{\mathbf{x}} : \mathbb{R}^N \rightarrow [0, 1] \quad (1.63)$$

$$\mathbf{s} \mapsto P(\times^N(-\infty, s)) \quad (1.64)$$

som karakteriserer fordelingen med den simultane kumulative fordelingen. For absolutt kontinuerlig fordelte variabler er det også en simultan tetthetsfunksjon $p(\cdot)$ som oppfyller

egenskapen

$$P(B) = \int_B p(\mathbf{s}) d\mathbf{s} \quad (1.65)$$

$$= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \mathbb{I}_B(s_1, \dots, s_N) p(s_1, \dots, s_N) ds_1 \dots ds_N \quad (1.66)$$

Den marginale fordelingen til variabel n i vektoren og dens marginale cdf er gitt ved

$$P_n(B) = P(\mathbb{R} \times \cdots \times \mathbb{R} \times B \times \mathbb{R} \times \cdots \times \mathbb{R}) \quad (1.67)$$

$$F_n(s) = P_n((-\infty, s)) \quad (1.68)$$

Den simultane kumulative fordelingen karakteriserer hele fordelingen. Med utgangspunkt i denne kan vi finne marginale kumulative fordelinger. Vi kan også finne betingede fordelinger ved å skalere simultanfordelinger med marginale. Det er derimot som oftest

Betinget fordeling

Har simultanfordeling til (x_1, \dots, x_N) . Kan finne betinget fordeling til delvektor med K komponenter for gitte realiseringer av de resterende $N - K$ komponentene. Det gir en funksjon $f : \mathbb{R}^K \rightarrow \mathbb{R}$. Eksempler:

$$f(s_2, \dots, s_N | x_1 = s_1) = \frac{f(\mathbf{s})}{f(s_1)} \quad (1.69)$$

$$f(s_1 | s_2, \dots, s_N) = \frac{f(\mathbf{s})}{f(s_2, \dots, s_N)} \quad (1.70)$$

Uavhengighet

Momenter

Forventningsverdi til $\mathbf{x} = (x_1, \dots, x_N)$ er bare en vektor der hver komponent er forventningsverdi til den tilhørende tilfeldige variabelen. Variansen er

$$var(\mathbf{x}) := \Sigma = \mathbb{E}[(\mathbf{x} - \mu)(\mathbf{x} - \mu)'] \quad (1.71)$$

$$= \mathbb{E}\mathbf{x}\mathbf{x}' - \mu\mu' \quad (1.72)$$

der elementene $\Sigma_{ij} = cov(x_i, x_j)$. Denne matrisen er positiv semi-definit og har egenskaper som følger av dette.

$$var(A\mathbf{x}) = A\Sigma A' \quad (1.73)$$

Kovariansen til to variabler er forventningsverdien til produktet av avviket fra forventningsverdi til hver av variablene

$$\text{cov}(X, Y) := \sigma_{X,Y} = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] \quad (1.74)$$

Kan få litt intuitjon av at kovarians er positiv hvis det er tendens til at positive avvik skjer samtidig i begge variablene. Fanger opp om det er lineær sammenheng. Kan skalere ved å dele på produktet av standardavvikene og få korrelasjonskoeffisient som er begrenset av $(-1,1)$

$$\rho := \frac{\sigma_{X,Y}}{\sigma_X \sigma_Y} \quad (1.75)$$

1.5.5 Tilfeldig utvalg

Vi kan definere et utvalg som N observasjoner av en variabel med en gitt fordeling, (X_1, \dots, X_N) der X_n har samme pdf $f(\cdot)$ for alle n . I praksis antar vi at det er en parametrisert fordeling $f(\cdot; \theta)$. Utvalget er tilfeldig dersom observasjonene er uavhengige.⁷ Da er simultanfordeling $f(x_1, \dots, x_N) = \prod_n f(x_n)$. I statistikk er utfordringen at vi ikke kjenner $f(\cdot)$ og vi forsøker å lære egenskaper ved denne fra utvalget...

Spørsmål om utvalg

Vi tenker gjerne at hver observasjon i utvalget er en person, men rammeverket er mer generelt. Kan være helt andre ting vi observerer. Med utgangspunkt i simultanfordelingen kan vi svare på spørsmål om sannsynlighet for ulike hendelser. Uavhengighet og identiske fordelinger gjør det ganske enkelt å jobbe med. Vi kan for eksempel finne sannsynlighet for at alle observasjonene tar høyere verdi enn a med

$$P(\{X_1 > a\} \cap \dots \{X_N > a\}) = \prod_n (1 - F(a)) \quad (1.76)$$

Sammendragsmål

I praksis jobber vi gjerne med sammendragsmål på utvalget. Vi konstruerer disse med såkalte *statistikker*(?) som er funksjoner fra utvalget $T : \mathbb{R}^N \rightarrow \mathbb{R}^k$.⁸ Vi bruker gjerne

⁷Tilfeldige utvalg kan betegnes som representative fordi store talls lov sikrer at andelen av observasjoner med gitte egenskaper i utvalget konvergerer mot andelen i populasjonen. Skal formalisere dette i kapittel om inferens.

⁸I praksis mapper til \mathbb{R} . Eneste restriksjon er at de ikke kan avhenge av parametre.

gjennomsnittet i utvalget som sammendragsmål på sentraltendensen.

$$\bar{X} := \frac{1}{N} \sum_n X_n \quad (1.77)$$

$$\mathbb{E}[\bar{x}] = \mathbb{E}\left[\frac{1}{N}(X_1 + \dots + X_N)\right] = \mu \quad (1.78)$$

$$\mathbb{V}[\bar{x}] = \mathbb{V}\left[\frac{1}{N}(X_1 + \dots + X_N)\right] = \frac{\sigma^2}{N} \quad (1.79)$$

der vi kun for siste resultat trenger uavhengighet. Så lenge observasjonene er fra samme fordeling kan vi i gjennomsnitt estimere sentraltendensen, men vi må ta hensyn til eventuell avhengighet for å vurdere presisjonen til målet på sentraltendensen...

Vi vil også ha et sammendragsmål på spredningen i utvalget,

$$s^2 = \frac{1}{N-1} \sum (X_n - \bar{X})^2 \quad (1.80)$$

$$\mathbb{E}[s^2] = \sigma^2 \quad (1.81)$$

der beviset er litt vanskelig.

1.5.6 Hierarki

En variabel har én fordeling, men det kan være lurt å bygge denne opp stegvis og bruke ulike fordelinger underveis. Dette gjør det enklere å motivere valg av endelig fordeling (som blir vår modell) og denne fordelingen blir enklere å tolke siden den gjerne avhenger av parametre i fordelingene over i hierarkiet. Jeg tenker at hierarkiske modeller blir viktige i statistisk modeller og at det kan være lurt å bruke bayesiansk statistikk siden det gjør det enklere å propagere usikkerheten assosiert med estimatene av de ulike parametrene i den hierarkiske strukturen. Jeg tar kort gjennomgang her siden med gitte parameterverdier så er det er rent sannsynlighetsproblem.

La oss illustrere med eksempel. Vi vil modellere hvor mange barn de samlet vil få. Vi kunne forsøkt å finne en parametrisert fordeling på dette direkte. Alternativt kan vi modellere antallet kvinner som i det hele tatt får barn og antallet barn dersom de gjør det..

$$Y \sim \text{bern}(p) \quad (1.82)$$

$$X|Y \sim \text{poisson}(\lambda) \quad (1.83)$$

$$f_X(x) = \sum_y f(x|y)f(y)dy \quad (1.84)$$

Kan forsøke å finne analytisk uttrykk for $f_X(x)$, men det kan være litt komplisert. Det i hvert fall enkelt å finne uttrykk for sentraltendens. Substituerer inn uttrykket for $f_X(\cdot)$

inn i $E[\cdot]$

$$E[X] = \int x f_X(x) dx \quad (1.85)$$

$$= \int x \int f(x|y) f(y) dy dx \quad (1.86)$$

$$= \int \left[x \int f(x|y) dx \right] f(y) dy \quad (1.87)$$

$$= \int E[X|y] f(y) dy \quad (1.88)$$

$$= E[E[X|Y]] \quad (1.89)$$

Mixture model

Tror dette er special case av hierarkisk fordeling der parameter i fordeling selv har fordeling.. hmhmm.

1.6 Transformasjon av tilfeldig variabel

Tilfeldige variabler er funksjoner som mapper elementer av utfallsrommet til tallinjen, $X : \Omega \rightarrow \mathbb{R}$. Vi har sett hvordan jeg kan utlede verdimengden til transformasjonen samt sannsynligheten til delmengder av denne verdimengden fra den eksplisitte formen til $X(\cdot)$ samt sannsynlighetsrommet $(\Omega, \mathcal{F}, \mathbb{P})$ den er definert på. Vi har også sett at dette inducerer et nytt sannsynlighetsrom $(\mathbb{R}, \mathbb{B}(\mathbb{R}), P)$ slik at vi kan abstrahere vekk fra den eksplisitte transformasjonen samt det opprinnelige, underliggende sannsynlighetsrommet.⁹ Vi skal nå se at vi kan definere en ny tilfeldig variabel $Y : \mathbb{R} \rightarrow \mathbb{R}$ på dette induuerte sannsynlighetsrommet og finne et eksplisitt uttrykk for fordelingen $f_Y(y)$ med utgangspunkt i $f_X(x)$ samt $Y(\cdot)$. Selv om Y egentlig er en funksjon vil jeg betegne den transformerte tilfeldige variabelen som $Y = g(X)$..

1.6.1 Fordeling til transformert variabel

Når vi gjør transformasjoner er det viktig å beholde oversikt over verdimengdene til transformasjonene, som samsvarer med *supporten* til fordeling. La $S(X) = \{x : f_X(x) > 0\}$ og $S(Y) = \{y : y = g(x), \exists x \in S(X)\}$ være support til henholdsvis X og Y . Vi kan finne fordelingen til Y hvis vi for alle delmengder $A \subset S(Y)$ kan finne en korresponderende delmengde $B = \{x : g(x) \in A\} \subset S(X)$. Vi kan gjøre dette mer operativt ved å definere en invers mapping $g^{-1} : \mathbb{B}(S(Y)) \rightarrow \mathbb{B}(S(X))$, som mapper delmengder av supporten og

⁹Det induerte sannsynlighetsrommet inneholder all informasjon om fordelingen til den tilfeldige variabelen, men kan medføre informasjonstap i forhold til opprinnelig sannsynlighetsrom.

er derfor definert selv om transformasjonen $g(\cdot)$ ikke er strengt monoton. Har da at

$$P(g(X) \in A) = P(x \in g^{-1}(A)) \quad (1.90)$$

Vi kan nå finne uttrykk for den kumulative fordeingen til $F_Y(t)$ ved å betrakte $A = (-\infty, t)$.

$$F_Y(t) = \int_{g^{-1}((-\infty, t))} f_X(x) dx \quad (1.91)$$

Det er enklere å jobbe med strengt monotone transformasjoner siden det er enklere å beskrive grenseverdiene i integrasjonen. Hvis strengt voksende reduserer problemet til

$$F_Y(t) = \int_{-\infty}^{g^{-1}(t)} f_X(x) dx = F_X(g^{-1}(t)) \quad (1.92)$$

Ved å betrakte $A = (-\infty, t)$ kan vi finne $f_Y(t)$ til transformasjon av kontinuerlig fordelt X .

1.6.2 Kontinuerlig

$$F_Y(y) = P(Y \leq y) = P(g(X) \leq y) = P(X \leq g^{-1}(y)) = F_X(g^{-1}(y)) \quad (1.93)$$

$$f_Y(y) = \frac{\partial}{\partial y} F_X(g^{-1}(y)) = f_X(g^{-1}(y)) \frac{dx}{dy} \quad (1.94)$$

der $\frac{dx}{dy} = \frac{d}{dx} g^{-1}(y)$. Hvis funksjonen derimot er monotont avtagende må vi snu ulikhets-tegnet,

$$F_Y(y) = P(Y \leq y) = P(g(X) \leq y) = P(X > g^{-1}(y)) = 1 - F_X(g^{-1}(y)) \quad (1.95)$$

$$f_Y(y) = \frac{\partial}{\partial y} [1 - F_X(g^{-1}(y))] = -f_X(g^{-1}(y)) \frac{dx}{dy} \quad (1.96)$$

Merk at $\frac{dx}{dy} < 0$ slik at uttrykket er positivt. Vi kan få felles uttrykk for begge tilfellene med

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{dx}{dy} \right| \quad (1.97)$$

Dette gir oss en formel med størrelser vi må plugge inn. Vi trenger inverse og den deriverte av den inverse. Eksempel: $f(x) = x^2/9, x \in (0, 3), y = g(x) = x^3$

$$g^{-1}(y) = y^{\frac{1}{3}}, \quad \frac{dx}{dy} = \frac{1}{3}y^{-\frac{2}{3}} \quad (1.98)$$

$$f_Y(y) = \frac{\left(y^{\frac{1}{3}}\right)^2}{9} \frac{1}{3}y^{-\frac{2}{3}} = \frac{1}{27}, \quad y \in (0, 27) \quad (1.99)$$

$$(1.100)$$

Husk at vi kan finne $f(x)$ ved å transformere tilbake, så ingen unnskyldning for å gjøre feil!

1.6.3 Diskret

For diskret variabler trenger vi ikke gå gjennom kumulativ fordelingsfunksjon..

$$p_Y(y) = P(Y = y) = P(g(X) = y) = P(X = g^{-1}(y)) = p_X(g^{-1}(y)) \quad (1.101)$$

La oss ta antall kast før første kron som eksempel. Ha da $p_X(x) = 0.5^x$. Anta nå at vi allerede vet at første kastet aldri gir treff. Vi vil derfor finne fordeling til $Y = g(X) = X+1$. Begynner med å finne $g^{-1}(c) = c - 1$. Det medfører da at

$$p_Y(y) = p_X(g^{-1}(y)) = 0.5^{g^{-1}(y)} = 0.5^{y-1} \quad (1.102)$$

Intuisjonen er at vi evaluerer pmf til X i det tallet som mapper til y -verdien vi er interessert i, altså i $g^{-1}(y)$.

1.6.4 Bivariat transformasjon

Vil utvide til å se på fordeling til en bivariat transformasjon $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, der¹⁰

$$g(X_1, X_2) = [g_1(X_1, X_2), g_2(X_1, X_2)] = (Y_1, Y_2) \quad (1.103)$$

Vi kan fortsatt i prinsippet finne

$$P((Y_1, Y_2) \in A) = P((X_1, X_2) \in \{(x_1, x_2) : [g_1(x_1, x_2), g_2(x_1, x_2)] \in A\}) \quad (1.104)$$

Antar at det er en invers transformasjon J der

$$J(y_1, y_2) = [g_1^{-1}(y_1, y_2), g_2^{-1}(y_1, y_2)] \quad (1.105)$$

¹⁰Sikkert god idé å sette det opp som kolonnevektorer..

Skalerer med determinant av jacobitil denne inverse transformasjonen og finner...

$$f_{Y_1, Y_2}(y_1, y_2) = f_{X_1, X_2}([g_1^{-1}(y_1, y_2), g_2^{-1}(y_1, y_2)]|J|) \quad (1.106)$$

Eksempel

Har to variabler $X_1, X_2 \sim NID(\mu, \sigma^2) \implies f(x_1, x_2) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x_1^2}{2}\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x_2^2}{2}\right)$. Transformasjonen $g(\cdot) := [g_1(\cdot), g_2(\cdot)] : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ er gitt ved

$$Y_1 = g_1(X_1, X_2) = X_1 + X_2 \quad (1.107)$$

$$Y_2 = g_2(X_1, X_2) = X_1 - X_2 \quad (1.108)$$

For å finne den inverse transformasjonen må jeg løse ligningsystemet med hensyn på X_1 og X_2 som gir

$$X_1 = h_1(Y_1, Y_2) = (Y_1 + Y_2)/2 \quad (1.109)$$

$$X_2 = h_2(Y_1, Y_2) = (Y_2 - Y_1)/2 \quad (1.110)$$

... trenger bare determinant til jacobitil transformasjon $h(\cdot) := [h_1(\cdot), h_2(\cdot)]$.

1.7 Momentgenererende funksjoner

Den momentgenererende funksjonen til en tilfeldig variabel er gitt ved

$$M(t) = \mathbb{E}[e^{tX}], \quad \text{for } t \in (-h, h) \quad (1.111)$$

Funksjonen er definert dersom uttrykket over konvergerer for verdier av t på et åpent intervall om 0. Kan ta et raskt eksempel på utledning av *mgf*. La $p(x) = \frac{1}{6} \left(\frac{5}{6}\right)^{x-1}$ være sannsynlighet for antall kast det tar å få sekser på terningen.

$$\mathbb{E}[e^{tX}] = \sum_{x=1}^{\infty} \frac{1}{6} \left(\frac{5}{6}\right)^{x-1} e^{tx} \quad (1.112)$$

$$= \frac{1}{6} e^t \sum_{x=1}^{\infty} \left(\frac{5e^t}{6}\right)^{x-1} \quad (1.113)$$

$$= \frac{e^t}{6} \left(1 - \frac{5e^t}{6}\right)^{-1} \quad (1.114)$$

som konvergerer dersom

$$\frac{5e^t}{6} < 1 \implies \log(t) < \log\left(\frac{6}{5}\right) \quad (1.115)$$

For fordelingene med definert *mgf* er det én-til-én korrespondanse mellom fordeling og *mgf*.¹¹ Med andre ord så vil

$$M_x(t) = M_y(t), \quad \forall t \in (-h, h) \quad (1.116)$$

$$\implies F_x(s) = F_y(s), \quad \forall s \in Z \quad (1.117)$$

Vi kan bruke denne alternative representasjonen til å beskrive egenskaper til fordelingen. Merk at¹²

$$M'(t) = \frac{\partial}{\partial t} \mathbb{E}[e^{tX}] \quad (1.118)$$

$$\begin{cases} \int_Z \frac{\partial}{\partial t} e^{tx} f(x) dx = \int_Z x e^{tx} f(x) dx \\ \sum_{x \in Z} \frac{\partial}{\partial t} e^{tx} p(x) = \sum_{x \in Z} x e^{tx} p(x) \end{cases} \quad (1.119)$$

som medfører at $M'(t)|_{t=0} = \mathbb{E}[x]$ og mer generelt at $M^{(k)}(t)|_{t=0} = \mathbb{E}[x^k]$. Det kan i praksis være litt vanskelig å finne *mgf* siden det er litt vanskelig å integrere og sånn. Men gitt at vi har en *mgf* kan vi også finne *mgf* til affine transformasjon $Y = a + bX$,

$$m_Y(t) = E[e^{(a+bX)t}] = E[e^{at} e^{bXt}] = e^{at} M_x(bt) \quad (1.120)$$

Dersom to variabler er uavhengige kan vi finne *mgf* til deres sum, $Z = X + Y$,

$$m_Z(t) = E[e^{(X+Y)t}] = E[e^{tX}] E[e^{tY}] = m_X(t) m_Y(t) \quad (1.121)$$

Flervariabel

Kan merke oss kort at *mgf* til tilfeldig variabel $\mathbf{x} = [x_1, \dots, x_N]$ er

$$M(\mathbf{t}) = \mathbb{E}[\exp\{\mathbf{t}'\mathbf{x}\}] = \int \dots \int \exp\{\mathbf{t}'\mathbf{x}\} f(\mathbf{x}) dx_1 \dots dx_N \quad (1.122)$$

og det følger at vi kan bruke dette til å finne marginal *mgf* til x_n ved evaluere den i $\mathbf{t} = (t_1, \dots, t_n, \dots, t_N) = (0, \dots, t_n, \dots, 0)$.

¹¹Det har noe sammenheng med laplace-transformasjon av funksjon. Det finnes også en såkalt karakteristisk funksjon som er en generalisering som er definert for alle fordelinger som har sammenheng med fourier-trasnformasjon.

¹²Har brukt at vi kan flytte derivasjonstegn inn og ut av sum og intergral. Er noen tekniske betingelser som må være oppfylt for at dette skal være gyldig operasjon (i betydning at uttrykkene har samme løsningsmengde).

1.8 Projeksjon av tilfeldige variabler

Det er mulig å konstruere vektorrom som består av andre objekter enn tradisjonelle vektorer (tupler av reelle tall). Skal nå konstruere vektorrom og utvikle analoge resultat for ortogonal projeksjon og minimering av avstand mellom objekter i rommet.

Vi betrakter en mengde av tilfeldige variabler. Det hadde vært veldig greit å ha et mål på avstand mellom objektene i mengden. Et naturlig mål er $RMSE(x, y) = \sqrt{E[(x - y)^2]}$. Hvis vi definerer indre produkt mellom objekter som $\langle x, y \rangle = E[xy]$ får vi det analoge resultat at $\|x\| = \sqrt{\langle x, x \rangle}$ og vi kan definere $x \perp y = 0 \iff E[xy] = 0$. På tilsvarende måte utgjør delmengder underrom dersom de er lukket for skalering og addisjon. La $\mathbb{1}_\Omega := \mathbb{1}$ være tilfeldig variabel som tar verdi 1 med sannsynlighet 1. Et eksempel på underrom er da

$$\text{span}(X) = \{\alpha \mathbb{1} + \beta x : \alpha, \beta \in \mathbb{R}\} \quad (1.123)$$

På samme måte som i tradisjonelle vektorrom har ortonormale basiser gode egenskaper. En mengde U er ortonormal basis for S hvis

$$E[u_j u_k] = \mathbb{1}\{j = k\} \quad \text{og} \quad \text{span}\{u_1, \dots, u_k\} = S \quad (1.124)$$

Jeg vil derfor lage en ortonormal basis for $S = \text{span}\{\mathbb{1}, x\}$. I tråd med gram schmidt algoritmen trekker jeg fra komponenten til x som går i retning til den konstante tilfeldige variabelen, nemlig $E[x] := \mu$. Da sitter jeg igjen med residualen $x - \mu$. Lengden til denne variabelen, i henhold til normen som ble definert over, er $\sigma_x := \sqrt{E[(x - \mu)^2]}$. Bruker dette til å skalere og har ortonormal basis

$$U = \left\{ \mathbb{1}, \frac{x - \mu}{\sigma_x} \right\}, \quad \text{span}(U) = S \quad (1.125)$$

Ser da at ved å standardisere variablene i data så er det utvalgsanalog til å lage ortogonal basis for de tilfeldige variablene. Gitt definisjonene over er det ortogonale projeksjonsteoremet helt analogt så gidder ikke gjenta dette. Det er veldig kult at vi kan overføre teori om vektor til R.V siden vi får veldig mye gratis og det gir oss geometrisk intuisjon. Tenk for eksempel at jeg vil projekte y på S som jeg har gitt en ortonormal basis. La

den være u_1, u_2 . Det følger da at

$$\mathbf{P}y = \langle y, u_1 \rangle u_1 + \langle y, u_2 \rangle u_2 \quad (1.126)$$

$$= E[y] + E\left[\frac{x - \mu}{\sigma_x} y\right] \frac{x - \mu}{\sigma_x} \quad (1.127)$$

$$= E[y] + \frac{\text{cov}(x, y)}{\text{var}(x)}(x - \mu) \quad (1.128)$$

$$= (E[y] - \beta\mu) + \beta x \quad (1.129)$$

merk at $E[(x - \mu)y] = \text{cov}(x, y)$ fordi $E[x - \mu] = 0$. Dette er populasjonsversjonen av den bivariate lineære regresjonslinjen der jeg brukte ortonormal basis. Skal nå generalisere til å finne $\mathbf{P}y = \text{proj}_S y = \mathbf{x}'\beta$. Merk at jeg alltid kan slenge inn en konstant tilfeldig variabel i mengden siden vi alltid "observerer denne uavhengig av hvilke data vi har.

$$E[(y - \mathbf{x}'\beta)\mathbf{x}] = 0 \quad (1.130)$$

$$\implies \beta = E[\mathbf{x}\mathbf{x}']E[\mathbf{x}y] \quad (1.131)$$

veldig nice. TODO: knytte til prediksjon, informasjonmengde og utvide til arbitrære funksjoner ved å definere underrom $L_2(X)$.

1.8.1 Projektering i L_2

Begynner med å definere L_2 som mengden av tilfeldige variabler med definert andre moment, $L_2 = \{x | \mathbb{E}(x^2) < \infty\}$. Vi kan definere et indre produkt mellom elementer i mengden, $\langle x, y \rangle = \mathbb{E}(xy)$ der $\mathbb{E}(xy) = 0 \implies y \perp x$. Normen til elementer i mengden er $\|x\| = \sqrt{\langle x, x \rangle}$. Merk at $\|x - y\| = \sqrt{\mathbb{E}[(x - y)^2]}$ som tilsvarer *root mean square error*. Delmengder utgjør et lineært underrom hvis de er lukket under skalering og addisjon, $x, y \in S \implies \alpha x + \beta y \in S, \forall \alpha, \beta \in \mathbb{R}$. For en mengde tilfeldige variabler $\mathbf{x} = (x_1, \dots, x_K)$ vil $\text{span}(\mathbf{x}) = \{\mathbf{x}'\mathbf{b} | \mathbf{b} \in \mathbb{R}^K\} \equiv S(\mathbf{x})$ være et underrom som består av alle lineære kombinasjoner av (x_1, \dots, x_K) . Mer generelt kan vi betrakte arbitrære deterministiske funksjoner av \mathbf{x} . En variabel z er *\mathbf{x} -measurable* hvis det finnes en funksjon $h : \mathbb{R}^K \rightarrow \mathbb{R}$, der $h(\mathbf{x}) = z \in L_2$. Mengden av disse variablene utgjør også et underrom som vi kan kalle $L_2(\mathbf{x})$.

Anta nå at vi for $y \in L_2$ og et underrom $S \subset L_2$ vil vi finne elementet i S som minimerer avstanden til y . Analogt til tradisjonelle vektorrom kan L_2 dekomponeres i et underrom S og dets ortogonale komplement S^\perp , og der $S^\perp = \{z | \langle z, x \rangle = 0, \forall x \in S\}$. Alle element i L_2 kan da skrives som en sum av et element i hvert underrom, $y = \hat{y} + \hat{u}$, der $\hat{y} \in S$ og $\hat{u} \in S^\perp$. Denne \hat{y} er løsningen på minimeringsproblemet $\arg \min_{z \in S} \|y - z\|$. Vi vet at løsningen eksisterer, er unik og har egenskapen $\langle (y - \hat{y}), x \rangle = 0, \forall x \in S$. Det eksisterer også en lineær transformasjon P slik at $P(y) = \hat{y}$. Denne transformasjonen

utfører den ortogonale projeksjonen av y på underrommet S .

Betrakt tilfellet der $S = S(\mathbf{x})$. Løsningen er da gitt ved $\hat{y} = \mathbf{x}'\mathbf{b}^*$, der $\langle x_k, y - \mathbf{x}'\mathbf{b}^* \rangle = 0, k = 1, \dots, K \iff \mathbb{E}(\mathbf{x}(y - \mathbf{x}'\mathbf{b}^*)) = \mathbf{0} \iff \mathbf{b}^* = \mathbb{E}(\mathbf{x}\mathbf{x}')^{-1}\mathbb{E}(\mathbf{x}y)$, hvis $\mathbb{E}(\mathbf{x}\mathbf{x}')$ er inverterbar.

Så langt er det helt analogt til tradisjonelle vektorrom, men vi kan nå knytte ortogonal projeksjon til forventning. Vi kan definere

$$\mathbb{E}(y|\mathbf{x}) \equiv \arg \min_{z \in L_2(\mathbf{x})} \|y - z\| \quad (1.132)$$

Den betingede forventningsverdien av y gitt \mathbf{x} er den \mathbf{x} -*measurable* variabelen som minimerer avstanden til y . Hvis vi definerer en konstant tilfeldig variabel $\mathbb{1}_\Omega \equiv \mathbb{I}\{\omega \in \Omega\}$, så vil ubetinget forventning være gitt ved

$$\mathbb{E}(y) \equiv \arg \min_{z \in L_2(\mathbb{1}_\Omega)} \|y - z\| \quad (1.133)$$

Den konstante tilfeldige variabelen som minimerer avstanden til y .

1.8.2 Payoff

Okay, hva er gevinsten ved å tenke på forventning som ortogonale projeksjoner?

1. $y = \mathbb{E}(y|\mathbf{x}) + u \implies \mathbb{E}[g(\mathbf{x})u] = 0, \forall g$ følger direkte av at $\mathbb{E}(y|\mathbf{x})$ er ortogonal projeksjon av y på $L_2(\mathbf{x})$. Alle andre variabler $g(\mathbf{x})$ ligger i underrommet $L_2(\mathbf{x})$ og u er derfor ortogonal med disse. Siden $\mathbb{E}(u) = 0$ er u ukorrelert med alle deterministiske funksjoner av \mathbf{x} .
2. $y = \mathbf{x}'\mathbf{b}^* + u \implies \mathbb{E}[\mathbf{x}'\gamma u] = 0, \forall \gamma \in \mathbb{R}^K$. Av samme argument er feilledd fra projeksjon på $S(\mathbf{x})$ ortogonal på alle lineære kombinasjoner av \mathbf{x} .
3. Har generelt at hvis underrommene $V_2 \subset V_1$ så vil det være ekvivalent om man projekterer y direkte på V_2 eller først projekterer på V_1 og deretter på V_2 . Ettersom $S(\mathbf{x}) \subset L_2(\mathbf{x})$ vil da $\mathbf{x}'\mathbf{b}^*$ være beste lineære tilnærming til $\mathbb{E}(y|\mathbf{x})$.
4. Kan bruke et tilsvarende argument for å utlede *law of iterated expectations*. Merk at $L_2(\mathbb{1}_\Omega) \subset L_2(\mathbf{x})$ slik at $\mathbb{E}[\mathbb{E}(y|\mathbf{x})] = \mathbb{E}(y)$.
5. Kan bruke ortogonal dekomponering av underrom, $S(\mathbf{x}) \subset L_2(\mathbf{x}) \subset L_2$, og pythagoras' setning til å dekomponere den forventede feilen ved å predikere y med $\mathbf{x}'\mathbf{b}$. La $\mathbb{E}(y|\mathbf{x}) \equiv f^*(\mathbf{x})$

$$\|y - \mathbf{x}'\mathbf{b}\|^2 = \|y - f^*(\mathbf{x})\|^2 + \|f^*(\mathbf{x}) - \mathbf{x}'\mathbf{b}^*\|^2 + \|\mathbf{x}'\mathbf{b}^* - \mathbf{x}'\mathbf{b}\|^2 \quad (1.134)$$

$$\mathbb{E}[(y - \mathbf{x}'\mathbf{b})^2] = \mathbb{E}[(y - f^*(\mathbf{x}))^2] + \mathbb{E}[(f^*(\mathbf{x}) - \mathbf{x}'\mathbf{b}^*)^2] + \mathbb{E}[(\mathbf{x}'\mathbf{b}^* - \mathbf{x}'\mathbf{b})^2] \quad (1.135)$$

6. Kan tilsvarende dekomponere *mean square error* mellom en estimator og parameter i varians og kvadrert bias ved å projekte på $L_2(\mathbb{1}_\Omega)$

$$\|\hat{\theta} - \theta\|^2 = \|\hat{\theta} - \mathbb{E}(\hat{\theta})\|^2 + \|\mathbb{E}(\hat{\theta}) - \theta\|^2 \quad (1.136)$$

7. Linearitet til forventning følger av linearitet til ortogonale projeksjoner.

Kapittel 2

Stokastiske prosesser

En stokastiske prosesser på \mathbb{R}^K er en samling tilfeldige vektorer $(\mathbf{x}_t)_{t \in T}$ definert på samme sannsynlighetsrom $(\Omega, \mathcal{F}, \mathbb{P})$. Vi betegner utfallsrommet T som indeksemengden og utfallsrommet til de tilfeldige vektorene er såkalt *state space*. Det at de lever i samme sannsynlighetsrom medfører at de har en simultanfordeling. Denne simultanfordelingen kan ofte være ganske komplisert fordi

$$f(x_1, \dots, x_N) = f(x_1)f(x_2|x_1)f(x_3|x_1, x_2) \cdots f(x_N|x_1, \dots, x_{N-1}) \quad (2.1)$$

$$= \prod f(x_n|\text{historie}_n) \quad (2.2)$$

der $\text{historie}_n = (x_1, \dots, x_{n-1})$. Stort sett har vi unngått dette problemet ved å anta uavhengighet, men vi skal utvikle vektøy for å håndtere avhengighet mellom observasjoner. Først litt motivasjon

1. Det åpner for avhengighet i sampling. Dette er spesielt relevant når vi observerer samme enhet på flere tidspunkt (panel/tidsserie), men det kan også være avhengighet i kryssseksjon. Tror det kan være avhengighet for observasjoner i samme gruppe (klasse, geografisk område, mm.), men veldig usikker på dette. Tror en tilnærming til dette er å skalere standardfeil med cluster.
2. Vi kan bruke det til å studere egenskap til estimator. For utvalg med gitt N kan vi utlede fordeling til estimator under gitt antagelser, men vi kan også være interessert i hvordan denne fordelingen endres når vi endrer N . Vi kan da betrakte det som en stokastisk følge der vi i hvert ledd observerer én ny realisering. Ved å undersøke egenskaper til følgen når N går mot uendelig kan vi utlede egenskaper under svakere antagelser og bruke dette som tilnærming for store utvalg.
3. Vi kan bruke det til å modellere systemer. Verden er dynamisk, ikke statisk. Mer om dette senere.

2.1 Asymptotisk teori

Jeg betrakter en følge av tilfeldige variabler $(X_1, X_2, \dots, X_n, \dots) = (X_n)_{n \in \mathbb{N}}$. Vi kan tenke at egenskapene til variablene avhenger av plassering i følgen, altså av n . Vi kan være interessert i egenskap for en gitt n eller for alle n som er større enn en gitt verdi. Vi sier at $\lim X_n$ har en egenskap dersom det eksisterer en N slik at alle de tilfeldige variablene X_n der $n \geq N$ har egenskapen... Det ble litt upresist. Uansett, grenseverdi har en formell definisjon. Følgen med tall konvergerer til et grenseverdien (som er et tall) dersom vi alltid kan finne en N der alle $x_n, n \geq N$ er i nabolaget til tallet, uansett hvor lite vi gjør nabolaget. Skal nå se på konvergens av følge med tilfeldige variabler. Grenseverdien er nå en tilfeldig variabel... vil på en måte lage et nabolag rundt denne tilfeldige variabelen og se om alle variablene i følgen for $n \geq N$ ligger i dette nabolaget. Er ikke helt opplagt hvordan vi måler avstand mellom tilfeldige variabler og konstruerer dette nabolaget, så med tilfeldige varabler har vi litt ulike former for konvergens

1. Konvergens i sannsynlighet: $X_N \xrightarrow{p} X \iff \lim_{N \rightarrow \infty} P(\|Z_n - Z\| > \epsilon) = 0, \forall \epsilon > 0$
2. Konvergens *almost surely*: $X_N \xrightarrow{a.s.} X \iff P(\lim_{N \rightarrow \infty} \|Z_n - Z\| > \epsilon) = 0, \forall \epsilon > 0$
3. Konvergens i fordeling: $X_N \xrightarrow{d} X \iff \lim_{N \rightarrow \infty} F_N(t) = F(t)$ for alle t der $F(\cdot)$ er kontinuerlig.
4. Konvergens i *mean square* (L2): $X_N \xrightarrow{m.s.} X \iff \lim_{N \rightarrow \infty} \mathbb{E}(X_n - X)^2 = 0$

Konvergens i sannsynlighet impliserer konvergens i fordeling. Tror jeg bruker konvergens i m.s. til å bevise konvergens i sannsynlighet. Merk at en konstant c bare er special case av R.V. X der $\mathbb{P}(X = c) = 1$. Det er et veldig fint resultat at konvergens er bevart av kontinuerlig transformasjoner $g(\cdot)$, slik at $Z_n \xrightarrow{p} Z \implies g(Z_n) \xrightarrow{p} g(Z)$. Noen regneregler:

Slutskys teorem

$X_n \xrightarrow{d} X$ og $A_n \xrightarrow{p} a$ impliserer at

1. $X_n + A_n \xrightarrow{d} X + a$
2. $X_n A_n \xrightarrow{d} Xa$

Kan generaliseres til tilfeldige vektorer og matriser slik at $A_N \mathbf{x}_N \xrightarrow{d} A\mathbf{x}$. Medfører at dersom $\mathbf{x} \sim N(\mathbf{0}, \Sigma)$ så er $A\mathbf{x} \sim N(\mathbf{0}, A\Sigma A')$

Kontinuerlig mapping teorem

Har et veldig greit resultat som sier at for uansett hvilken konvergensmåte vi har, så vil konvergens medfører at den samme konvergens holder for kontinuerlige transforma-

sjon av både grenseverdien og følgen,

$$X_n \xrightarrow{p} X \implies g(X_n) \xrightarrow{p} g(X) \quad (2.3)$$

2.1.1 Store talls lov

Store talls lov som sier at utvalgsmoment fra *iid* prosess konvergerer i sannsynlighet til populasjonsmoment. Så lenge første moment er definert så vil

$$\mathbb{E}_{P_N}[X_n] = \frac{1}{N} \sum X_n \xrightarrow{p} \mathbb{E}[X] \quad (2.4)$$

For å knytte dette til definisjonene over vil jeg bruke litt annen notasjon. Vi har en følge $(X_n)_{n \in \mathbb{N}}$ som er *iid*. Med utgangspunkt i denne følgen kan vi konstruere en ny følge med gjennomsnitt av de N første variablene, $(\bar{X}_N)_{N \in \mathbb{N}}$. Vi kan nå vise at denne nye følgen konvergerer i sannsynlighet til $\mathbb{E}[X]$ som vi kan betrakte som en tilfeldig variabel der $\mathbb{P}(\{\mathbb{E}[X]\}) = 1$.

Bevis

For å bevise store talls lov vil jeg først utlede Markovs og Chebyshevs ulikheter. La X være ikke-negativ tilfeldig variabel og $g(\cdot)$ er transformasjon som flytter tyngde nedover

$$Y = g(X) = \begin{cases} a, & \text{hvis } X \geq a \\ 0, & \text{ellers} \end{cases} \quad (2.5)$$

Det følger da at

$$\mathbb{E}[X] \geq aP(X \geq a) \quad (2.6)$$

$$\implies P(X \geq a) \leq \frac{\mathbb{E}[X]}{a} \quad (2.7)$$

som er Markovs ulikhet. Kan utlede Chebyshevs ulikhet som spesialtilfelle der $X = (\bar{X}_N - \mu)^2$ og $a = \epsilon^2$.

$$P[(\bar{X}_N - \mu)^2 \geq \epsilon^2] \leq \frac{\mathbb{E}[(\bar{X}_N - \mu)^2]}{\epsilon^2} = \frac{\mathbb{V}[\bar{X}_N]}{\epsilon^2} \quad (2.8)$$

$$\implies P(|\bar{X}_N - \mu| \geq \epsilon) \leq \frac{\sigma^2}{N\epsilon^2} \rightarrow 0 \quad (2.9)$$

når $N \rightarrow \infty$. Okay, nå var jo beviset ferdig uten at jeg definerte Chebyshevs ulikhet... Kunne også tatt beviset fra *mean square*. Uansett, LLN er et veldig fundamentalt resultat fordi det enkelt kan utvides. Det gjelder også for utvalgsmoment til kontinuerlige

funksjoner av X .

$$\frac{1}{N} \sum g(X_n) \xrightarrow{p} \mathbb{E}[g(X)]. \quad (2.10)$$

Det kan også brukes til å bevise at relativ andel i utvalg konvergerer til sannsynlighet

$$\frac{1}{N} \sum I\{X_n \in B\} \xrightarrow{p} \mathbb{E}[I\{X \in B\}] = P(B) \quad (2.11)$$

2.1.2 Sentralgrenseteoremet

Store talls lov sier at i uendelig store utvalg er hele tyngden av fordelingen til utvalgsmomentet konsentrert på de populasjonsmomentet. Det har sammenheng med at empirisk fordeling konvergerer til teoretisk fordeling. Det er et viktig teoretisk resultat, men i praksis har vi aldri uendelig store utvalg så vi vil også vite noe om hvor raskt tyngden til utvalgsfordeling konvergerer: mer presist ønsker vi å angi sannsynlighet for avvik mellom utvalgsgjennomsnitt og forventningsverdi. Her kommer sentralgrenseteoremet oss til unnsetning. Det sier at så lenge $\{X_n\}$ er *iid* og $\mathbb{E}[X^2] < \infty$

$$\sum_{n=1}^N \frac{X_n - \mu}{\sigma^2} \xrightarrow{d} N(0, 1) \quad (2.12)$$

$$\frac{(X_1 + \dots + X_N) - N\mu}{\sqrt{N}\sigma^2} \xrightarrow{d} N(0, 1) \quad (2.13)$$

$$\frac{\bar{X}_N - \mu}{\sqrt{\mathbb{V}[\bar{X}_N]}} \xrightarrow{d} N(0, 1) \quad (2.14)$$

$$\sqrt{N}(\bar{X}_N - \mu) \xrightarrow{d} N(0, \sigma^2) \quad (2.15)$$

Ser at dette er veldig nyttig siden differansen av utvalgsmoment og populasjonsmoment (oppskalert med rate of convergence) blir normalfordelt uavhengig av underliggende fordeling.

Det er veldig praktisk at vi kan bruke Slutsky til å finne asymptotisk fordeling til lineære transformasjoner av normalfordelte estimatorer og at de fortsatt er normalfordelte. Kan bruke delta-metode til å generalisere dette resultatet til differensierbare funksjoner..

2.1.3 Delta-metoden

Husker at $A_N \xrightarrow{p} A$ og $\mathbf{x}_N \xrightarrow{d} \mathbf{x} \implies A_N \mathbf{x}_N \xrightarrow{d} A\mathbf{x}$. Dette er veldig nice siden CLT kan gi oss at $\mathbf{x} \sim N(\mathbf{0}, \Sigma)$. Jeg vet allerede hvordan jeg finner fordeling til lineær transformasjon. Nå skal jeg også kunne finne fordeling til transformasjoner som er lokalt

lineære (ie. kontinuerlige..).

$$Y_N \xrightarrow{d} N(\mu, \frac{\sigma^2}{N}) \implies g(Y_N) \xrightarrow{d} N\left(g(\mu), (g'(\mu))^2 \frac{\sigma^2}{N}\right) \quad (2.16)$$

2.1.4 LNN og CLT med flere variabler

2.1.5 LLN og CLT med avhengighet mellom observasjoner

Hvis prosessen er *iid* så er $\mathcal{L}(\mathbf{x}_t) = P, \forall t$ og simultanfordelingen er produkt av marginal. Slike følger er veldig greie å jobbe med siden vi kan bruke LLN og CLT, men det er litt restriktivt siden det kan være litt persistens i størrelse over tid. Hvis vi vil jobbe med tidsserier er det derfor relevant å få finne en større klasse av stokastiske prosesser som *nesten* er *iid* og har de samme gode egenskapene. Det viser seg at LLN holder for prosesser som er stasjonære og ergodiske. Stasjonærhet medfører at simultanfordeling til del-tupler av simultanfordelingen ikke endres av å forskyves. Altså:

$$\mathcal{L}(\mathbf{x}_{t1}, \dots, \mathbf{x}_{tk}) = \mathcal{L}(\mathbf{x}_{1t+m}, \dots, \mathbf{x}_{tk+m}) \quad (2.17)$$

Det finnes litt ulike og kompliserte definisjoner av ergodisitet. Jeg velger å si at en stasjonær prosess er ergodisk dersom LLN holder. Dette flytter målposten til å si noe om tilstrekkelige betingelser for ergodisitet. En uformell definisjon på ergodisitet er at gjennomsnitt av observasjon over tid omtrent samsvarer med gjennomsnitt på et gitt tidspunkt.

Det er et poeng at vi kan få CLT som kun krever at prosess er martingale difference sequence. Jeg skal forsøke å utlede litt mer formelt hva dette er for noe. Jeg vet ikke hvor relevant det er, men jeg kjører på og håper at det blir litt payoff senere. Trenger først å introdusere noen konsepter.

2.2 Markov-kjeder

Markovkjeder har diskret tilstandsrom og sannsynlighet for ulike tilstander i en periode kun avhenger av tilstand i perioden før,

$$\mathbb{P}(X_n = x | X_1, \dots, X_{n-1}) = \mathbb{P}(X_n = x | X_{n-1}) \quad (2.18)$$

dette forenkler uttrykket for simultanfordelingen

$$f(x_1, \dots, x_N) = f(x_1)f(x_2|x_1)f(x_3|x_1, x_2) \cdots f(x_N|x_1, \dots, x_{N-1}) \quad (2.19)$$

$$= f(x_1) \prod_{n=2}^N f(x_n|x_{n-1}) \quad (2.20)$$

Markovkjeder gir et rammeverk med nok struktur til at det er mulig å utlede teoretiske resultat og samtidig har det tilstrekkelig fleksibilitet til å beskrive systemer i virkeligheten.

2.2.1 Overgangssannsynlighet

De sentrale størrelsene i en markovkjede er tilstandsrommet og sannsynlighet for å bevege seg mellom tilstander. Hvis markovkjeden er overgangssannsynlighetene uavhengig av n slik at vi kan definere

$$p_{ij} := \mathbb{P}(x_{n+1} = j | x_n = i). \quad (2.21)$$

og organisere disse i en matrise P . Sentrale spørsmål: $p_{ij}(n)$ og konvergens.. mer om dette senere.

2.3 Annet

En filtrasjon er en økende følge av informasjonsmengder, altså $(\mathcal{F}_t) = (\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_t, \dots)$ der $\mathcal{F}_t \subset \mathcal{F}_{t+1} \forall t$. Merk at en informasjonsmengde bare er en mengde av tilfeldige variabler. Filtrasjonen som er generert av (x_t) er

$$\mathcal{F}_0 = \emptyset \text{ og } \mathcal{F}_t = \{x_1, \dots, x_t\} \text{ for alle } t \geq 1 \quad (2.22)$$

En stokastisk prosess (z_t) er *adapted* til en filtrasjon (\mathcal{F}_t) hvis z_t er \mathcal{F}_t -measurable for alle t . Altså at vi kan regne ut z_t fra størrelsene i informasjonsmengden som blir realisert på samme tidspunkt. Prosessen (z_t) er i tillegg en martingale wrt (\mathcal{F}_t) hvis

$$\mathbb{E}[z_{t+1} | \mathcal{F}_t] = z_t \quad (2.23)$$

La differansen $d_t = z_t - z_{t-1}$. Dette definerer en stokastisk prosess (d_t) som er *martingale difference sequence* wrt (\mathcal{F}_t) hvis

$$\mathbb{E}[d_{t+1} | \mathcal{F}_t] = 0 \quad (2.24)$$

okay... får håpe det blir noe payoff en dag LOL. Tror jeg kommer til å studere dette til høsten så forhåpentligvis blir det en søt syntese. Tidsserier og sånn.. vi får se. Men er glad for at jeg fikk sannsynlighet på litt tryggere grunn. Litt usikker på hva jeg skal gjøre fremover... får sjå.

Kapittel 3

Noen kjente fordelinger

3.1 Normalfordeling

Normalfordelingen dukker opp ofte som følge av sentralgrenseteoremet. Summen av tilfeldige avvik har en tendens til å jevne seg ut slik at tyngden blir konsentrert rundt sentraltendensen. Det er likevel litt fascinerende at det asymptotisk har eksakt normalfordeling uavhengig av den underliggende fordelingen til hvert avvik. Det er en familie med fordelinger som er unikt bestemt av to parametre. Tetthet til standardnormal er:

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}z^2\right\}. \quad (3.1)$$

Vi har ingen *closed form* for cdf $\Phi(z) = \int_{-\infty}^z \phi(z)dz$. Vi kan konstruere andre medlemmer av familien gjennom *affine* transformasjon.

$$x = \mu + \sigma z \sim N(\mu, \sigma^2) \quad (3.2)$$

$$z = \frac{x - \mu}{\sigma} \sim N(0, 1) \quad (3.3)$$

Normalfordelingen har gode egenskaper

1. Alle lineære kombinasjoner av normalfordelinger er også normalfordelte.
2. Fravær av korrelasjon impliserer uavhengighet.
3. Betinget forventingsfunksjon (CEF) er lineær.

3.1.1 Normalfordelt utvalg

Har nå sett på egenskaper som holder for alle *iid* utvalg fra fordeling så lenge relevante moment er definert. For å få hele utvalgsfordelingene må vi enten bruke asymptotisk tilnærming for store utvalg eller anta at $f(\cdot)$ er kjent slik at vi i praksis kan utlede analytiske

resultat for vilkårlig N . Den helt klart vanligste antagelsen er at X er normalfordelt. Vi kan da vise noen klassiske resultat om sammendragsmålene over,

1. \bar{X} og s^2 er uavhengige
2. $\bar{X} \sim N(\mu, \frac{\sigma^2}{N})$
3. $\frac{N-1}{\sigma^2} \sim \chi^2(N-1)$

Bevisene for (1) og (3) er vanskelige. Kan ta bevis for (2) med mgf. Merk først at

$$M_{\bar{X}}(t) = \mathbb{E} \left[\exp \left(t \left(\frac{X_1 + \dots + X_N}{N} \right) \right) \right] \quad (3.4)$$

$$= \mathbb{E}[\exp(tX_1/N)] \cdot \dots \cdot \mathbb{E}[\exp(tX_N/N)] \quad (3.5)$$

$$= \Pi_n M_X(t/N) = m_X(t/N)^N \quad (3.6)$$

som innebærer at

$$M_{\bar{X}}(t) = \exp \left(\mu \frac{t}{N} + \frac{\sigma^2(t/N)^2}{2} \right)^N \quad (3.7)$$

$$= \exp \left(\mu t + \frac{(\sigma^2/N)t^2}{2} \right) \quad (3.8)$$

som betyr at $\bar{X} \sim N(\mu, \frac{\sigma^2}{N})$.

3.1.2 Truncated normalfordeling

Generelt er en truncated fordeling betinget av at variabelen tar utfall i et gitt intervall. Vi vil for eksempel finne betinget tetthet til en variabel y gitt at vi vet at $y > c$. Det kan vises at dette bare er ubetinget tetthet til y skalert med sannsynligheten for å være i intervallet slik at det integrerer til 1,

$$f(y|y > c) = \frac{f(y)}{P(y > c)} = \frac{f(y)}{1 - F(c)} \quad (3.9)$$

som gir mening siden

$$\int_c^\infty f(y)dy = 1 - \int_{-\infty}^c f(y)dy = 1 - F(c) \quad (3.10)$$

Dersom y er standardnormalfordelt kan vi utlede enkle uttrykk for hvordan momentene til den avkuttete fordelingen til y avhenger av c ,

$$\mathbb{E}[y|y > c] = \frac{\phi(c)}{1 - \Phi(c)} \quad (3.11)$$

og har også et uttrykk for varians. Disse vil jeg utlede og det vil også være enkelt å utvide andre normalfordelte variabler.

3.1.3 Multivariat normalfordeling

En vektor kan også være normalfordelt. Normalfordeling i én dimensjon (skalar) er bare en special case, så alle resultatene under har analog for skalarer.

$$\mathbf{x} = \mathbf{A}\mathbf{z} \implies V(\mathbf{x}) = E[\mathbf{x}\mathbf{x}'] = \mathbf{A}E[\mathbf{z}\mathbf{z}']\mathbf{A}' = \mathbf{A}\mathbf{A}' \quad (3.12)$$

3.2 Fordelinger assosiert med normalfordeling

I statistisk inferens dukker normalfordeling ofte opp fordi vi har mange uavhengige observasjoner. Vi bruker testobservatorer til å undersøke om det er rimelig at de observerte data kan ha blitt generert av en avgrenset modell. Disse testobservatorene har ofte en fordeling som har sammenheng med normalfordelte størrelser i utvalget... hm.

3.2.1 χ^2 -fordeling.

Den kvadrerte lengden til en normalfordelt vektor,

$$y = \|\mathbf{z}\|^2 = \mathbf{z}'\mathbf{z} = \sum_{m=1}^M z_m^2 \sim \chi^2(M) \quad (3.13)$$

I praksis har vi ofte normalfordelinger uten standardisert varians, $\mathbf{x} \sim N(\mathbf{0}, \Sigma)$ der $\Sigma = \mathbf{A}\mathbf{A}'$.¹ Vi kan likevel utlede χ^2 fordeling fra dette ved å bruke kvadratisk form. I praksis skalerer det ned vektoren før vi tar lengden..

$$\mathbf{x}'\Sigma^{-1}\mathbf{x} = \mathbf{x}'(\mathbf{A}\mathbf{A}')^{-1}\mathbf{x} \quad (3.14)$$

$$= \mathbf{x}'(\mathbf{A}')^{-1}\mathbf{A}^{-1}\mathbf{x} \quad (3.15)$$

$$= (\mathbf{A}^{-1}\mathbf{x})'\mathbf{A}^{-1}\mathbf{x} \quad (3.16)$$

$$= \mathbf{z}'\mathbf{z} \sim \chi^2(M) \quad (3.17)$$

der $\mathbf{A}^{-1}\mathbf{x} = \mathbf{z}$ følger av at lineære kombinasjoner av normal er normal og at

$$V(\mathbf{A}^{-1}\mathbf{x}) = E[\mathbf{A}^{-1}\mathbf{x}(\mathbf{A}^{-1}\mathbf{x})'] \quad (3.18)$$

$$= \mathbf{A}^{-1}E[\mathbf{x}\mathbf{x}'](\mathbf{A}')^{-1} \quad (3.19)$$

$$= \mathbf{A}^{-1}\mathbf{A}\mathbf{A}'(\mathbf{A}')^{-1} \quad (3.20)$$

¹Det eksisterer en slik dekomponering av positiv definit (analog til kvadratroten), kan bruke Cholesky til å finne.

3.2.2 t-fordeling

hm

3.2.3 F-fordeling

hm

3.3 Fordelinger fra bernoulli-prosess

3.3.1 Binomialfordeling

Sannsynlighet for k treff av N uavhengige $Y_n \sim \text{Bernoulli}(P)$ er

$$P_x(k) = \binom{N}{k} p^k (1-p)^{(N-k)} \quad (3.21)$$

Merk at vi kan visualisere utfallene med et tre (graf) siden det deler i to i hvert steg. Forventningsverdi er

$$\mathbb{E}X := \mu_X = \mathbb{E}g(Y_1, \dots, Y_N) = \mathbb{E} \sum Y_n = \sum \mathbb{E}Y_n = np \quad (3.22)$$

Variansen er

$$\mathbb{V}X := \sigma_X^2 = \mathbb{V}g(Y_1, \dots, Y_N) = \sum \mathbb{V}Y_n = np(1-p) \quad (3.23)$$

3.3.2 Negativ binomialfordeling

Sannsynligheten for x antall forsøk før vi oppnår r treff er

$$p(x; r, p) = \binom{x-1}{r-1} p^{r-1} (1-p)^{x-r} p \quad (3.24)$$

Intuisjonen er at vi må ha $r-1$ treff på $x-1$ forsøk og deretter treffe på siste.

3.3.3 Geometrisk fordeling

Dette er bare spesialtilfelle av negativ binomialfordeling der $r = 1$,

$$p(x; p) = p(1-p)^{x-1} \quad (3.25)$$

Formelen blir enklere fordi det bare er én måte vi kan bomme $x-1$ ganger på $x-1$ forsøk. Da bommer vi hver eneste gang inntil vi treffer. Tror jeg vil vise at det er en fordeling uten hukommelse...

3.3.4 Multinomialfordeling

Generalisering av binomialfordeling der det er K kategorier i stedet for bare 2. Kan tenke på det som sannsynlighet for antall baller med ulike farger fra en urne etter N trekk med andeler p_1, p_2, \dots, p_K . Tenker at det er en tilfeldig vektor (X_1, X_2, \dots, X_K) som angir antall i hver kategori og at vi kan beskrive pmf til denne.

3.4 Fordelinger fra poisson-prosess

3.4.1 Poissonfordeling

Poissonfordeling gir oss sannsynlighet for antall treff per tidsenhet,

$$p(x; \lambda) = \frac{\lambda^x}{x!} e^{-\lambda} \quad (3.26)$$

der λ er en intensitetsparameter. Vi kan utvide til å finne antall treff på lengre intervall ved å skalere intensiteten, $\lambda' = \lambda \cdot t$, der t er antallet tidsenheter i det nye intervallet.

Kan motivere som kontinuerlig tilnærming til binomialfordeling når p er lav og N er høy..

For å utlede egenskaper for vi bruk for at e^x har alternativ representasjon som sum av uendelig følge

$$e^x = \sum_{i=0}^{\infty} \frac{x^i}{i!} \quad (3.27)$$

Kan bruke dette til å vise at det er en gyldig pmf,

$$\sum_{x=0}^{\infty} \frac{\lambda^x}{x!} e^{-\lambda} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^{-\lambda} e^{\lambda} = 1 \quad (3.28)$$

3.4.2 Eksponentialfordeling

Lengde mellom treff i poissonfordeling. Bruke til å modellere varighet av noe.

$$f(x; \theta) = \theta e^{-\theta x} \quad (3.29)$$

3.5 Andre fordelinger

3.5.1 Uniformfordeling

Det er en kontinuerlig fordeling med like stor sannsynlighet for utfall i alle delmengder som er like store. I én dimensjon kan vi spesifisere $X \sim U(a, b)$, uniform på intervallet

$[a, b]$. I flere dimensjoner kan vi være litt mer kreative med geometriske objekt, f.eks. disk eller kradrat. Tetthetsfunksjonen vil uansett være en konstant siden den ikke avhenger av hvor i mengden vi er. For å finne denne konstanten i én dimensjon kan vi bruke

$$\int_a^b k dx = k(b - a) = 1 \implies k = \frac{1}{b - a} \quad (3.30)$$

Dette følger også av at vi skal ha et rektangel der ene siden er bredde er k , lengde er $b - a$ og areal skal være 1. Forventningsverdi er

$$\mathbb{E}X = \int_a^b x f(x) dx = \frac{1}{b - a} \int_a^b x dx = \frac{a^2 - b^2}{2(b - a)} = \frac{a + b}{2} \quad (3.31)$$

Kan også merke oss at cdf til enhetsuniform er $F(x) = \int f(x) = \int 1 = x$. Hvis jeg har N uavhengige uniforme variabler og jeg vil finne forventningsverdi til $Y = \max\{X_1, \dots, X_N\}$ kan jeg bruke at

$$F(Y_n = y) = P(X_1 < y, \dots, X_N < y) = \Pi F_X(y) = y^N \quad (3.32)$$

$$\mathbb{E}Y = \int_0^1 y f(y) dy = \int_0^1 N y^{N-1} y dy = \int_0^1 N y^N dy \quad (3.33)$$

$$= \frac{N}{N + 1} \quad (3.34)$$

3.5.2 Gammafordeling

Vi har en funksjon $\Gamma(a)$ som er en slags generalisering av factorial til positive reelle tall,

$$\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt \quad (3.35)$$

Denne funksjonen har mange nyttige egenskaper som det sikkert er mulig å utlede,

- $\Gamma(1) = 1$
- $\Gamma(n) = (n - 1)!$
- $\Gamma(n + 1) = n\Gamma(n)$

Vi kan bruke denne funksjonen til å konstruere en *pdf*,

$$f(t; \alpha) = \frac{t^{\alpha-1} e^{-t}}{\Gamma(\alpha)} \quad (3.36)$$

I praksis bruker vi en annen representasjon siden vi også vil ha parameter som justerer skala til fordelingen. Vi oppnår dette ved å definere en ny variabel $x := t/\beta$ og se på

fordelingen til denne

$$f(x; \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta} \quad (3.37)$$

Dette er en veldig fleksibel parametrisk familie og jeg kan vise at vi kan få andre fordelinger som spesialtilfelle ved å spesifisere verdi på parametre...

3.5.3 Betafordeling

Det er en betafunksjon,

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \quad (3.38)$$

og vi kan bruke dette til å konstruere en *pdf*,

$$f(x; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} \quad (3.39)$$

Denne funksjonen er bare definert på begrenset intervall, $x \in (0, 1)$. Dette er veldig praktisk siden vi kan bruke fordelingen til å beskrive sannsynlighet til parametre med begrenset intervall. Det er også en veldig fleksibel parametrisk familie som kan ha mange ulike former avhengig av parameterverdier...

3.5.4 Hypergeometrisk

Ligner litt på binomialfordeling fordi vi er interessert i sannsynlighet for x antall treff på gitt antall trekk, men relativ andel som er gunstige i populasjonen blir påvirket av observasjonene som blir trukket. Vi kan betrakte det som er urne-modell uten tilbakelegging. Sannsynlighet for x treff på K trekk fra en populasjon der M av N er gunstige er

$$p(x; K, M, N) = \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}} \quad (3.40)$$

Denne funksjonen er litt vanskelig å jobbe med. Foretrekker kontinuerlige funksjoner så vi unngår slitsom kombinatorikk.

Kapittel 4

Inferens

I sannsynlighetsteori tar vi utgangspunkt i et veldefinert stokastisk forsøk som er beskrevet av sannsynlighetsrommet $(\Omega, \mathcal{F}, \mathbb{P})$ og bruker dette til å beregne sannsynlighet for ulike hendelser.¹ Dette forutsetter at fordelingen er kjent. I statistikk er fremgangsmåten omvendt; vi observerer utfall og med utgangspunkt i dette vil vi si noe om egenskapene til fordelingen som genererte utfallene. De realiserte utfallene gir ikke tilstrekkelig informasjon til å entydig bestemme egenskaper ved fordelingen P ; ulike fordelinger kan generere samme observasjoner og hvis vi observerer nye realiseringer fra samme fordeling vil det være tilfeldig variasjon. Det er derfor vesentlig å kvantifisere usikkerheten til våre mål på egenskapene til P .

4.1 Motivasjon

Bedrifter, myndigheter og andre organisasjoner samler inn store mengder av data. Det er i hovedsak tabeller med verdier til observasjoner langs ulike egenskaper. Dette kan betraktes som et råstoff som må bearbeides og analyseres for å skape verdi. Det overordnede målet er å lære fra data. For å gjøre det må vi skape alternative representasjoner av tabellen med tall.² Vi kan grovt skille mellom to typer representasjoner ut fra om de er enkle eller komplekse.

Enkle representasjoner

Dette er representasjoner som kan tolkes av mennesker og bidra til å skape innsikt og gi bedre informasjonsgrunnlag for å fatte beslutninger. Eksempler på dette er tabeller med sammendragsmål, visuelle fremstillinger eller enkle regresjonsmodeller. Vi kan her bruke at tabulær data definerer en (empirisk) fordeling. Det er en utfordring at hele simultanfordelingen er et veldig komplisert objekt og at vi har begrenset med dimensjoner

¹I praksis jobber vi ofte med $\Omega = \mathbb{R}$ slik at vi kan bruke funksjoner $f : \mathbb{R} \rightarrow \mathbb{R}$ til beregne sannsynlighet for ulike hendelser gitt fordelingen $P : \mathcal{B}(\mathbb{R} \rightarrow \mathbb{R})$

²Representasjon brukes her som en deterministisk transformasjon av tabellen til et annet objekt.

til rådighet. Vi kan visualisere univariate fordelinger og betingede fordelinger. Vi kan også visualisere hvordan sentraltendensen i de betingede fordelingene avhenger av verdien vi betinger av. Generelt er vi ofte opptatt av hvordan et utfall avhenger av andre observerte egenskaper. Regresjonsmodeller gir et sammendragsmål på denne sammenhengen.

Komplekse representasjoner

Det kan også være hensiktsmessig å bruke mer komplekse representasjon som er mer fleksible og kan avdekke mer subtile mønstre i data. Disse er ofte ikke så enkle for mennesker å tolke og vil ikke på samme måte bidra til å skape innsikt og forståelse for sammenhengene mellom variablene i datasettet, men de kan blant annet være bedre egnet til å predikere verdien av en utfallsvariabel når vi kun observerer en delmengde av egenskapene til en gitt observasjon. Dette kan betegnes som maskinlæring.

4.1.1 Generalisering

Representasjoner av (den empiriske fordelingen i) utvalget er ikke hele fortellingen. Generelt er vi ikke (bare) interessert i å lære om data som foreligger, men å lære om egenskaper til fordelingen som genererte dataene. Dette kan være fordi vi kun har data om et utvalg av observasjoner fra en veldefinert populasjon av potensielle observasjoner. I så fall vil vi gjerne være interessert i fordelingen i denne populasjonen. Mer generelt kan vi betrakte observasjonene i datasettet som realiserte verdier fra en datagenereringsprosess som også vil generere nye verdier i fremtiden. Vi må derfor også ta hensyn til at egenskapene til den empiriske fordelingen vi observerer kan avvike fra den teoretiske fordelingen som genererte datasettet. Dette kan vi betegne som *stokastisk usikkerhet*. Denne usikkerheten kan vi behandle på systematisk måte under gitte antagelser. Videre vil det være *induktiv usikkerhet* knyttet til om data som foreligger og representasjonene vi bruker faktisk gir oss svar på det vi spør om.

Stokastisk usikkerhet

Dette er usikkerhet i estimering av modellen gitt at den er riktig spesifisert. Denne usikkerheten skyldes at vi har begrenset antall observasjoner slik at vi ikke er helt sikre på egenskapene i prosessen; nye utvalg av samme størrelse ville gitt annen føyning og det er en kvantifiserbar usikkerhet som følge av dette. Vi kan i prinsippet håndtere denne usikkerheten på en konsekvent måte og dette skal jeg beskrive i resten av kapitlet.

Induktiv usikkerhet

Det vil også være mer generell usikkerhet om i hvilken grad det gitte utvalget og modellen gjør at vi kan svare på spørsmål vi er interessert i. Her kreves det kontekst-spesifikk

kunnskap.

1. Hvordan er data generert; er det skjevheter i utvalget? Er det missing values? Målefeil?
2. Er det riktig spesifisering av modell? Her er det jo mye fleksibilitet, slik at andre valg kan føre til andre konklusjoner. Viktig å være transparent og gjøre goodness of fit test på antagelser i modell som påvirker konklusjon. Konklusjoner er mer troverdige dersom de er robuste.³
3. Med prediksjon kan ekstern validitet i prinsippet testes direkte med kryssvalidering, men det er fortsatt vesentlig å forstå hvilke variabler/egenskaper ved observasjon som gjør at hypotesefunksjonen trekker konklusjoner om y for å bygge kredibilitet og sikre at det kan generaliseres til nye setninger. En klassifiserer som skiller mellom ulv og husky ut i fra om det er snø på bildet kan ha gode resultater for et gitt datasett, men ikke nødvendigvis så nyttig for alle setninger.
4. I økonometri forsøker vi å beskrive en egenskap til den sanne prosessen som genererte utfallet y . Dette kan avvike fra simultanfordeling mellom variablene vi observerer. Vi kan tenke at det er en sann, deterministisk prosess $y = f(x_1, \dots, x_K)$ som bestemmer utfallet, men vi observerer bare en delmengde av disse variablene. Utfordring blir da å si noe om den kausale $\frac{\partial f}{\partial x_k}$ gitt den doble utfordringen at vi kun har delmengde av variabler og begrenset antall observasjoner. Teori om inferens kan brukes til å håndtere det siste. For førstnevnte trenger vi forskningsdesign med tilfeldig variasjon x_k .

4.2 Formelt rammeverk

Målet er å lære fra data. Vi vil lære om egenskaper til den prosessen som genererte de observerte verdiene i datasettet. For å formalisere dette slik at vi kan anvende sannsynlighetsteori vil jeg innføre noen begreper og notasjon,

- Vi observerer realiserte verdier $\mathbf{z}_n \in Z = \mathbb{R}^d$, der Z er *utfallsrommet* til observasjonene.
- *Utvalget* består av N observasjoner $Z_D = (\mathbf{z}_1, \dots, \mathbf{z}_N) \in Z_D = \times_{n=1}^N Z = \mathbb{R}^{N \times d}$, der Z_D er *utvalgsrommet*.
- I praksis velger vi ofte å dekomponere observasjonene i $\mathbf{z}_n = (\mathbf{x}_n, y_n)$ og undersøke hvordan (den betingede fordelingen til) *utfallet* y avhenger av *input* \mathbf{x} .⁴

³Robust i betydningen at de ikke er sensitive for antagelser og andre (mer eller mindre) vilkårlig valg.

⁴Merk at det finnes veldig mye ulike terminologi for å betegne disse variablene...

Observasjonene er realiserte verdier fra en fordeling $\mathcal{L}(\mathbf{z}_n) = P_0$.⁵ Denne sanne fordelingen er ukjent. Selv om fordelingen er ukjent vil vi ofte gjøre antagelser om egenskaper den oppfyller. Med andre ord avgrenser vi oss til å betrakte en delmengde av mulige fordelinger $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$. Dette er *modellen* vår.

4.2.1 Modell

Dersom den inneholder den sanne fordelingen, $P_0 \in \mathcal{P}$, er modellen riktig spesifisert. (I praksis ikke siden modell er forenkling, men mange god grunner til å påføre struktur. Kan også være fordel å være agnostisk.. siden konklusjon avhenger av antagelse.. Kanskje flytte diskusjon av modell til annet avsnitt og knytte til bias varians tradeoff..?). Rammeverket er generelt og omfatter både parametriske og ikke-parametriske tilnærminger

- Parametrisk hvis fordeling kan blir indeksert med parametervektor θ med endelig dimensjon.
- Ikke-parametrisk hvis θ trenger uendelig dimensjon for å karakterisere fordeling.
- Semi-parametrisk hvis vi kan dekomponere $\Theta = \Theta_0 \times \Theta_1$, der kun den første rommet har endelig dimensjon og inneholder parameter vi er interessert i.

Valg av struktur

... Motivere med eksempler. Generalisere til delmengder av utfallsrommet der vi har få observasjoner. Bruker informasjon vi har fra før. Kombinere data med prior kunnskap.

Identifiserbarhet

En modell er identifiserbar hvis det eksisterer en injektiv (én-til-én) funksjon $g : \theta \mapsto \mathbb{P}_\theta$. Hvis den ikke er injektiv kan vi ikke vite parameteren selv om vi observerer selve fordelingen som har generert observasjonene. Eksempler:

- $X \sim \text{bern}(g(\theta))$, kan bare lære θ dersom $g : \theta \mapsto \rho$ er injektiv.
- $Y = I\{X > a/2\}$ og $X \sim U(0, a)$. Vi observerer bare Y ; kan vi lære a fra cdf til Y ?
Nei, fordi $P(Y = 1) = 1 - P(X < a/2) = 1 - \int_0^{a/2} \frac{1}{a} dx = 1 - (\frac{a/2}{a}) = 1/2$ for alle a .

4.2.2 Lære egenskaper til fordeling

Generelt kan vi betrakte egenskaper til fordeling P som alle størrelser som vi kan beregne når fordelingen er kjent. Med andre ord, alle $\gamma(P)$, der $\gamma : \mathcal{P} \rightarrow S$. Eksempler på egenskaper er mål på sentraltendens som forventningsverdi og median, samt mål på spredning

⁵Vi kan bruke subskript 0 for å betegne sanne verdier.

som standardavvik. I multivariate fordelinger er vi ofte interessert i å lære den betingede forventningen av én variabel som funksjon av de andre. Vi kan også være interessert i å lære hele fordelingen slik at γ er identitetsfunksjonen, eller vi kan ønske å lære sannsynlighet for spesifikk hendelse. Merk at hvis vi har avgrenset \mathcal{P} til en spesifikk parametrisk klasse så vil det være tilstrekkelig å lære $\gamma(P) = \theta$ siden alle andre egenskaper vil være en deterministisk funksjon av P_θ .

Hvordan kan vi lære om egenskapene når P er ukjent? Vi må ta utgangspunkt i utvalget av realiserte verdier fra P og finne en representasjon som er informativ om den egenskapen vi vil lære om. Denne representasjonen er verdien av en funksjon av utvalget, $T(\mathbf{z}_D)$. Vi betegner funksjoner som er definert på utvalgsmengden som *statistikker*.⁶ Hvis funksjonen mapper til den samme mengden S som egenskapen befinner seg kan vi betegne det som en *estimator* for $\gamma(P)$. Vi håper å bruke dette til å lære om $\gamma(P)$, og vi bruker notasjonen $T := \hat{\gamma}$ for å gjøre sammenhengen eksplisitt. Merk forøvrig at det er ingenting i definisjonen av en estimator som krever at estimatoren faktisk er noe informativ om egenskapen vi er interessert i. Vi skal bruke den såkalte *utvalgsfordelingen* til estimatoren til å kvantifisere hvor informativ den er om egenskapen den estimerer.

4.2.3 Utvalgsfordelingen til estimatorer

For et gitt utvalg tar en estimator en gitt verdi. Samtidig vet vi at i andre utvalg ville vi hatt andre observasjoner og estimatoren ville tatt en annen verdi. Utvalgsfordelingen er den teoretiske fordelingen av verdier estimatoren tar i uendelig gjentatte forsøk.⁷ Dette er i praksis et tankeeksperiment siden vi kun observerer vårt ene utvalg, men denne teoretiske konstruksjonen lar oss kvantifisere hvor informativt estimatoren er om egenskapen vi er interessert i.

Hele utvalget \mathbf{z}_D er en tilfeldig variabel med fordeling $\mathcal{L}(\mathbf{z}_D) := P_D$. Det observerte utvalget utgjør bare én realisering fra denne fordelingen. Hvis de ulike observasjonene er uavhengige og fra identitisk fordeling P kan vi skrive fordelingen til utvalget på produktform⁸

$$P_D = \Pi_n \mathcal{L}(\mathbf{z}_n) = \Pi_n P \quad (4.1)$$

Estimatoren er en deterministisk transformasjon av utvalgsrommet. Siden utvalget er en tilfeldig variabel er den også en tilfeldig variabel med en fordeling $\mathcal{L}(\hat{\gamma})$. Evaluert på en mengde $A \subset S$ gir det

$$\mathcal{L}(\hat{\gamma})(A) = P_D\{\mathbf{z}_D \in Z_D : \hat{\gamma}(\mathbf{z}_D) \in A\}. \quad (4.2)$$

⁶Tror det også er krav om at de ikke må avhenge av parameter...

⁷Merk analogien til stokastiske forsøk som ble beskrevet i kapittel 1.

⁸Hadde vært enklere å gjøre dette mer eksplisitt hvis jeg brukte pdf; vurder å omskrive.

Merk at fordelingen er entydig bestemt av P_D og T , og det er i prinsippet mulig å utlede denne analytisk. Utfordringen er at P og dermed P_D er ukjente. Vi skal nå se hvordan vi kan estimere utvalgsfordelingen med utgangspunkt i det éne realiserte utvalget.

4.2.4 Estimere utvalgsfordelingen

Vi har tre ulike fremgangsmåter. Den første er å anta en parametrisk form på fordelingen P . Da kan vi få en form på P_D som gjør at vi kan utlede utvalgsfordelingen analytisk for vilkårlig utvalgsstørrelse N . Den andre måten er å estimere P_D gjennom å bruke den empiriske fordelingen P_N som estimator for P . Da kan vi være mer agnostisk om funksjonell form til fordelingen, men det krever større utvalg for at det skal gi godt estimat på utvalgsfordelingen. Den tredje måten er å bruke asymptotisk teori.

Anta parametrisk klasse

Ta bernoulli til binomial.. for gjennomsnitt.

Ta normalfordeling og knytt til t-fordeling fordi vi må estimere σ^2 .

Bootstrap

Metoden over krever ganske sterke antagelser. I praksis er modellen ofte feilspesifisert. Hvor store konsekvenser det har avhenger i hvilken grad modellen er tilstrekkelig fleksibel til å tilnærme den sanne funksjonelle formen til fordelingen av P . En alternativ fremgangsmåte er å anta at observasjonene er *iid* og bruke den empiriske fordelingen fra de N realiseringene i utvalget som estimator på P . Den empiriske utvalgsfordelingen er

$$\mathcal{L}_{\hat{P}_N}(\hat{\gamma}) := \text{fordelingen til } \hat{\gamma}(z_1, \dots, z_N) \text{ når } z_1, \dots, z_N \stackrel{iid}{\sim} \hat{P}_N \quad (4.3)$$

Denne fordelingen vil avvike fra den sanne fordelingen $\mathcal{L}_P(\hat{\gamma})$. Den er også litt vanskelig å jobbe med analytisk, blant annet fordi fordelingen er diskret. Løsningen er at vi kan sample fra fordelingen. Dette er veldig enkelt fordi fordelingen \hat{P}_N legger lik sannsynlighetstyngde på hver av realiseringene (z_1^0, \dots, z_N^0) slik at vi kan sample med tilbakelegging fra utvalget vårt. Algoritmen ser da slik ut

1. \hat{P}_N = empirisk fordeling av observert data (z_1^0, \dots, z_N^0)
2. for m in $1, \dots, M$ do:
3. trekk (x_1^b, \dots, x_N^b) fra \hat{P}_N
4. set $\hat{\gamma}_m^b = \hat{\gamma}(x_1^b, \dots, x_N^b)$
5. end for, returner utvalget $\hat{\gamma}_1^b, \dots, \hat{\gamma}_M^b$

Vi har da M realiseringer av estimatoren fra den empiriske fordelingen. Vi kan bruke standardavvik og forventningsverdi til denne empiriske fordelingen som estimat for de tilsvarende egenskapene til den sanne fordelingen til estimatoren.

Merk at denne tilnærmingen innebærer to steg av tilnærminger gjennom empirisk fordeling. Den første er at vi bruker $\mathcal{L}_{\hat{P}_N}(\hat{\gamma})$ som estimat på $\mathcal{L}_P(\hat{\gamma})$. Denne tilnærmingen er begrenset av hvor mange observasjoner vi har i utvalget. Den andre tilnærmingen er i bootstrap-samplingen av $\mathcal{L}_{\hat{P}_N}(\hat{\gamma})$. Her kan vi utgangspunktet velge antall samples m vilkårlig høyt slik at vi får vilkårlig god tilnærming.⁹

Parametrisk bootstrap

I ikke-parametrisk bootstrap resampler vi direkte fra utvalget. Dette er greit siden vi ikke trenger å gjøre noen antagelser. Alternativt kan vi ha modeller med mer struktur der vi har gjort antagelser om hvordan data i utvalget har blitt generert, for eksempel ved å spesifisere en parametrisert modell. Gitt estimatene som er beregnet fra det gitte utvalget kan vi da generere nye observasjoner fra denne modellen. Vi genererer verdier fra fordeling med estimerte parametre. Jeg tror dette blant annet kan være aktuelt i regresjonsmodeller, men litt usikker på hvordan jeg gjennomfører det i praksis og hva som er fordelene med denne fremgangsmåten.

Asymptotisk teori

Fordelingen til estimatorer avhenger av størrelsen på utvalget. Flere observasjoner gir mer presise estimat. Vi kan betrakte en følge av estimatorer, $(\hat{\theta}_N)_{N \in \mathbb{N}} = (\hat{\theta}_1, \dots, \hat{\theta}_N, \dots)$, for å se på hvordan fordelingen endrer seg når utvalget vokser. I asymptotisk teori ser vi på fordelingen i grensetilfellet der $N \rightarrow \infty$ der vi kan utlede eksakt fordeling under svakere antagelser enn det som er nødvendig for å utlede fordeling som gjelder for vilkårlig N . Denne asymptotiske fordelingen vil være en god tilnærming så lenge vi har tilstrekkelig stort utvalg...

4.3 Egenskaper til estimatorer

Nå som vi har betegnet utvalgsfordelingen og sett på hvordan vi kan estimere den fra et utvalg, kan vi gå videre til å kvantifisere hvor informativ en estimator er om en egenskap. Vi vil at tyngden i utvalgsfordelingen skal være konsentrert om egenskaper $\gamma(P)$. Dette medfører at det er lite sannsynlig at vi observerer store avvik mellom det observerte estimatet og den sanne egenskapen.

Sentrale begrep for å beskrive dette er bias og varians. Den forventningsrette estimatoren med lavest varians omtales ofte som den *effektive* estimatoren, men det er ikke

⁹I praksis kan det være ikke-trivielle kostnader i form av tid og *computational power*.

nødvendigvis så godt begrunnet at forventningsrette estimatorer skal ha så privilegert status. I praksis er vi ofte bare interessert i å minimere forventet avvik (RMSE) og at det er en trade-off mellom varians og bias. Vi kan gjøre dette litt mer formelt ved å innføre notasjon

$$\|\hat{\gamma} - \gamma\|^2 = MSE(\hat{\gamma}, \gamma) = E[(\hat{\gamma} - \gamma)^2] \quad (4.4)$$

Kan vise at dette kan dekomponeres i kvadratet bias og varians ved å skrive ut uttrykket og eliminere ledd. Men dette følger av at $\gamma \in L_2(\mathbb{I}_\Omega)$. Hvis jeg projeksjoner $\hat{\gamma}$ ned på denne mengden finner jeg $E[\hat{\gamma}]$ som er konstanten som minimerer avstand til $\hat{\gamma}$. Denne kvadrerte avstanden er variansen. Denne konstanten avviker fra parameteren som ligger i samme mengde. Denne avstanden er biasen. Disse to differansene er ortogonale fordi den ene er i $L_2(\mathbb{I}_\Omega)$ og den andre i det ortogonale komplementet til denne mengden. Det følger da fra pythagoras at den samlede kvadrerte avstanden er summen av kvadratene til katenene. TODO: Ta bias variance trade-off her slik at jeg er ferdig med det.

4.4 Estimering

Vi har innført det formelle rammeverket med estimatorer med utvalgsfordelinger som har tyngde sentrert rundt den egenskapen vi vil lære om. Hvordan skal vi kommunisere hva vi lærer fra én realisering av utvalgsfordelingen? Vi har tre generelle fremgangsmåter: punktestimat, konfidensmengder og hypotesetester.

4.4.1 Punktestimat

Punktestimat er den realiserte verdien av $\hat{\gamma}(\mathbf{z}_D)$ for utvalget. Dette utgjør gjerne vår beste gjetning på egenskapen $\gamma(P)$. Det er dessuten enkelt å representere og jobbe med siden det er element i samme mengde som egenskapen befinner seg i. Det betyr at vi kan bruke punktestimatet som en proxy for egenskapen. På en annen side vet vi at det vil være avvik mellom punktestimatet og egenskapen. Vi vet ikke hvor stort avviket er for vårt spesifikke utvalg, men vi kan beregne fordelingen til avviket. Dette kan vi bruke til å knytte en slags feilmargin til punktestimatet, for eksempel ved å rapportere standardfeil til estimatoren.¹⁰ Da kan vi fortsette å benytte punktestimatet samtidig som leseren har et inntrykk av usikkerheten knyttet til dette.

Merk forøvrig at at punktestimat ikke nødvendigvis er tall eller vektorer. Egenskapen kan være hele funksjoner, for eksempel når vi er interessert i tetthetsfunksjoner eller betinget forventningsfunksjon. Da vil den estimerte funksjonen være et punktestimat. Hvis kun er interessert i ren prediksjon vil det i noen sammenhenger være tilstrekkelig

¹⁰Som for forventningsrette estimatorer tilsvarer standardavviket til avviket.

å bruke \hat{f} som proxy for $f := E[y|\cdot]$, men i mange andre sammenhenger vil vi være interessert i å kvantifisere usikkerheten knyttet til prediksjonene.¹¹

4.4.2 Konfidensmengder

En alternativ tilnærming er å konstruere en tilfeldig mengde som med stor sannsynlighet inneholder den sanne verdien av egenskapen vi er interessert i. Konfidensmengden er en funksjon fra utvalget til mengden av delmengder av S . Vi kan spesifisere sannsynligheten for at $\gamma(P)$ er element i realiseringen av denne mengden. Jo større sannsynlighet for at den skal inneholde egenskapen, jo større (og dermed mindre informativ) må konfidensmengden være. Vi spesifiserer et nivå $\alpha \in (0, 1)$ og sier at $C(\mathbf{z}_D)$ er en $100(1 - \alpha)\%$ konfidensmengde for $\gamma(P)$ dersom

$$\mathbb{P}\{\gamma(P) \in C(\mathbf{z}_D)\} = P_D\{\mathbf{z}_D \in Z_D : \gamma(P) \in C(\mathbf{z}_D)\}. \quad (4.5)$$

Merk at dette er et argument som bruker fordelingen til $C(\mathbf{z}_D)$ og at $\gamma(P)$ er konstant. For et gitt realisert utvalg er det ikke meningsfullt å snakke om sannsynlighet for at det inneholder parameteren, i hvert fall ikke i streng tolkning av klassisk statistikk. Det vi kan si er at i mange gjentatte forsøk vil $100(1 - \alpha)\%$ av konfidensmengdene som blir konstruert på samme måte inneholde den sanne parameteren.

Litt intuisjon

Jeg tenker at det er poeng at utvalgsfordelingen til estimator gir fordeling til avviket dersom forventningsrett for parameter,

$$\hat{\theta} \sim N\left(\theta, \frac{\sigma^2}{N}\right) \quad (4.6)$$

$$\implies (\hat{\theta} - \theta) := e \quad (4.7)$$

$$\implies \hat{\theta} = \theta + e \quad (4.8)$$

der $e \sim N\left(0, \frac{\sigma^2}{N}\right)$. Dette er en latent variabel som vi ikke kan observere siden vi ikke kjenner θ . Men vi kjenner fordelingen. Den er symmetrisk og sentrert rundt 0, slik at større avvik blir gradvis mindre sannsynlig. Avvikene kan i teorien ofte være vilkårlig store, men vi kan avgrense til å kun betrakte rimelige eller plausible avvik. Dette kan vi kvantifisere ved å spesifisere et intervall av avvik som blir realisert i $100(1 - \alpha)\%$ av utvalg. Vi vil finne et tall e_{max} der $\mathbb{P}(-e_{max} < e < e_{max}) = 1 - \alpha$. Da vil konfidensmengden være

¹¹Da må vi både ta hensyn til usikkerheten til estimatoren \hat{f} samt fordeling til avvik fra betinget gjennomsnitt for hver observasjon. Dette er ikke helt trivielt å verken utlede eller representere...

bestå av det gitte realiserte punktestimatet $+-$ de største avvikene vi vil underholde,

$$C(\mathbf{z}_D) = (\hat{\theta} - e_{max}, \hat{\theta} + e_{max}) \quad (4.9)$$

Utlede konfidensmengde

Anta at $P = N(\mu, \sigma)$, at σ er kjent og at vi vil finne konfidensmengde for μ .

$$\mathcal{L}(\bar{x}_N) = N\left(\mu, \frac{\sigma^2}{N}\right) \quad (4.10)$$

$$\mathcal{L}\left(\sqrt{N}\frac{\bar{x}_N - \mu}{\sigma}\right) = N(0, 1) \quad (4.11)$$

$$(4.12)$$

dette medfører at:

$$\mathbb{P}\left\{\left|\sqrt{N}\frac{\bar{x}_N - \mu}{\sigma}\right| \leq z_{\alpha/2}\right\} \quad \text{når} \quad z_{\alpha/2} := \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \quad (4.13)$$

som kan brukes til å vise at

$$C(\mathbf{x}) = (\bar{x}_N - e, \bar{x}_N + e) \quad (4.14)$$

for $e = \frac{\sigma}{\sqrt{N}}z_{\alpha/2}$.

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \quad (4.15)$$

$$\implies P\left(z_{\alpha/2} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < z_{1-\alpha/2}\right) = 1 - \alpha \quad (4.16)$$

$$\implies P\left(z_{\alpha/2}\frac{\sigma}{\sqrt{n}} - \bar{X} < -\mu < z_{1-\alpha/2}\frac{\sigma}{\sqrt{n}} - \bar{X}\right) = 1 - \alpha \quad (4.17)$$

$$\implies P\left(\bar{X} - z_{\alpha/2}\frac{\sigma}{\sqrt{n}} < \mu < \bar{X} - z_{1-\alpha/2}\frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha \quad (4.18)$$

$$(4.19)$$

Er ikke happy med konfidensmengder altså. De er vanskelig å konseptualisere og representere dersom θ er i mer enn én dimensjon. Samtidig kan det være et problem at punktestimat ikke tar hensyn til usikkerheten. Skal nå se på den tredje fremgangsmåten som både er konseptuelt enkel og tar hensyn til usikkerhet.

4.4.3 Hypotesetester

Vi vil se om data gir tilstrekkelig bevis for å avkrefte at data generert fra modell med $\theta \in \Theta_{H_0} \subset \Theta$.

Anta at vi observerer N realiseringer fra $P_{\theta_0} \in \{P_{\theta} : \theta \in \Theta\}$. Dette gir ikke tilstrekkelig informasjon til å finne den sanne θ_0 , men vi kan jo ha lyst til å avgrense Θ der det er rimelig at den ligger. I stedet for å representere en slik konfidensmengde direkte, så kan vi spesifisere en hypotese apriori om at $\theta_0 \in \Theta_{H_0}$.¹² Dette er den såkalte nullhypotesen. Vi kan tenke at den representerer status quo eller på andre måter er en konservativ påstand. I medisinske forsøk kan de for eksempel være at en eksperimentell behandling ikke har noen effekt. Deretter skal vi se om data gir sterkt nok bevis for at vi kan forkaste hypotesen om at den sanne parameteren $\theta_0 \in \Theta_{H_0} \subset \Theta$.

En naiv fremgangsmåte vil være å bruke en punktestimator $\hat{\theta}$ og forkaste hypotesen dersom $\hat{\theta} \notin \Theta_{H_0}$. Problemet er at $\hat{\theta} \neq \theta_0$ slik at vi ikke kan trekke slutningen $\hat{\theta} \notin \Theta_{H_0} \implies \theta_0 \notin \Theta_{H_0}$. Estimatoren er en tilfeldig variabel som varierer mellom ulike utvalg, slik at den kan ta andre verdier selv om hypotesen er sann. Vi trenger derfor en buffer for at det skal gi tilstrekkelig bevis til å forkaste.

Formelt rammeverk

En *test* er en funksjon $\psi : Z_D \rightarrow \{0, 1\}$ der verdi 1 medfører at hypotesen forkastes. Verdien av testen avhenger altså av hvilke data som blir realisert. Det er litt vanskelig å konseptualisere delmengder av $\mathbb{R}^{N \times p}$, så i praksis jobber vi med en testobservator $T : Z_D \times \Theta \rightarrow \mathbb{R}$ der $T(\mathbf{z}_D, \theta_0)$ har kjent fordeling. Med andre så vet vi fordelingen gitt at nullhypotesen er riktig og kan forkaste dersom vi observerer en urimelig verdi av testobservator. Testen kan representeres som $\psi(T(\mathbf{z}_D))$...

Hvor sannsynlig det er at vi forkaster hypotese når vi ser data avhenger av parameteren i fordeling som genererer data. Dette er beskrevet med *styrkefunksjonen* $\beta_{\psi}(\theta) := \mathbb{P}_{\theta}\{\psi(\mathbf{z}_D) = 1\}$. Ideelt sett så vil vi jo at $\beta_{\psi}(\theta) = 0$ for $\theta \in \Theta_{H_0}$, slik at testen ikke forkaster dersom data blir generert av fordeling som er konsistent med nullhypotesen. Tilsvarende skulle vi kanskje ønske at $\beta_{\psi}(\theta) = 1$ for $\theta \notin \Theta_{H_0}$ slik at hypotesen alltid ble forkastet dersom hypotesen om parameteren ikke stemmer. I praksis er ikke dette mulig siden samme utvalg kan bli generert fra fordelinger med ulike parameterverdier. Det vil derfor alltid være to typer feil vi kan gjøre:

1. Type I: Forkast hypotesen selv om $\theta \in \Theta_{H_0}$.
2. Type II: Ikke forkast hypotesen selv om $\theta \notin \Theta_{H_0}$.

Vi er mest redd for å gjøre type I feil. Hypotesen representerer status quo og bevisbyrden ligger på de som vil forkaste den. Vi vil derfor utforme testen slik at det er lite sannsynlig å gjøre den feilen. Merk at det er en tradeoff her: ved å gjøre vanskeligere å forkaste vil

¹²Hvis Θ_0 er ett enkelt element (singleton) er den en enkel hypotese. Hvis det er en mengde er den en sammensatt hypotese. Det kan entent være fordi vi har en-sidet test eller fordi det er flere parametre og vi kun vil teste avgrensing på én av de.

vi øke sannsynlighet for å beholde selv om feil. Vi spesifiserer nivået α til testen er den største sannsynligheten for å gjøre type I feil, $\alpha_\psi(\theta) = \sup_{\theta \in \Theta_0} \beta_\psi(\theta)$.

Old school hypotesetestning

I praksis konstruerer vi tester med tre steg:

1. Velger nivå (α). I praksis 0.01 eller 0.05, avhenger litt av konsekvens hvis man tar feil.
2. Finner en såkalt testobservator T med kjent fordeling gitt θ .
3. Finner kritisk verdi c slik at $\sup_{\theta \in \Theta_0} \mathbb{P}_\theta\{T(\mathbf{z}_D) > c\} = \alpha$

Dette gir en test som er beskrevet med (T, c) og der $\psi(\mathbf{z}_D) = \mathbb{I}\{T(\mathbf{z}_D) > c\}$. Fremgangsmåten som ble beskrevet over er litt *old school* og mer er relevant for klassiske eksperiment i naturvitenskap enn det jeg holder på med i praksis.

Moderne hypotesetesting

I stedet for å formulere en test apriori er det mulig å ta utgangspunkt i den realiserte verdien av testobservatoren og se med hvilket nivå α det ville vært mulig å forkaste hypotesen. Formelt kan vi definere *p-verdi* til den (implisitte) testen som

$$p(\mathbf{z}_D) := \text{nivået } \alpha \text{ slik at } c(\alpha) = T(\mathbf{z}_D) \quad (4.20)$$

Altså hva nivået på testen var dersom vi lot den kritiske verdien være den realiserte verdien slik at hypotesen akkurat blir forkastet. Dette tilsvarer sannsynlighet for å se en observasjon som er minst like ekstremgitt at hypotesen er sann.

Konstruksjon av test

Anta at vi har X_1, \dots, X_N som er *iid* fra $\mathcal{L}(X) = N(\mu, \sigma)$ med kjent σ . Jeg vil teste om vi kan forkaste $H_0 : \mu < 0$. For å konstruere testen må jeg ta utgangspunkt i en størrelse som sier noe om μ gitt observerte data. Jeg bruker \bar{X}_N .

$$\psi(X_1, \dots, X_N) = \mathbb{I}\{\bar{X}_N > c\} \quad (4.21)$$

$$\beta_\psi(\mu) = \mathbb{P}_\mu(\bar{X}_N > c) \quad (4.22)$$

$$= \mathbb{P}_\mu\left(\sqrt{N}\frac{\bar{X}_N - \mu}{\sigma} > \frac{\sqrt{N}(c - \mu)}{\sigma}\right) \quad (4.23)$$

$$= 1 - \Phi\left(\frac{\sqrt{N}(c - \mu)}{\sigma}\right) \quad (4.24)$$

Ser litt indirekte at styrken øker når μ øker og reduseres når c øker. Valg av c avhenger av nivået til testen

$$\alpha_\psi = \sup_{\theta \in \Theta_0} \beta_\psi(\mu) = \beta_\psi(0) = 1 - \Phi\left(\frac{\sqrt{N}c}{\sigma}\right) \quad (4.25)$$

$$\implies c = \frac{\sigma\Phi^{-1}(1 - \alpha)}{\sqrt{N}} \quad (4.26)$$

Dette medfører at vi forkaster når

$$\bar{X}_N > \frac{\sigma\Phi^{-1}(1 - \alpha)}{\sqrt{N}} \iff \frac{\sqrt{N}\bar{X}_N}{\sigma} > \sigma\Phi^{-1}(1 - \alpha) := z_\alpha \quad (4.27)$$

I praksis bruker vi i stedet en standardnormalfordelt testobservator, men det er jo poeng at vi kan oversette til kritisk verdi av fordelingen til \bar{X} for å få det i samme måleenhet. Merk koblingen til konfidensmengder: vi forkaster dersom μ_0 ikke er i $1 - \alpha$ konfidensmengde rundt $\hat{\mu}$.

Hvorfor ikke akseptere nullhypotese?

Falsifiseringsprinsipp i vitenskap.. data kan ikke bevise at hypotese er sann (endelig antall observasjon av hvite sauer beviser ikke at alle sauer er hvite..)

Analogi til rettsal. hm.

Hypotesetest kan være feilspesifisert på måte som gir lavere enn oppgitt styrke. Får ikke forkastet selv om hypotese er feil.

Hypotesetest for å håndtere målefeil

Anta at vi er interessert i x , men observerer $\tilde{x} = x + u$ der $u \sim N(0, \sigma^2)$. Det kan for eksempel være promilletest der vi vil sjekke om noen har $x > 0.8$. Vi vil ikke straffe noen som er uskyldig, så vil at maks 5% sannsynlig at straff dersom $x \leq 0.8$. Hvor må vi da sette grense for observert \tilde{x} ? Vel, vi tar utgangspunkt i $x = 0.8$ og finner $P\left(\frac{\tilde{x}-0.8}{\sigma} < c\right) = 1 - \alpha = 0.95$ og finner $c = z_\alpha$ og kritisk målt promille som $\tilde{x}' = z_\alpha\sigma + 0.8$.

4.4.4 Modellering

Generelt er modeller forenklede representasjoner av virkeligheten som er enklere å tolke og manipulere, og dermed kan brukes til å analysere spesifikke mekanismer og svare på gitte spørsmål. Hvorvidt en modell er sann eller ikke er derfor litt *besides the point*; det avgjørende er hvorvidt det er egnet til sitt formål. Modeller har en struktur - altså, den består av størrelser og relasjoner mellom disse - som gjør at vi kan manipulere de på logiske konsekvente metoder og trekke entydige konklusjoner. Dette er mulig fordi vi

er eksplisitt om premissene. Nedsiden med dette er at premissene nødvendigvis innebærer forenklinger og ikke er oppfylt i virkeligheten. Konklusjonene vi trekker fra modell avhenger av disse premissene. De er derfor sanne for modellen, men det er ikke gitt at konklusjonen kan overføres på virkeligheten. Her kreves det dømmekraft og tester for å vurdere om antagelsene gir en rimelig tilnærming av virkeligheten.

Modellen påfører en struktur på prosessen som genererte utvalget. Dette kan for eksempel være at regresjonslinjen $E[Y|X = x] = h(x)$ er glatt (eller enda sterkere: lineær). En fordel med denne forenklete strukturen er at den resulterende modellen blir både enkel å tolke og å manipulere. Det gjør også at vi kan bruke data fra hele utvalget til å beregne verdier av $h(x)$ for ulike x . Siden vi kan bruke informasjon fra flere observasjoner blir det mindre variasjon i verdiene av $h(x)$ fra ulike utvalg enn om vi kun bruker gjennomsnitt fra lokale observasjoner. Anta for eksempel at vi vil se på sammenheng mellom gjennomsnittlig høyde og alder. For hver alder kan vi dekomponere høyde til gitt observasjon som gjennomsnitt for den alderen (signal) pluss avvik fra gjennomsnittet (støy). Hvis vi har få observasjoner for hver alder vil gjennomsnitt per alder i utvalget være sensitivt for mengden støy i det gitte utvalget vi observerer. Dersom vi bruker en parametrisk funksjonell form, for eksempel $E[høyde|alder = x] = h(x) = \beta_0 + \beta_1 x + \beta_2 x^2$ vil estimatet for hver enkelt alder bruke informasjon fra hele utvalget slik at det blir mindre sensitivt for variasjon mellom utvalg. Mer formelt er det en bias-varianse tradeoff, der vi kan forbedre de estimatene ved å påføre (potensielt feil) struktur fordi det reduserer variasjon mellom utvalg. I valg av struktur vil det være en avveining som avhenger av hvor mye data vi har i utvalg og hvor mye kunnskap vi har om fordelingen fra før.¹³

4.5 David og Mac

Jeg tenker at det er poeng at utvalgsfordelingen til estimator gir fordeling til avviket dersom forventningsrett for parameter,

$$\hat{\theta} \sim N\left(\theta, \frac{\sigma^2}{N}\right) \quad (4.28)$$

$$\implies (\hat{\theta} - \theta) := \hat{\gamma} \quad (4.29)$$

$$\implies \hat{\theta} = \theta + \hat{\gamma} \quad (4.30)$$

der $\hat{\gamma} \sim N\left(0, \frac{\sigma^2}{N}\right)$. Vi kjenner fordelingen til det tilfeldige avviket, men vi kjenner ikke fordelingen til $\hat{\theta}$ siden den avhenger av ukjent θ . Hvis vi kommer med et forslag eller hypotese $H_0 : \theta = \theta_0$ får vi hele fordelingen. Kan da forkaste nullhypotesen dersom esti-

¹³Jeg tror det er litt problemet med å forsøke å lære om struktur til fordeling fra gitt utvalg... kan finne struktur som passer til det gitte realiserde utvalget, men ikke nødvendigvis beskriver fordeling. Omtalt som p-hacking...

matet impliserer et usannsynlig stort tilfeldig avvik. I praksis jobber vi med testobservator standardiert fordeling,

$$z = \frac{\hat{\theta} - \theta_0}{se(\hat{\theta})} \sim N(0, 1) \quad (4.31)$$

merk at hvis sann θ er θ_1 er $\hat{\theta} = \theta_1 + \hat{\gamma}$. Dette impliserer at

$$z = \frac{\theta_1 + \hat{\gamma} - \theta_0}{se(\hat{\theta})} = \frac{\theta_1 - \theta_0}{se(\hat{\theta})} + \frac{\hat{\gamma}}{se(\hat{\theta})} \sim N(\lambda, 1) \quad (4.32)$$

der $\lambda := \frac{\theta_1 - \theta_0}{se(\hat{\theta})}$. Hvor mye fordelingen til testestimatoren flytter på seg avhenger både av forskjellen mellom sann parameter og nullhypotesen samt variansen til det tilfeldige avviket. Vi vil at den skal flytte seg mye slik at stor sannsynlighet for at den realiserde testobservatoren havner i forkastningsregionen når $\theta_1 \neq \theta_0$. Kan knytte dette til styrken av testen så blir funksjon av θ .

En teststatistikk er *pivotal* dersom den har samme fordeling for alle fordelingene/DGP som samsvarer med nullhypotese. Hvis det er en enkel hypotese er testobservatoren pivotal per konstruksjon, men med sammensatte hypoteser kan fordelingen avhenge av andre egenskaper ved fordelingen..

Kan utlede eksakte tester under sterke antagelser, asymptotiske tester under svakere antagelser. Teorien her blir fort veldig vanskelig... foretrekker å bruke simulering; veldig generelt og enkelt rammeverk og kan ha bedre egenskaper.

Vi bruker tester til å avgjøre om en gitt restriksjon er kompatibel med data vi observerer i et utvalg. Alternativt kan vi bruke konfidensmengder til å beskrive avgrensinger som er kompatible. For en enkelt estimator blir dette punkttestimat pluss minus kvantiler til fordelingen av avviket mellom estimator og sann parameter. Mer generelt må vi ta utgangspunkt i testobservator og vurdere om den kan forkaste nullehypoteser $\theta = \theta_0$ for ulike verider av θ . Alt den ikke kan forkaste med vårt utvalg blir da med i konfidensmengden. Hadde vært greit med fremgangsmåte for å konstruere disse. Tror bootstap er fremtiden.

Kapittel 5

Momentestimatorer

5.1 Utvalgsanalogprinsippet

Den empiriske sannsynlighetsfordelingen gir tyngde $\frac{1}{N}$ til hver av observasjonene i utvalget, $\hat{P}_N(B) \equiv \frac{1}{N} \sum \mathbb{I}\{\mathbf{z}_n \in B\}$. Det er et sentralt resultat at dersom observasjonene er *iid* vil $\hat{P}_N(B) \xrightarrow{P} P(B)$. Dette følger av store talls lov. La $h(\mathbf{z}_n) \equiv \mathbb{I}\{\mathbf{z}_n \in B\}$ og merk at forventningsverdi til en indikatorfunksjon tilsvarer sannsynligheten.

$$\frac{1}{N} \sum h(\mathbf{z}_n) \xrightarrow{P} \mathbb{E}[h(\mathbf{z})] \quad (5.1)$$

$$\implies \frac{1}{N} \sum \mathbb{I}\{\mathbf{z}_n \in B\} \xrightarrow{P} P(B) \quad (5.2)$$

For spesialtilfelle der $B = (-\infty, s]$ så følger det at $\hat{F}_N \xrightarrow{P} F$, der $\hat{F}_N(s) = \frac{1}{N} \sum \mathbb{I}\{X_n < s\}$. Dette motiverer utvalgsanalogprinsippet. I mange tilfeller kan vi finne gode estimatorer ved å evaluere γ på den empiriske sannsynlighetsfordelingen, $\hat{\gamma} = \gamma(\hat{P}_N)$.

$$\gamma(P) = \mathbb{E}_P(x) \quad (5.3)$$

$$\gamma(\hat{P}_N) = \mathbb{E}_{\hat{P}_N}(x) = \frac{1}{N} \sum x_n \equiv \bar{x} \quad (5.4)$$

For å gjøre dette litt mer operativt kan vi merke at egenskaper ofte er funksjon av kumulativ fordeling, $\gamma = \gamma(F)$. Vet ikke hvor vesentlig det poenget var. Uansett, dette er en fleksibel ikke-parametrisk fremgangsmåte der vi erstatter F med empirisk CDF \hat{F}_N og egenskapene vi måtte være interessert i. Lurer på om vi har noen generelle fremgangsmåter til å kvantifisere usikkerhet til disse estimatene uten å pålegge mer struktur..

5.1.1 Motivere OLS som utvalgsanalog

Vi kan nå bruke dette rammeverket til å studere relasjonen mellom inputvektor \mathbf{x} og avhengig variabel y . Den betingede forventningen $\mathbb{E}[y|\mathbf{x}]$ gir et godt sammendragsmål

på relasjonen, men den kan være vanskelig å estimere og å tolke. I praksis estimerer vi ofte den lineære populasjonsregresjonsfunksjonen, $\gamma(P) = \mathbf{b}^*$. Hvorvidt dette gir en god tilnærming avhenger om CEF er tilnærmet lineær. Det er likevel ganske fleksibelt siden vi kan transformere inputvektor til et *feature space*, $\Phi : \mathbb{R}^K \rightarrow \mathbb{R}^L$. Da kan vi ofte få tilnærmet lineær relasjon i forhold til den transformerte inputvektoren, $\mathbb{E}[y|\Phi(\mathbf{x})]$. Kan eventuelt også transformere y .

Dersom $\mathbb{E}(\mathbf{x}(y - \mathbf{x}'\mathbf{b}^*)) = \mathbf{0}$, observasjonene er *iid* og $\mathbb{E}(\mathbf{x}\mathbf{x}')$ er inverterbar gir utvalgsanalogsprinsippet en konsistent estimator for $\mathbf{b}^* \equiv \beta$

$$\hat{\beta} = \mathbb{E}_{\hat{P}_N}(\mathbf{x}\mathbf{x}')^{-1} \mathbb{E}_{\hat{P}_N}(\mathbf{x}y) = \left(\frac{1}{N} \sum \mathbf{x}_n \mathbf{x}_n' \right)^{-1} \frac{1}{N} \sum \mathbf{x}_n y_n \quad (5.5)$$

Utvalgsanalogsprinsippet er intuitivt. Det er naturlig å bruke den relative andelen av observasjoner som havner i en mengde som estimat på sannsynlighet for den mengden i populasjonen. Asymptotisk kan vi da observere P og lære egenskaper ved prosessen uten å måtte gjøre antagelser. Vi kan la data snakke for seg selv. Hvorfor trenger vi andre måter å utlede estimatorer? For det første har vi aldri uendelig store utvalg. Hele poenget er å generalisere fra utvalg og da må vi håndtere det faktum at $\hat{P}_N \neq P$. For det første er \hat{P}_N alltid diskret selv om fordelingen er absolutt kontinuerlig. For det andre kan vi få estimatorer med bedre egenskaper ved å påføre struktur a priori. Dette motiverer empirisk risikominimering som gir oss et rammeverk for å kombinere utvalgsanalog med struktur.

5.2 Momentestimator

Den enkleste momentestimatoren følger direkte fra utvalgsanalogsprinsippet

$$\theta = \mathbb{E}[X] = \int x f(x) dx \quad (5.6)$$

$$\hat{\theta} = \mathbb{E}_{\hat{P}_N}[X] = \sum x_n f_N(x_n) = \frac{1}{N} \sum x_n \quad (5.7)$$

Mer generelt kan anta parametrisert form $f(x; \theta)$ der $\theta = g(\mathbb{E}X)$. Fremgangsmåten er da å finne moment og løse ligningen med hensyn på parameter for å finne estimator. Eksempel: $X \sim \text{geo}(p)$, der $p = 1/\mathbb{E}X$.

$$g(p) = \mathbb{E}[X] = \int x f(x) dx \quad (5.8)$$

$$g(\hat{p}) = \mathbb{E}_{\hat{P}_N}[X] = \sum x_n f_N(x_n) = \frac{1}{N} \sum x_n \quad (5.9)$$

$$\hat{p} = \frac{N}{\sum x_n} \quad (5.10)$$

Kan utvide til å estimere flere parametre som da gir oss et ligningssystem. Jeg er litt usikker på hvilken notasjon jeg ønsker å bruke.

5.2.1 Egenskaper

Gitt regularitetsbetingelser er estimatorene konsistente og asymptotisk normale med varians som det er mulig å beregne.

$$\sqrt{N}(\hat{\theta} - \theta) \xrightarrow{d} N(0, Avar(\hat{\theta})) \quad (5.11)$$

der

$$Avar(\hat{\theta}) = g\mathbb{E}[xx']g' \quad (5.12)$$

tror det generelt har en sandwich form og at g består av derivater som kan gis enkel form dersom momentbetingelse er lineær, men ser på dette senere når jeg får notasjon på plass.

5.3 GMM

TODO: motivere GMM. Tror jeg vil motivere momentestimatorer i samme slengen. GMM er utvidelse av momentestimator til overdeterminerte ligningssystem der vi har flere instrument enn endogene variabler. Fordelen med å inkludere flere instrument er at vi får mer effektiv estimator (lavere asymptotisk varians) og kan bruke overidentifikasjon for å teste gyldighet til instrument siden vi kan observere u med $L - 1$ instrument. Det er mye notasjon, så det er viktig å være ryddig. Tenker vi har en prosess som genererer observasjoner til utvalget vårt, $\{y_n, \mathbf{x}_n, \mathbf{z}_n\} = \{\mathbf{w}_n\}$ som er ergodisk og stasjonær. Jeg har en momentbetingelse som definerer verdien på parameteren jeg vil finne. I praksis er det ofte ortogonalitetsbetingelse. Kan være greit å tenke på hvordan jeg kan relatere dette til L_2 , men annen gang. Parameter sier noe om relasjon mellom variabler i prosessen

$$g_n(\delta) = g(\mathbf{z}_n u_n(\delta)) = g(\mathbf{z}_n(y - \mathbf{x}_n' \delta)) \quad (5.13)$$

$$\mathbb{E}[g_n(\delta)] = \mathbf{0} \quad (5.14)$$

der vi kan tenke på $g_n(\delta)$ som en tilfeldig variabel. Det korresponderene empiriske momentet er

$$\mathbb{E}_{P_N}[\mathbf{g}_n(\tilde{\delta})] = \frac{1}{N} \sum \mathbf{g}_n(\tilde{\delta}) \equiv \mathbf{g}_N(\tilde{\delta}) \quad (5.15)$$

$$= \frac{1}{N} \sum \mathbf{z}_n(\mathbf{y}_n - \mathbf{x}_n' \tilde{\delta}) \equiv \mathbf{S}_{zy} - \mathbf{S}_{zx} \tilde{\delta} \quad (5.16)$$

Hvis eksakt identifisert kan jeg løse dette for $\tilde{\delta}$ som da blir min estimator $\hat{\delta}$. Hvis overidentifisert er ikke \mathbf{S}_{zx} inverterbar. Det er ikke mulig å få alle utvalgsmomentene lik 0. En naturlig løsning er å minimere det samlede avviket fra 0, altså minimere lengden av vektoren $\mathbf{g}_N(\tilde{\delta})$. Merk at $\|\mathbf{x}\|^2$ er $\mathbf{x}'\mathbf{x}$. Men vi får lavere asymptotisk varians ved å vekte momentene, slik at moment med lavere varians får høyere vekt. Litt analogt til vektet minste kvadrat. Uansett, i stedet for å minimere lengden av vektoren direkte setter vi opp en kvadratisk form

$$J(\tilde{\delta}, \hat{\mathbf{W}}) = \mathbf{g}_N(\tilde{\delta})' \hat{\mathbf{W}} \mathbf{g}_N(\tilde{\delta}) \quad (5.17)$$

der vi kan finne closed form løsning på minimeringsproblemet som gir et eksplisitt uttrykk for GMM-estimatoren.

$$\hat{\delta}_{GMM} = \arg \min_{\tilde{\delta}} J(\tilde{\delta}, \hat{\mathbf{W}}) \quad (5.18)$$

$$= \left(\mathbf{S}_{zx} \hat{\mathbf{W}} \mathbf{S}_{zx} \right)^{-1} \mathbf{S}_{zx}' \hat{\mathbf{W}} \mathbf{S}_{zy} \quad (5.19)$$

Har nå funnet estimatorene. Neste steg blir å utlede den asymptotiske fordelingen. Fremgangsmåten er å substituere inn for y og bruke dette til å få uttrykk for utvalgsfeilen $\hat{\delta} - \delta$.

$$\mathbf{S}_{zy} = \frac{1}{N} \sum \mathbf{z}_n \mathbf{y}_n = \frac{1}{N} \sum \mathbf{z}_n (\mathbf{x}_n' \delta + u_n) \quad (5.20)$$

$$= \mathbf{S}_{zx} \delta + \bar{\mathbf{g}} \quad (5.21)$$

der $\bar{\mathbf{g}} \equiv \frac{1}{N} \sum \mathbf{g}_n(\delta) = \frac{1}{N} \sum \mathbf{z}_n u_n$. Det følger da at

$$\hat{\delta} = \delta + \left(\mathbf{S}_{zx}' \hat{\mathbf{W}} \mathbf{S}_{zx} \right)^{-1} \mathbf{S}_{zx}' \hat{\mathbf{W}} \bar{\mathbf{g}} \quad (5.22)$$

slik at $\hat{\delta} - \delta = \mathbf{A}_N \bar{\mathbf{g}}$. Den er da konsistent hvis $\mathbf{A}_N \xrightarrow{p} \mathbf{A}$ og $\bar{\mathbf{g}} \xrightarrow{p} \mathbb{E}[\mathbf{g}_n(\delta)] = \mathbf{0}$. Den asymptotiske fordelingen er da

$$\sqrt{N}(\hat{\delta} - \delta) = \mathbf{A}_N \sqrt{N} \bar{\mathbf{g}} \xrightarrow{d} N(\mathbf{0}, \mathbf{A} \mathbf{S} \mathbf{A}') \quad (5.23)$$

hvis $\{\mathbf{g}_n\}$ er *mds* slik at at $\sqrt{N} \bar{\mathbf{g}} \xrightarrow{d} N(\mathbf{0}, \mathbf{S})$ og der $\mathbf{S} = \text{var}(\mathbf{g}_n) = \mathbb{E}[\mathbf{g}_n \mathbf{g}_n']$. Dette ser ganske ryddig ut, men \mathbf{A} skjuler masse dritt.

$$\mathbf{A} = (\Sigma'_{ZX} \mathbf{W} \Sigma_{ZX})^{-1} \Sigma'_{ZX} \mathbf{W} \quad (5.24)$$

Kan få ryddet opp hvis $\hat{\mathbf{W}} \xrightarrow{p} \mathbf{S}^{-1}$ som er den asymptisk effektive estimatoren. Følger da at

$$Avar(\hat{\delta}(\hat{\mathbf{S}}^{-1})) = (\boldsymbol{\Sigma}'_{ZX} \mathbf{S}^{-1} \boldsymbol{\Sigma}_{ZX})^{-1} \boldsymbol{\Sigma}'_{ZX} \mathbf{S}^{-1} \mathbf{S} \mathbf{S}^{-1} \boldsymbol{\Sigma}_{ZX} (\boldsymbol{\Sigma}'_{ZX} \mathbf{S}^{-1} \boldsymbol{\Sigma}_{ZX})^{-1} \quad (5.25)$$

$$= (\boldsymbol{\Sigma}'_{ZX} \mathbf{S}^{-1} \boldsymbol{\Sigma}_{ZX})^{-1} \quad (5.26)$$

Dette kan vi finne dersom vi har konsistent estimator for \mathbf{S}

$$\mathbf{S} = \mathbb{E}[\mathbf{g}_n \mathbf{g}_n'] = \mathbb{E}[\mathbf{z}_n u_n (\mathbf{z}_n u_n)'] \quad (5.27)$$

$$\hat{\mathbf{S}} = \frac{1}{N} \sum \hat{u}_n^2 \mathbf{z}_n \mathbf{z}_n' \quad (5.28)$$

Det er ikke helt opplagt hvorfor denne estimatoren fungerer, men tror dette er white sin hetero-robuste greie som jeg ser på senere. Problemet nå er at vi trenger $\hat{\delta}$ for å finne \hat{u}_n fordi

$$\hat{u}_n = y_n - \mathbf{x}_n' \hat{\delta} \quad (5.29)$$

men vi trenger vektematrise for å finne $\hat{\delta}$! Kan vår *catch 22* med to-steps estimator.

1. Velger en default vektematrise, som oftest $\hat{\mathbf{W}} = \mathbf{I}$ eller $\hat{\mathbf{W}} = \mathbf{S}_{zz}^{-1}$. Bruker dette til å finne $\delta(\hat{\mathbf{W}})$. Bruker dette til å finne $\hat{\mathbf{S}}$.
2. Bruker dette til å finne $\delta(\hat{\mathbf{S}}^{-1})$

Tror jeg alternativt jeg kunne brukt en algoritme som gjentar prosess til konvergens. Uansett, jeg har nå et veldig fleksibelt rammeverk som lar meg utlede asymptotisk effektive estimatorene for en stor mengde av DGPer. Hvorfor sitter ikke alle og kjører GMM? I praksis har vi ofte upresis estimering av \mathbf{S}^{-1} slik at det blir lite gevinst i forhold til 2sls. Kan påføre litt ekstra antagelser og utlede de vanlig estimatorerene som special case av asymptotisk effektive GMM og slipper da 2-steps opplegget over.

5.3.1 2SLS

Innfører antagelse om betinget homoskedastisitet

$$\mathbb{E}[u_n^2 | \mathbf{z}_n] = \sigma^2 \quad (5.30)$$

Det følger da at

$$\mathbf{S} = \mathbb{E}[\mathbb{E}(\mathbf{z}_n u_n^2 \mathbf{z}_n' | \mathbf{z}_n)] = \sigma^2 \mathbb{E}[\mathbf{z}_n \mathbf{z}_n'] \quad (5.31)$$

$$\hat{\mathbf{S}} = \hat{\sigma}^2 \frac{1}{N} \sum \mathbf{z}_n \mathbf{z}_n' = \hat{\sigma}^2 \mathbf{S}_{zz} \quad (5.32)$$

Slenger dette inn i GMM-estimatoren og får 2SLS

$$\hat{\delta}(\hat{\mathbf{S}}^{-1}) = (\mathbf{S}_{zx}(\hat{\sigma}^2 \mathbf{S}_{zz})^{-1} \mathbf{S}_{zx})^{-1} \mathbf{S}_{zx}' (\hat{\sigma}^2 \mathbf{S}_{zz})^{-1} \mathbf{S}_{zy} \quad (5.33)$$

$$= (\mathbf{S}_{zx}(\mathbf{S}_{zz}^{-1} \mathbf{S}_{zx})^{-1} \mathbf{S}_{zx}' \mathbf{S}_{zz}^{-1} \mathbf{S}_{zy} \quad (5.34)$$

$$= \hat{\delta}(\hat{\mathbf{S}}_{zz}^{-1}) = \hat{\delta}_{2SLS} \quad (5.35)$$

Det at estimatoren er konsistent og asymptotisk normalfordelt følger av at det er special case av GMM. Kan finne asymptotisk varians med homoskedastisitet

$$Avar(\hat{\delta}_{2SLS}) = (\mathbf{\Sigma}'_{zx}(\sigma^2 \mathbf{S}_{zz})^{-1} \mathbf{\Sigma}_{zx})^{-1} \quad (5.36)$$

$$\widehat{Avar}(\hat{\delta}_{2SLS}) = \hat{S}^2 (\mathbf{\Sigma}'_{zx} \mathbf{S}_{zz}^{-1} \mathbf{S}_{zx})^{-1} \quad (5.37)$$

Kan jo også utvide til robuste standardfeil her. Kan bruke 2SLS (og OLS) selv om antagelse om homoskedastisitet ikke er oppfylt, men da må vi leve med at estimatorene ikke er asymptotisk effektive. For spesialtilfelle med eksakt identifisering er det tilstrekkelig å bruke at $\mathbf{\Sigma}_{zx}$ er inverterbar til å utlede IV og OLS fra GMM.

Kapittel 6

Maximum likelihood

Vi ønsker å lære om en fordeling P med utgangspunkt i realiserte verdier fra fordelingen. Dette er det generelle utgangspunktet i statistisk inferens. I likelihood tilnærmingen gjør vi sterke antagelser ved å anta at fordelingen tilhører en parametrisk klasse $P := P_{\theta_0} \in \mathcal{P} = \{f(\cdot; \theta) : \theta \in \Theta \subset \mathbb{R}^k\}$ der $f(\cdot; \cdot)$ har kjent form.¹ Med andre ord antar vi at fordelingen er kjent opp til en ukjent parameter og at denne parameteren fullt ut beskriver hele fordelingen som genererer data vi observerer. Funksjonen $f(\cdot; \theta)$ evaluert i en gitt verdi $\theta = \theta'$ er en sannsynlighetstetthetsfunksjon.² Hvis vi derimot holder den realiserte verdien konstant og betrakter det som en funksjon av θ kan vi betegne det som likelihoodfunksjonen, $f(\cdot; z) := L(\cdot; z) := L(\cdot)$. Vi kan betrakte z enten som en gitt realisert verdi eller som en tilfeldig variabel $z := z(\omega)$ som tar en konstant verdi når tilstanden ω blir avslørt. I den første tolkningen har likelihoodfunksjonen en ganske konkret tolkning som relativ sannsynlighet for at en fordeling med de ulike parameterverdiene kan ha generert den observerte verdien z .³ Hvis vi derimot betrakter situasjonen før verdi er realisert er likelihoodfunksjonen utfallet av $g : \Omega \rightarrow S = \{L(\cdot; z) : z^{-1}(\omega) \in \Omega\}$. Utfallet er en funksjon! Men hvis vi bare evaluerer dette i gitte verdier av θ blir det en tilfeldig variabel med fordeling som vi kan beskrive og analysere på vanlig måte.

Denne tilnærmingen krever sterke antagelser, men har til gjengjeld en veldig rik teori. Estimatoren vil også kunne ha gode egenskaper selv om modellen er feilspesifisert.

¹Denne mengden burde kanskje bestått av fordelinger, ikke tetthetsfunksjoner. Men så lenge modell er identifiserbar eksisterer det injektive funksjoner $\theta \mapsto P_\theta$ og $P_\theta \mapsto f_\theta$ slik at de ulike representasjonene av fordelingen er kjent når θ er kjent.

²Eller en pmf. Distinksjonen er ikke viktig og hvis jeg kunne litt measure theory tror jeg fremstillingen kunne blitt gjort mer ryddig.

³Vet ikke om tolkningen er helt presis. Merk at det ikke er en sannsynlighetsfordeling siden det ikke oppfyller aksiom. Siden den kun betegner relativ sannsynlighet er den bare unik opp til en multiplikativ konstant siden slike skaleringer inneholder samme informasjon. Det er praktisk siden det medfører at likelihoodfunksjonen er invariant til valg av måleenhet på observasjonen.

6.1 Begreper

Likelihoodfunksjonen for en gitt observasjon har en konkret tolkning. Formen på funksjonen beskriver i hvilken grad ulike parameterverdier korresponderer med den realiserte verdien vi observerer. Det er uansett begrenset hvor vi kan lære fra én enkelt realisering. Heldigvis er det enkelt å kombinere informasjon fra ulike realiseringer fra samme fordeling så lenge disse er uavhengige. Anta at vi observerer (z_1, \dots, z_n) der $\mathcal{L}(z_n) = P_{\theta_0}$ for $n = 1, \dots, N$. Vi kan da kalle likelihoodfunksjonen for hver av observasjonene, $L_n(\cdot; z_n)$ for *likelihood contribution* til observasjon n . Vi kan kombinere informasjonen ved å betrakte hele utvalget som én realisering fra simultanfordelingen $\mathcal{L}(z_1, \dots, z_n) = \pi_n \mathcal{L}(z_n)$ slik at $L(\cdot; z_1, \dots, z_n) = \pi_n L_n(\cdot; z_n)$. Den samlede likelihoodfunksjonen er da

$$L : \Theta \rightarrow [0, \infty) \quad (6.1)$$

$$: \theta \mapsto \Pi f(z_n; \theta) = \Pi f(z_n; \theta) \quad (6.2)$$

Likelihoodfunksjonen er riktignok ikke unik; alle skaleringer med positiv konstant inneholder akkurat like mye informasjon om θ siden vi kun kan vurdere relativ sannsynlighet for ulike parameterverdier gitt observert utvalg. Dette har litt sammenheng med at vi ønsker at likelihood skal være invariant for én-til-én transformasjoner av data, for eksempel valg av måleenhet. La $y = g(x)$ og $x = g^{-1}(y) := x(y)$. Da er

$$F_Y(y) = P(Y < y) = P(g(X) < y) = P(X < x(y)) = F_x(x(y)) \quad (6.3)$$

$$f_Y(y) = \frac{\partial}{\partial y} F_x(x(y)) = f_x(x(y)) \left| \frac{dx}{dy} \right| \quad (6.4)$$

Det følger derfor at $L(\theta|x) = f_X(\theta; x)$ og $L(\theta|y) = f_Y(\theta; y) = L(\theta|x) \left| \frac{dx}{dy} \right|$. Disse er ulike med positiv skalar, men relativ likelihood evaluert i to ulike verdier θ_1 og θ_2 vil være den samme for begge funksjonene og er dermed begge like gode likelihoodfunksjon for $\mathcal{L}(X)$.

I praksis er det enklere å jobbe med logaritmen av likelihood når vi kombinerer informasjon fra ulike kilder siden

$$\log L(\cdot; z_1, \dots, z_n) = \log[\pi_n L_n(\cdot; z_n)] \quad (6.5)$$

$$= \Sigma_n \log[L_n(\cdot; z_n)] \quad (6.6)$$

$$= \Sigma_n \log L_n(\cdot; z_n) \quad (6.7)$$

$$(6.8)$$

Logaritmen er en positive monoton transformasjon som bevarer $\arg \max$, gjør det enklere å kombinere informasjon fra avhengige observasjoner og optimere numerisk. Funksjonsverdiene har ikke like umiddelbar tolkning som i likelihoodfunksjonen, men vi skal se at de sentrale teoretiske størrelsene er knyttet til denne såkalte *loglikelihood-funksjonen*. Merk

også at verdien for hver θ er summen av N uavhengige realiseringer fra P_{θ_0} slik at hvis vi skalerer det med $1/N$ så vil det konvergere mot⁴

$$\mathbb{E}_{\theta_0}[\log L(\theta, z)] := \int_Z \log L(\theta, z) f(z; \theta_0) dz \quad (6.9)$$

6.1.1 Score

Helningen⁵ til loglikelihood-funksjonen betegnes som dens *score*,

$$S_n(\theta) = \frac{\partial}{\partial \theta} \log L_n(\theta) \quad (6.10)$$

$$\implies S(\theta) = \frac{\partial}{\partial \theta} \log L(\theta) = \Sigma_n S_n(\theta) \quad (6.11)$$

Hvis vi betrakter denne størrelsen som en funksjon av θ for gitt realisert verdi av z betegnes det som *score-funksjonen* og hvis vi derimot betrakter hvordan verdien for en gitt θ avhenger av tilfeldig z betegnes det som *score-statistic*.

Det gir mening at P_{θ_0} asymptotisk generer et utvalg som korresponderer med θ_0 i betydningen at av alle kandidat-verdier av θ så er det mest sannsynlig at det ble generert fra fordeling med θ_0 . Det medfører at forventningsverdien til log-likelihoodfunksjonen er størst når den er evaluert i θ_0 og tilsvarende at forventningsverdi til score-statistikken i θ_0 er 0,

$$\mathbb{E}_{\theta_0} S_n(\theta_0) := \int S_n(\theta_0) f_{\theta_0}(x) dz \quad (6.12)$$

$$= \int \frac{\partial}{\partial \theta} \log L_n(\theta_0) f_{\theta_0}(x) dz \quad (6.13)$$

$$= \int \frac{\frac{\partial}{\partial \theta} L_n(\theta_0)}{L_n(\theta_0)} f_{\theta_0}(x) dz \quad (6.14)$$

$$= \int \frac{\partial}{\partial \theta} L_n(\theta_0) dz \quad (6.15)$$

$$= \frac{\partial}{\partial \theta} \int L_n(\theta_0) dz = 0 \quad (6.16)$$

der vi har brukt at $L(\theta_0) := L(\theta_0; z) := f_{\theta_0}(z) := f(z; \theta_0)$. Utvalgsanalogen til dette er å velge

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \mathbb{E}_{\hat{P}_N} [S_n(\theta)] \quad (6.17)$$

$$= \arg \max_{\theta \in \Theta} \frac{1}{N} \sum S_n(\theta, z_n) = 0 \quad (6.18)$$

som vi i enkle eksempler kan finne en *closed form* løsning på slik at vi får et eksplisitt ut-

⁴Det er direkte resultat av store talls lov: $\mathbb{E}_{\hat{P}_N}[g(z)] \xrightarrow{P} \mathbb{E}[g(z)]$.

⁵Skal generalisere til parametervektor med flere dimensjoner senere. Da vil det være gradienten.

trykk $\hat{\theta}_{MLE} = g(z_1, \dots, z_N)$ og som ofte vil tilsvare momentestimatoren.⁶ Vi vil maksimere forventen log-likelihood, men vi observerer det ikke så vi maksimerer i stedet gjennomsnitt loglikelihood fra observasjonen generert av sann fordeling og lener oss på at størrelsene konvergerer asymptotisk.

6.1.2 Informasjon

Hvor mye lærer vi om θ_0 fra å observere én realisering fra P_{θ_0} ? Som nevnt lener vi oss på at $\mathbb{E}_{\theta_0} S_n(\theta_0) = 0$ og at gjennomsnittet i mitt utvalg bestående av N iid observasjoner konvergerer mot denne sentraltendensen. Men jeg har et begrenset antall observasjoner og vil derfor være interessert i mål på spredningen til den tilfeldige variabelen $S_n(\theta_0)$. Dette angir den teoretiske *fisher-informasjonen* til den ukjente fordelingen P_{θ_0} ,

$$I(\theta) := \mathbb{V}_{\theta_0}[S(\theta_0)] \quad (6.19)$$

$$= \mathbb{E}_{\theta_0}[S(\theta_0)^2] \quad (6.20)$$

$$:= \int_Z S(\theta_0)^2 f_{\theta_0}(z) dz \quad (6.21)$$

$$:= \int_Z S(\theta_0, z)^2 f(z; \theta_0) dz \quad (6.22)$$

$$(6.23)$$

Denne spredningen avhenger av hvor spiss toppen til $\mathbb{E}_{\theta_0}[\log L(\theta)]$ er i $\theta = \theta_0$. For at det skal være mye informasjon om θ i hver observasjon av z vil vi at den skal ha en spiss topp i θ_0 . Det viser seg at vi kan bruke denne intuisjonen til å finne en alternativ utledning av fisher-informasjonen. Generelt vil hesse-matrisen beskrive krumming til score-funksjon,

$$H(\theta) = \frac{\partial}{\partial \theta} S(\theta) = \frac{\partial^2}{\partial \theta \partial \theta'} \log L(\theta). \quad (6.24)$$

Den teoretiske fisher-info tilsvarende den negative verdien av den forventede hessematrisen evaluert i θ_0 ,

$$I(\theta_0) = -\mathbb{E}_{\theta_0} \left[\frac{\partial}{\partial \theta} S(\theta_0) \right] \quad (6.25)$$

$$= -\mathbb{E}_{\theta_0} [H(\theta_0)] \quad (6.26)$$

$$= - \int_Z H(\theta_0, z) f(z; \theta_0) dz \quad (6.27)$$

Mer informasjon gir bedre presisjon av $\hat{\theta}_{MLE}$. Noe kontra-intuitivt er det derfor bedre med høyere varians til $S(\theta_0)$. Et veldig sentralt resultat, som jeg forhåpentligvis kommer

⁶Kanskje finne noe måte å koble de sammen.

tilbake til senere, er at for alle MLE-estimatorer vil

$$\hat{\theta}_{MLE} - \theta_0 \xrightarrow{d} N(0, I(\theta)^{-1}) \quad (6.28)$$

Så langt har jeg betraktet fisher-informasjon som en ren teoretisk størrelse ved P_{θ_0} , men informasjon i utvalget avhenger i tillegg av antall observerte verdier fra fordelingen. For et utvalg tar jeg utgangspunkt i den samlede log-likelihoodfunksjonen som er sum av log-likelihood-contribution fra $n = 1, \dots, N$. Det medfører at utvalgsstørrelse blir bakt inn i fisher-informasjonen slik at det ikke er en eksplisitt N med i uttrykket.

6.1.3 Alternativ utledning

Log-likelihood contribution for gitt parameterverdi, $\log L_n(\theta; z_n)$, er en tilfeldig variabel med forventningsverdi $\mathbb{E}_{\theta_0}[\log L_n(\theta)] := \int_Z \log[f(z; \theta)] f(z; \theta_0) dz$ som er en skalar. For å understreke at $\mathbb{E}_{\theta_0}[\log L_n(\cdot)]$ er en helt vanlig funksjon som mapper $\mathbb{R}^d \rightarrow \mathbb{R}$ vil jeg betegne den som g .⁷ Denne funksjonen g kan utledes fra fordelingen P_{θ_0} og vi kan finne egenskaper ved denne helt streite funksjonen som dermed er egenskaper ved den sanne fordelingen.

- Score: $\frac{\partial}{\partial \theta} g(\theta)|_{\theta_0}$
- Fisher-informasjon: $\frac{\partial^2}{\partial \theta \partial \theta} g(\theta)|_{\theta_0}$

Vi kan ikke observere P_{θ_0} og kjenner dermed ikke g . Men vi kan bruke utvalgsanaloget til å tilnærme oss funksjonen siden $\mathcal{L}(z_n) = P_{\theta_0}$ medfører at

$$\mathbb{E}_{\hat{P}_N}[\log L_n(\theta)] := \frac{1}{N} \sum_n \log L_n(\theta; z_n) \quad (6.29)$$

$$\xrightarrow{p} \mathbb{E}_{\theta_0}[\log L_n(\theta)] = g(\theta) \quad (6.30)$$

Vi kan konsistent estimere $g(\cdot)$ med gjennomsnitt i utvalget og bruke det til å beregne egenskaper til funksjonen. Den teoretiske fisher-informasjonen gir oss for eksempel et mål på hvor mye vi lærer om θ_0 fra én observasjon. For å finne den estimerte fisher-informasjonen i utvalget trenger vi bare å skalere gjennomsnittet med N observasjoner. Dette gir oss tilbake hessematrisen fra $\log L(\cdot)$ i utvalget.

⁷Den mapper fra parametermengden som generelt er delmengde av \mathbb{R}^d .

6.2 Eksempler

6.2.1 Bernoulli

Likelihoodfunksjonen kan skrives kompakt på én linje,

$$L_n(\rho) = \rho_n^x (1 - \rho)^{1-x_n} \quad (6.31)$$

$$\log L_n(\rho) = x_n \log(\rho) + (1 - x_n) \log(1 - \rho) \quad (6.32)$$

Vi deriverer med hensyn på parameter og finner score contribution,

$$\frac{\partial}{\partial \rho} \log L_n(\rho) := S_n(\rho) = \frac{x_n}{\rho} - \frac{1 - x_n}{1 - \rho} \quad (6.33)$$

$$= \frac{x_n - \rho}{\rho(1 - \rho)} \quad (6.34)$$

Vi kan anta at $\mathcal{L}(x_n) = \text{bernoulli}(\rho_0)$. Forventningsverdi til score contribution er da

$$\mathbb{E}_{\rho_0} [S_n(\rho)] = \mathbb{E}_{\rho_0} \left[\frac{x_n - \rho}{\rho(1 - \rho)} \right] \quad (6.35)$$

$$= \frac{\rho_0 - \rho}{\rho(1 - \rho)} \quad (6.36)$$

som medfører at $\mathbb{E}_{\rho_0} [S_n(\rho) | \rho_0] = 0$. Vi kan også finne variansen til score contribution evaluert i sann parameter

$$\mathbb{V}_{\rho_0} [S_n(\rho_0)] = \mathbb{E}_{\rho_0} [S_n(\rho_0)^2] \quad (6.37)$$

$$= \frac{1}{(\rho_0(1 - \rho_0))^2} \mathbb{E}_{\rho_0} [(x_n - \rho_0)^2] \quad (6.38)$$

$$= \frac{1}{(\rho_0(1 - \rho_0))^2} \mathbb{V}_{\rho_0} [x_n] \quad (6.39)$$

$$= \frac{1}{\rho_0(1 - \rho_0)} \quad (6.40)$$

Kan også vises at dette tilsvarer den negative forventningsverdien til hessen evaluert i sann parameter, men hessen er ganske stygg siden vi må bruke kvotientregel. I stedet utvider jeg til å se på utvalg som består av summen av N contributions.

$$\log L(\rho) = \sum_{n=1}^N \log L_n(\rho) = \sum_{n=1}^N [x_n \log(\rho) + (1 - x_n) \log(1 - \rho)] \quad (6.41)$$

$$S(\rho) = \sum_{n=1}^N S_n(\rho) = \frac{\sum_{n=1}^N [x_n - \rho]}{\rho(1 - \rho)} \quad (6.42)$$

Dette medfører at $\mathbb{E}[\log L_n(\rho)] = \mathbb{E}[\frac{1}{N} \log L(\rho)]$. Forventningsverdi til størrelsene er den samme, men med flere observasjoner så gir evaluering av forventningsverdi med hensyn på empirisk fordeling en bedre tilnærming. *Dette er reflektert i høyere fisher-informasjon.* Vi finner $\hat{\rho}_{MLE}$ ved å løse

$$\mathbb{E}_{\hat{P}_N} [S(\hat{\rho})] = 0 \quad (6.43)$$

$$\implies \hat{\rho} = \bar{x}_N \quad (6.44)$$

For å finne variansen til denne punktestimatoren må vi beregne fisher-informasjonen i utvalget. Med uavhengige variabler kan vi enkelt summere variansene slik at

$$\mathbb{V}_{\rho_0} [S(\rho_0)] = N \cdot \mathbb{V}_{\rho_0} [S_n(\rho_0)] \quad (6.45)$$

$$= \frac{N}{\rho_0(1 - \rho_0)} := I(\rho_0) \quad (6.46)$$

Kunne kanskje evaluert variansen med hensyn på empirisk fordeling. Tror vi får samme resultat av å plugge inn punktestimator. Ble litt usikker. Uansett har vi nå at

$$(\hat{\rho} - \rho_0) \xrightarrow{d} N(0, I(\rho_0)^{-1}) = N\left(0, \frac{\rho_0(1 - \rho_0)}{N}\right) \quad (6.47)$$

og vår estimasjon av denne fordelingen fra vårt éne realiserte utvalg er

$$N\left(0, \frac{\hat{\rho}(1 - \hat{\rho})}{N}\right) \quad (6.48)$$

Dette er vårt beste forsøk på å tilnærme oss den ukjente asymptotiske fordelingen som igjen uansett bare vil være en tilnærming for vårt endelige utvalg. Med endelig antall observasjoner kan $\hat{\rho}$ bare ta et endelig antall verdier, så den eksakte fordelingen kan ikke være kontinuerlig. Vi kan vise at fordelingen til $N \times \hat{\rho}$ er $\text{binom}(N, \rho)$ og brukt dette resultatet i stedet, men er jo kjekt at MLE gir et generelt rammeverk til å finne asymptotisk fordelingen til stor klasse av estimatorer med gode asymptotiske egenskaper!

6.2.2 Normalfordeling med kjent varians

$$f(x; \mu) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\} \quad (6.49)$$

$$\log L_n(\mu) = -\frac{(x - \mu)^2}{2\sigma^2} + \dots \quad (6.50)$$

$$S_n(\mu) = \frac{\partial}{\partial \mu} \log L_n(\mu) = \frac{x - \mu}{\sigma^2} \quad (6.51)$$

$$\mathbb{E}[S_n(\mu)] = 0 \implies \mu = \mathbb{E}[x] \quad (6.52)$$

$$I(\mu) := -\mathbb{E}\left[\frac{\partial}{\partial \mu} S_n(\mu)\right] = \mathbb{E}[\sigma^{-2}] \quad (6.53)$$

$$Avar(\hat{\mu}) = I(\mu)^{-1} = \sigma^2 \quad (6.54)$$

$$\implies \sqrt{N}(\hat{\mu} - \mu) \xrightarrow{d} N(0, \sigma^2) \quad (6.55)$$

6.2.3 Uniform

Davidson og MacKinnon gjør et poeng av at det er to typer ML estimatorer,

- Type 1: $\hat{\theta} = \arg \max_{\theta \in \Theta \subset \mathbb{R}^K} \log L(\mathbf{z}_d, \theta)$
- Type 2: Implisitt definert av $S(\mathbf{z}, \hat{\theta}) = 0$, der $S(\mathbf{z}, \hat{\theta})$ er score-vektor med typisk element $S_k = \frac{\partial \log L(S(\mathbf{z}, \hat{\theta}), \theta)}{\partial \theta_k}$.

De fleste ML estimatorer er begge typer, men det er ikke alltid vi finner optimum gjennom første ordens betingelse fordi likelihoodfunksjonen ikke er differensierbar. Eksempel på det er uniformfordeling...

6.2.4 Andre hendelser

Vi observerer utfall fra P_θ og på bakgrunn av dette vil vi kvantifisere relativ sannsynlighet for ulike verdier av den ukjente θ . Likelihoodfunksjonen fikser dette og kan håndtere ulike typer utfall. La oss ta normalfordeling med kjent $\sigma = 1$ som eksempel

1. Observere utfall direkte, f. eks. $X = 1.3$, $L(\theta) = \phi(1.3 - \theta)$
2. Observere at utfall er i et intervall, f.eks. $X \in (1, 3)$, $L(\theta) = P_\theta(X \in (1, 3)) = \Phi(3 - \theta) - \Phi(1 - \theta)$
3. Observere en funksjon av utfall, f.eks. $Y = g(X) = \max(X_1, \dots, X_N) := X_N = 4$, $L(\theta) = P(g(X) = 4) = N\Phi(4 - \theta)^{N-1}\phi(4 - \theta)$

der siste følger av at $P(X_N < s) = P(X_1 < s, \dots, X_N < s) = \Phi(s - \theta)^N$ og $L(\theta) = f_\theta(s) = \frac{\partial}{\partial s} F(s)$. Ser generelt at jeg vil evaluere tettheter og eventuelt integral over tettheter dersom jeg ikke har eksakt verdi. Kan også enkelt kombinere informasjon fra ulike kilder så lenge

observasjonen er uavhengige. Med log-likelihood er det bare å summere opp funksjonene. Det kan enten være enkeltobservasjoner eller fra ulike utvalg. Trenger ikke justere for antall observasjoner som ingikk for å konstruere funksjonen, siden all informasjon er oppsumert i selve funksjonen.

6.3 Oppsummere informasjon fra likelihoodfunksjonen

For et gitt utvalg vil likelihoodfunksjonen angi relativ sannsynlighet for ulike parameterverdier. Dette gir både informasjon om hvilke verdier som er mest sannsynlig og samt hvor sikre vi er på at den sanne, ukjente parameteren er i ulike intervall. Det er en utfordring at det kan være vanskelig å kommunisere denne informasjonen, spesielt hvis parameteren er en vektor slik at likelihood blir funksjon av flere variabler. Det kan dessuten være litt vanskelig å jobbe med funksjoner. Vi vil derfor ønske å finne alternative måter å oppsummere informasjon i likelihoodkurven. Ser generelt at det er både enklere å jobbe med log-likelihood numerisk og at analytiske resultat bruker denne formen.

Vi vet for det første at maksimumsverdien gir det best punkttestimatet. Videre vil spissheten til funksjonen rundt $\hat{\theta} = \arg \max L(\theta)$ gi et mål på hvor sikre vi er på at den gitte $p_{\hat{\theta}}$ har generert utvalget. Hvis det er flatt på toppen er det mange ulike kandidater som er omtrent like sannsynlige, eg. kan korrespondere med observasjonene vi har observert, slik at det er mye usikkerhet knyttet til $\hat{\theta}$.

6.3.1 Kvadratisk tilnærming

Vi kan bruke størrelsene over til å beregne en andre ordens taylor ekspansjon av $\log L(\theta)$ i $\theta = \hat{\theta}$.

$$\log L(\theta) \approx \log L(\hat{\theta}) + S(\hat{\theta})(\theta - \hat{\theta}) - \frac{I(\hat{\theta})^{-1}}{2}(\theta - \hat{\theta})^2 \quad (6.56)$$

$$= \log L(\hat{\theta}) - \frac{I(\hat{\theta})^{-1}}{2}(\theta - \hat{\theta})^2 \quad (6.57)$$

Hele formen på funksjonen er da beskrevet av punktet $(\hat{\theta}, \log L(\hat{\theta}))$ samt den estimerte fisher-informasjonen. Det er mulig å vise at dette gir en eksakt beskrivelse av loglikelihoodfunksjonen til forventningsverdien av normalfordeling. For andre likelihoodfunksjoner vil det være en god tilnærming dersom de er såkalt *regulære*. Det kan vises at de fleste loglikelihoodfunksjoner konvergerer mot denne kvadratiske formen når antall observasjoner øker og dette har litt sammenheng med CLT. Jeg skal nå utvikle teoretiske resultat som har utgangspunkt i denne forenklete representasjonen. Disse resultatene vil holde eksakt for gjennomsnitt av normalfordeling og være en asymptotisk tilnærming for andre fordelinger. Først skal jeg bare utlede to enkle sammenhenger til. Det første er en forenklet

representasjon av normalisert likelihood.

$$\log \left(\frac{L(\theta)}{L(\hat{\theta})} \right) = \log L(\theta) - \log L(\hat{\theta}) \quad (6.58)$$

$$\approx \frac{I(\hat{\theta})^{-1}}{2} (\theta - \hat{\theta})^2 \quad (6.59)$$

Den andre sammenhengen er bare første ordens taylor ekspansjon av score funksjonen som kan brukes til å visualisere om den kvadratiske tilnærmingen er god ved å se om sammenhengen under faktisk er lineær.

$$S(\theta) \approx S(\hat{\theta}) - I(\hat{\theta})^{-1}(\theta - \hat{\theta}) \quad (6.60)$$

$$= -I(\hat{\theta})^{-1}(\theta - \hat{\theta}) \quad (6.61)$$

6.3.2 Konfidensintervall

Jeg kan ha lyst til å konstruere et intervall $\Theta_c \subset \Theta$ der det virker rimelig at $P \in \{P_\theta : \theta \in \Theta_c\}$ kan ha generert det utvalget jeg observerer. Et greit utgangspunkt kan være å betrakte mengden

$$\Theta_c = \left\{ \theta : \frac{L(\theta)}{L(\hat{\theta})} > c \right\} \quad (6.62)$$

Spørsmålet nå er hvordan vi skal velge c . Vil forsøke å knytte det til noe som i prinsippet er en observerbar sannsynlighet. Vil manipulere uttrykket slik at jeg får en størrelse med kjent fordeling. Begynner med å observere at

$$\log \left(\frac{L(\theta)}{L(\hat{\theta})} \right) = \frac{\sigma^2}{2N} (\bar{x} - \theta)^2 \quad (6.63)$$

for gjennomsnitt av normalfordeling og at dette kan være god tilnærming for andre regulære likelihoods. Jeg vet at

$$\bar{x} \sim N \left(\mu, \frac{\sigma^2}{N} \right) \implies 2 \cdot \frac{L(\hat{\theta})}{L(\theta)} \sim \chi^2(1) \quad (6.64)$$

Manipulerer begge sider av ulikheten over og får

$$P \left(2 \cdot \frac{L(\hat{\theta})}{L(\theta)} < -2 \log(c) = \chi^2_{1-\alpha} \right) = 1 - \alpha \quad (6.65)$$

6.4 Likelihood i flere dimensjoner

Har nå sett at jeg kan anta at observasjonene i utvalget mitt er *iid* realiseringer fra $P_\theta \in \mathcal{P}_\theta$ og bruke parameterverdien som gjør det mest sannsynlig å observere verdiene i utvalget mitt som estimat. Dette kan vi generalisere til observasjoner $\mathbf{z} \in \mathbb{R}^d$ der P_θ nå blir en simultanfordeling. Det er en utfordring at simultanfordelinger er veldig komplekse objekter. Da skal angis sannsynlighet for utfall i ganske vilkårlige delmengder av \mathbb{R}^d . Det er både vanskelig å estimere og beskrive. I praksis vil vi ofte heller si noe om betinget sannsynlighet.

6.4.1 Betinget likelihood

I praksis vil vi ofte dekomponere $\mathbf{z} = (\mathbf{x}, y)$ og se på hvordan \mathbf{x} påvirker fordeling av y . Da får vi bruk for at

$$f(\mathbf{x}, y) = f(y|\mathbf{x})f(\mathbf{x}) \quad (6.66)$$

der vi er interessert i $f(y|\mathbf{x})$. Vi kan parametrisere tetthetene over slik at

$$f(\mathbf{x}, y; \theta, \gamma) = f(y|\mathbf{x}; \theta)f(\mathbf{x}; \gamma) \quad (6.67)$$

Hvis vi tar log-likelihood får vi

$$\log L(\theta, \gamma) = \log(f(y|\mathbf{x}; \theta)) + \log(f(\mathbf{x}; \gamma)) \quad (6.68)$$

Hvis vi bare er interessert i θ så er andre leddet en uvesentlig konstant. Får samme estimat ved å kun betrakte første del som om vi betraktet likelihood til hele simultanfordelingen.⁸

6.4.2 Generell fremgangsmåte til å finne likelihood til betinget fordeling

Vi har en regresjonsmodell

$$y = \mathbf{x}'\beta_0 + u, \quad u|\mathbf{x} \sim N(0, \sigma_0^2) \quad (6.69)$$

⁸Gitt at det ikke er funksjonell relasjon mellom parameterene (θ, γ) ... i praksis vil vi neppe ønske å modellere dette.

Vi begynner med å finne betinget kumulativ sannsynlighet

$$P(Y \leq y|\mathbf{x}) = F(\mathbf{x}'\beta_0 + u < y|\mathbf{x}) \quad (6.70)$$

$$= F(u < y - \mathbf{x}'\beta_0|\mathbf{x}) \quad (6.71)$$

$$= F\left(\frac{u}{\sigma_0} < \frac{y - \mathbf{x}'\beta_0}{\sigma_0}|\mathbf{x}\right) \quad (6.72)$$

$$= \Phi\left(\frac{y - \mathbf{x}'\beta_0}{\sigma_0}\right) \quad (6.73)$$

der $F(c|\mathbf{x}) := \int_{-\infty}^c yf(y|\mathbf{x})dy$. Merk at selv om $\mathbf{x}'\beta$ er tilfeldig så kan vi behandle det som en konstant når vi betinger av \mathbf{x} . Vi kan deretter enkelt finne tetthet ved å derivere

$$f(y|x) = \frac{\partial}{\partial y} \Phi\left(\frac{y - \mathbf{x}'\beta_0}{\sigma_0}\right) \quad (6.74)$$

$$= \frac{1}{\sigma_0} \phi\left(\frac{y - \mathbf{x}'\beta_0}{\sigma_0}\right) \quad (6.75)$$

Denne fremgangsmåten bruker eksplisitt antagelse om betinget fordeling til feilledd i stedet for å modellere betinget fordeling til y direkte.. Tror det er litt ulike måter man kan gjøre dette på.

6.4.3 Betinget normal

Et konkret eksempel er $y|\mathbf{x} \sim N(\mu_x, \sigma_x^2) = N(g(\mathbf{x}), h(\mathbf{x})^2)$ og spesifisere hvordan parametrene avhenger av \mathbf{x} . Vanlig valg er $g(\mathbf{x}) = \mathbf{x}'\beta$ og $h(\mathbf{x}) = \sigma$, altså at vi varians ikke avhenger av \mathbf{x} . Det er ingenting i veien for at vi modeller hvordan varians avhenger av \mathbf{x} , men ofte er vi bare interessert i betinget forventningsverdi.⁹ Vi kan da skrive opp likelihood-funksjonen.

$$L(\beta, \sigma) = \prod_n f(\mathbf{x}_n, y_n) = \prod_n \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y_n - \mathbf{x}_n'\beta)^2}{2\sigma^2}\right\} \cdot f(\mathbf{x}_n) \quad (6.76)$$

Kan merke da at parameter i betinget fordeling ikke avhenger av $f(\mathbf{x}_n)$ slik at vi kan se bort i fra dette.. og vise at OLS maksimerer likelihood. hm.

6.4.4 Betinget bernoulli

Vi vil se på hvordan ulike egenskaper x til en person påvirker sannsynlighet for at hen deltar i arbeidsmarkedet.

$$q(s) = \mathbb{P}\{y = 1|x = s\} \quad (6.77)$$

⁹Tror dette er eksempel på semi-parametrisk estimering der σ er såkalt nuisance-parameter.

For at funksjonen skal tilfredstille aksiom til sannsynlighetsfunksjoner må $q(s) \in [0, 1] \forall s \in \mathbb{R}^k$. En type funksjoner som tilfredstiller det kaller vi kumulative sannsynlighetsfunksjoner. I praksis bruker vi derfor cdf med lineær parametrisering, $q(s) = F(s'\beta)$. Kan ikke estimere det med OLS siden det er en ikke-lineær funksjon mhp parametrene. For å gjøre MLE operativt må vi ha en spesifisert log likelihood funksjon som vi kan optimere. Første steg er betinget pmf. I utgangspunktet er det en piecewise funksjon, men jeg kan bruke triks for å få det på én linje:

$$P(y = i | x = s) = F(s'\beta)^i (1 - F(s'\beta))^{1-i} \quad (6.78)$$

$$\implies \log L(\beta) = \sum y_n \log(F(s'\beta)) + \sum (1 - y_n) \log(1 - F(s'\beta)) \quad (6.79)$$

Dette kan jeg løse og få logit eller probit avhengig av valg av F . Ble litt ukomfortabel notasjon fordi jeg ikke vil bruke store bokstaver, men skal helst sikkert se nærmere på dette senere.

6.5 Prinsipp for å utlede tester

Vi vil ofte teste om data i utvalg gir tilstrekkelig bevis til at vi kan forkaste påstand om at $\theta \in \Theta_0$. Vi forkaster dersom det er lite sannsynlig at de faktiske observasjonene i utvalget har blitt generert fra en fordeling med parameter fra nullhypotesen. Videre gjør vi ofte avgrensinger av hvilke fordelinger \mathcal{P} vi vil betrakte. Det er da nyttig å gjøre spesifikasjonstester for se om vi kan forkaste at $\mathbb{P} \in \mathcal{P}$ slik at modell er feilspesifisert, for eksempel ved at feilledd er heteroskedastisk. Har tre ulike prinsipper for å utlede testobservator fra likelihoodfunksjon som er asymptotisk ekvivalente og alle gir χ^2 -fordelte testobservator. Tror at at *t-ogF-fordeling* bare er justering som tar hensyn til at utvalg er begrenset, men litt usikker på dette.

6.5.1 Greier fra DM

Prinsipp for å utlede testobservator som asymptotisk er $\chi^2(r)$ -fordelt, der r er antallet restriksjoner. Husk at nullhypotese er restriksjon av hypoteserommet og impliserer en avgrenset modell. Tror kanskje vi kan bruke parametrisert bootstrap som alternativ..

Wald bruker bare uavgrenset modell. Gitt at hypotesen er sann er $r(\hat{\theta}) \sim N(0, V(r(\hat{\theta})))$ slik at $r(\hat{\theta})'[V(r(\hat{\theta}))]^{-1}r(\hat{\theta}) \sim \chi^2(r)$.

6.5.2 Wald-test

Denne fremgangsmåten tar utgangspunkt i en (asymptotisk) normalfordelt MLE estimator $\hat{\theta}$ og bruker at

$$\mathbf{R}\hat{\theta} - \mathbf{q} \xrightarrow{d} N(\mathbf{0}, \mathbf{R}\Sigma\mathbf{R}') \quad (6.80)$$

hypotesen $\mathbf{R}\theta = \mathbf{q}$ er sann.. Kan da forkaste nullhypotese hvis testestimator gir stort avvik fra 0, der vi bruker kvantiler av normalfordeling til å konkretisere hva som utgjør tilstrekkelig stort avvik for våre formål. Med enkelt parameter har testen form

$$W = \frac{\hat{\delta} - \delta_0}{\hat{se}} \xrightarrow{d} N(0, 1) \quad (6.81)$$

Med flere parametre tror jeg den fremgangsmåten gir flervariabel normalfordeling, men det er mye greiere å få tilbake et enkelt tall slik at vi kan forkaste hvis langt fra null. Alle tre prinsippene for å utlede asymptotiske tester gir oss derfor generelt testobservatorer som er χ^2 -fordelt. For den enkle testen over gir det

$$\xi_w = Z^2 = (\hat{\delta} - \delta_0)[v\hat{ar}]^{-1}(\hat{\delta} - \delta_0) \xrightarrow{d} \chi^2(1) \quad (6.82)$$

I praksis er vi ofte interessert i forskjell mellom parametre siden vi ikke har et kjent benchmark vi tester mot. Det er relativt greit dersom utvalgene er uavhengige av hverandre siden varians til differansen av estimatorer er sum av variansen til hver av de. Et eksempel er forskjell i parameter i bernoulli-fordelt variabel. Vi har $\bar{X}_1 \sim \widehat{binom}(p_1, n_1)$ og $\bar{X}_2 \sim binom(p_2, n_2)$. Da er $\hat{\delta} = \hat{p}_2 - \hat{p}_1$. Gjenstår bare å finne $\hat{se} := \widehat{se}(\hat{\delta})$. Vet at \hat{p}_j gjennomsnitt. Vet at varians til gjennomsnitt er $\frac{\sigma^2}{N}$. Vet at $\sigma = p(1-p)$ i bernoulli, som er fordeling til X . Dette er tilstrekkelig til å finne \hat{se} , men gidder ikke skrive. Ta det som oppgave når du leser dette. Tilsvarende kan vi teste differanse mellom normalfordelte. En utvidelse et t-test i stedet for z-test.

6.5.3 Likelihood ratio

Likelihood ratio tar utgangspunkt i forskjellen i maksimum av log-likelihood fra ubetinget og betinget optimering.

$$\lambda = 2 \log \left(\frac{\sup_{\theta \in \Theta} L(\theta)}{\sup_{\theta \in \Theta_0} L(\theta)} \right) = 2 \left(\frac{L(\hat{\theta})}{L(\hat{\theta}_0)} \right) \quad (6.83)$$

der $\hat{\theta}$ er MLE-estimatoren og $\hat{\theta}_0$ er MLE-estimator avgrenset til $\Theta = \Theta_0$. Dette har visst en χ^2 fordeling. Intuisjon for dette er at hvis forskjellen er stor så er det lite sannsynlig at avgrensingen ikke er bindene, altså at lite sannsynlig at $\theta \in \Theta_0$

6.5.4 Lagrange multiplier

Nullhypotesen medfører en restriksjon av parameterrommet. Undersøker i hvilken grad restriksjon er bindene ved å se på lagrangemultiplier assosiert med restriksjonen av de ulike parameterne, skyggepris. Hvis stor skyggepris er lite sannsynlig at sann parameter i Θ_0 som vi har avgrenset til å velge løsning innenfor..

6.6 Egenskaper ved feilspesifikasjon

Det er ganske sterk antagelse at $P_0 \in \mathcal{P}$, så kjekt at ikke alt rakner dersom denne antagelsen er feil.

6.6.1 Total variation distance og KL-divergence

Jeg vil ha et mål på avstand mellom to probability measures $D(\mathbb{P}_{\theta_0}, \mathbb{P}_{\theta})$. Et ganske naturlig mål er *Total variation distance*

$$TV(\mathbb{P}_{\theta_0}, \mathbb{P}_{\theta}) = \max_{A \subset \Omega} |\mathbb{P}_{\theta_0}(A) - \mathbb{P}_{\theta}(A)| \quad (6.84)$$

Dette gir et tall i intervallet $[0, 1]$ og tilfredstiller egenskapene til en avstand (symmetrisk, trekantulikheter..). Hvis vi ser på fordelinger på \mathbb{R} kan vi beregne størrelsen med

$$TV(P_{\theta_0}, P_{\theta}) = \begin{cases} \frac{1}{2} \sum |p_{\theta_0}(x) - p_{\theta}(x)| \\ \frac{1}{2} \int |f_{\theta_0}(x) - f_{\theta}(x)| dx \end{cases} \quad (6.85)$$

Dette tilsvarer areal av mellom kurvene i området der den ene er større enn den andre. Det er symmetri siden areal under begge kurvene summerer til 1. Dette er et naturlig mål med gode egenskaper, men det litt vanskelig å gjøre operativt. Dette motiverer Kullback-Leibler (KL) divergence som har noe av de samme gode egenskapene, men som vi kan estimere fra utvalg med observasjoner fra P_{θ_0} .

$$KL(P_{\theta_0}, P_{\theta}) = \begin{cases} \sum p_{\theta_0}(x) \log \left(\frac{p_{\theta_0}(x)}{p_{\theta}(x)} \right) \\ \int f_{\theta_0}(x) \log \left(\frac{f_{\theta_0}(x)}{f_{\theta}(x)} \right) dx \end{cases} \quad (6.86)$$

Dette målet er ikke symmetrisk og tilfredstiller ikke triangelulikheter, men er i likhet med TVD 0 når funksjonene er like og vokser når avstanden øker. Selve tallet har ikke naturlig tolkning. Merk at dette er forventningsverdi av en funksjon med hensyn på P_{θ_0} .

$$KL(P_{\theta_0}, P_{\theta}) = E_{\theta_0}[\log(f_{\theta_0})] - E_{\theta_0}[\log(f_{\theta})] \quad (6.87)$$

Et naturlig valg av θ er den verdien som minimerer den empiriske analogen til KL-divergence, $\widehat{KL}(P_{\theta_0}, P_\theta)$. Merk at første ledd er konstant som ikke påvirker arg min.

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \widehat{KL}(P_{\theta_0}, P_\theta) \quad (6.88)$$

$$= \arg \min_{\theta \in \Theta} \{E_{\theta_0}[\log(f_{\theta_0}(x_n))] - E_{\hat{P}_N}[\log(f_\theta(x_n))]\} \quad (6.89)$$

$$= \arg \min_{\theta \in \Theta} \text{konstant} - \frac{1}{N} \sum \log(f_\theta(x_n)) \quad (6.90)$$

$$= \arg \max_{\theta \in \Theta} \log(\Pi_n f_\theta(x_n)) \quad (6.91)$$

$$= \arg \max_{\theta \in \Theta} \log(\Pi_n L_n(\theta; x_n)) \quad (6.92)$$

Løsningen på dette optimeringsproblemet tilsvarer MLE-estimatoren. Så langt har vi antydnet at den sanne fordelingen tilhører den parametriske klassen som vi søker over, men siden $E_{\theta_0}[\log(f_{\theta_0}(x_n))]$ uansett bare er en konstant som vi kan se bort i fra i optimeringen, så sier dette resultatet oss at MLE asymptotisk gir oss θ som korresponderer med den fordelingen innenfor \mathcal{P} som minimerer avstand til den sanne P , både i betydningen av TVD og KL. Kunne jo selvsagt kommet nærmere ved å søke over en riktig spesifisert \mathcal{P} , men det er jo uansett et greit resultat.

6.6.2 MLE fra empirisk risikominimering

Kan vise at MLE-estimatoren kan utledes som spesialtilfelle av empirisk risikominimering.¹⁰ Anta at vi vil estimere en ukjent tetthetsfunksjon $q(\cdot)$ med utgangspunkt i observerte realiseringer fra fordeling med den tettheten. Definerer tapsfunksjon til kandidat $p(\cdot)$ ved $L(p, x) := -\log(p(x))$. Hvis vi observerer realisering $x = s$ så vil det realiserste tapet være større desto lavere verdi av $p(s)$, altså jo lavere tyngde vår kandidat plasserer på den realiserste verdien. Risikofunksjonen er dermed gitt ved

$$R(p) = \mathbb{E}_q[L(p, x)] = - \int \log(p(s)) ds \quad (6.93)$$

For å knytte dette til MLE avgrensers vi til å betrakte et parametrisert hypoteserom $\mathcal{P}_\theta = \{P_\theta : \theta \in \Theta\}$. Antar at modellen er identifisert slik at $\theta \mapsto P_\theta$ er bijektiv (én-til-én)

¹⁰I hvert fall for estimering av tetthetsfunksjon... skal se om jeg kan utvide til regresjon.

slik at vi ekvivalent kan løse minimeringsproblemet med hensyn på θ .

$$P_{\hat{\theta}} = \arg \min_{p \in \mathcal{P}_{\theta}} R_{emp}(p) \quad (6.94)$$

$$\implies \hat{\theta} = \arg \min_{\theta \in \Theta} - \sum \log(p_{\theta}(s_n)) \quad (6.95)$$

$$\implies \hat{\theta} = \arg \max_{\theta \in \Theta} \sum \log(p(\theta; s_n)) \quad (6.96)$$

Ser at løsningen tilsvare $\hat{\theta}_{MLE}$ når vi bruke $-\log(p(s))$ som tapsfunksjonen. Vi kan også dekomponere risikoen assosiert med denne tapsfunksjonen for å knytte det til KL-divergence,

$$R(p) = - \int \log(p(s)) q(s) ds \quad (6.97)$$

$$= \mathbb{E}_q[(-\log(p(s)) + \log(q(s)) - \log(q(s)))] \quad (6.98)$$

$$= \mathbb{E}_q \left[\log \left(\frac{q(s)}{p(s)} \right) \right] - \mathbb{E}_q[\log(q(s))] \quad (6.99)$$

der første ledd er KL-divergence og andre ledd er entropy. Ser at entropy tilsvare risiko når fordelingen er kjent. Får tilbake igjen MLE-estimatoren ved å minimere utvalgsanalogen til KL-divergence med hensyn på parametrisert p som beskrevet i seksjonen over.

Liten digresjon, må endres eller slettes

Kan utlede estimering av logistisk regresjon m.m. ved å bruke såkalt logistisk tap,

$$-\{y \log h_{\theta}(x) + (1 - y) \log(h_{\theta}(x))\} \quad (6.100)$$

men denne tapsfunksjonen kan jeg jo uansett utlede fra loglikelihood til bernoulli-fordelingen. Det er jo litt poeng at jeg kan motivere dette uten MLE, men er uansett bedre å gjøre det innenfor.

6.6.3 Kvasi-MLE

Forventningsverdien til score-funksjonen evaluert i sann parameter er 0.

$$\mathbb{E}_{\theta_0} S(\theta_0) := \int S(\theta_0) f(x; \theta_0) dx = 0 \quad (6.101)$$

$$(6.102)$$

Utvalgsanalogen er å finne $\hat{\theta}$ som gjør at $\mathbb{E}_{\hat{P}_N}[S(\theta)] = 0$. Dette gir oss et ligningssystem av momentbetingelser som vi kan løse og estimatoren kan betraktes som en momentestimator. Hvis modellen er riktig spesifisert er estimatorene ekvivalente, men egenskapene til

momentestimatoren vil være gyldig for en større klasse av fordelinger som har de samme første-ordens betingelsene. Hvis vi bruker F.O.B fra MLE til å betrakte fordelingen til estimatoren fra denne utvidede mengden kan vi betegne det som kvasi-likelihood. Fra store talls lov vet vi at estimatoren er konsistent og fra CLT får vi normalfordelingen, men asymptotisk varians er ikke lenger $I(\theta)^{-1}$. Vi må bruke såkalt sandwich estimator,

$$(\hat{\theta}_{QMLE} - \theta_0) \xrightarrow{d} N(0, V) \quad (6.103)$$

der

$$V = \dots \quad (6.104)$$

Vil knytte dette til robust standardfeil i regresjon..S

6.6.4 Extremum estimators

Kan betrakte en klasse av estimator som er løsning på et optimeringsproblem av en objektfunksjon $Q_N(\cdot)$ som har verdimengde i \mathbb{R} ,

$$\hat{\theta} = \arg \max_{\theta \in \Theta} Q_N(\theta; \mathbf{z}_d) \quad (6.105)$$

En viktig underklasse er M-estimatorer der $Q_N(\cdot)$ har formen

$$Q_N(\theta; \mathbf{z}_d) = \frac{1}{N} \sum_n m(\theta; \mathbf{z}_n). \quad (6.106)$$

Et annet eksempel er GMM der

$$Q_N(\theta; \mathbf{z}_d) = -\frac{1}{2} g_N(\theta; \mathbf{z}_d)' \hat{W} g_N(\theta; \mathbf{z}_d) \quad (6.107)$$

der $g_N(\theta; \mathbf{z}_d) := \frac{1}{N} g(\theta; \mathbf{z}_n)$ og $g(\cdot)$ er momentbetingelse. Tror vi kan bruke denne mer generelle klassen av estimatorer til å dra koblinger mellom MLE og GMM, blant annet type *Limited information maximum likelihood* og finne asymptotiske tester for GMM. Jeg mistenker at de greiene der må bli i neste liv.

Kapittel 7

Lineær regresjon

Noen av fordelene med å avgrense til lineære funksjoner er at vi kan få parametre med enkle tolkninger og det empiriske risikominimeringsproblemet med kvadratisk tap har en analytisk løsning (MKM). Kjenner teoretisk egenskap, stabil, prediksjonsrisiko. Stor klasse av additative modeller... tenker det finnes parametrisk representasjon av splines og lignende... litt usikker på hva jeg sier om dette. Nært knyttet til MKM, men dette er bare én av flere måter å estimere koeffisientene. (Vekte observasjonene i tapsfunksjonen ut fra varians... legge til regularisering... si noe om IV?).

Lineær regresjon bruker i ulike fagfelt, for ulike formål og kan motiveres på ulike måter. Beste tilnærmede løsning på et overdeterminert ligningssystem. I økonometri motiveres det gjerne av såkalt *conditional independence assumption* der behandling er tilfeldig fordelt innad i delgrupper og kan betraktes som analog til *matching estimator*. I statistisk modellering (MLE/Bayes) gjør vi eksplisitte antagelser om parametrisert struktur til fordeling som generer data. I ren prediksjonssetting kan vi være mer agnostisk om fordeling og bare løser risikominimeringsproblem. Jeg begynner med siste tilnærming.

7.1 Egenskap ved simultanfordeling

beskrive egenskap ved simultanfordeling. ting vi kan forsøke å lære fra realiserte observasjoner

Vi kan betegne den betingede forventningsfunksjonen $E[Y|X] := \int yf(y|X)dy$ som (populasjons) regresjonsfunksjonen.¹ Dette er en tilfeldig variabel som for hver $X = x$ angir forventningsverdi til den betingede fordelingen av y . Det er den ortogonale projeksjonen av y ned på underrommet som består av alle tilfeldige variabler som kan skrives som en deterministisk funksjon av X . Dette medfører at det minimerer forventet avvik

¹Merk at populasjonsregresjonsfunksjonen (PRF) også brukes som betegnelse på den beste lineære tilnærmingen. Jeg tenker det er bedre å betegne dette som den lineære populasjonsregresjonsfunksjonen.

og at vi kan dekomponere

$$Y = E[Y|X] + Y - E[Y|X] = E[Y|X] + U \quad (7.1)$$

der feilledet U per konstruksjon er ortogonal på alle funksjoner av X , altså $E[g(X)U] = 0$. I praksis bruker vi små bokstaver av notasjonell konvensjon.

Den lineære populasjonsregresjonsfunksjonen tilsvarer den ortogonale projeksjonen av y ned på mengden av tilfeldige variabler som kan skrives som en lineær funksjon av x . Det medfører at vi kan dekomponere

$$y = \beta'x + y - \beta'x = \beta'x + u \quad (7.2)$$

der feilledet u per konstruksjon er ortogonal på alle lineære funksjoner av x , altså $E[(b'x)u] = b'E[xu] = 0$. Dersom det inkluderer et konstantledd så medfører det også at $E[1u] = E[u] = 0$ og $cov(x_k, u) = 0$. Vi kan også motivere dette som beste lineære tilnærming til $E[y|x]$ som er det vi egentlig er interessert i. Husk at dette er en tilfeldig variabel, så det er ingenting i verien for å betrakte det som den avhengige variabelen i regresjonen.

7.1.1 Projeksjon

7.1.2 Dekomponering av varians

knytte noe til avstand, geometri.. mest mulig analog

7.1.3 Tolkning av feilledd

avvik, $E[u|x]$, strukturelt eller ikke. tolkning av parameter.

7.2 Numeriske egenskaper

Minste kvadrats metode har en rekke egenskaper som gjelder for alle datasett. Disse algebraiske egenskapene følger av projeksjoneringen og har geometriske tolkningner.

7.2.1 Ortogonal projeksjon

Minimeringsproblemet med kvadratisk tapsfunksjon er

$$h_{\hat{b}} = \arg \min_{h_b \in \mathcal{H}_l} \mathbb{E}_{P_{\hat{N}}}[(y_n - h_b(\mathbf{x}_n))^2] \quad (7.3)$$

$$\hat{\mathbf{b}} = \arg \min_{b \in \mathbb{R}^K} \frac{1}{N} \sum_n (y_n - \mathbf{x}'_n \mathbf{b})^2 \quad (7.4)$$

Generelt må vi minimere tapsfunksjonen numerisk ved å bruke algoritme som søker over parameterromet. Her kan vi løse det analytisk.² Vi begynner med å sette det opp på matriseform ved å stappe input-vektorene,

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_N \end{bmatrix} \quad (7.5)$$

slik at $col_k(\mathbf{X})$ gir verdi av feature k til hver av de N observasjonene i utvalget. Omskriver tapsfunksjon³,

$$\frac{1}{N} \sum_n (y_n - \mathbf{x}'_n \mathbf{b})^2 = \frac{1}{N} (\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b}) \quad (7.6)$$

Kan knytte dette til avstand og ortogonal projeksjon. Uansett, finner

$$\hat{\mathbf{b}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (7.7)$$

7.2.2 Frisch-Waugh-Lovell

Vi projekterer \mathbf{y} på $S(\mathbf{X})$.

$$\mathbf{y} = \mathbf{X}\hat{\beta} + \mathbf{M}\mathbf{y} \quad (7.8)$$

$$= \mathbf{X}_1\hat{\beta}_1 + \mathbf{X}_2\hat{\beta}_2 + \mathbf{M}\mathbf{y} \quad (7.9)$$

FWL-theoremet sier at vi får samme $\hat{\beta}_2$ ved å projekte \mathbf{y} på residualen av projeksjonen av \mathbf{X}_2 på $S(\mathbf{X}_1)$. Merk at matrise kan transformere flere vektorer om gangen, men når jeg snakker om residualer så begrenser jeg implisitt til å se på én vektorer.. Kan få noe intuisjon ved å tenke på bivariat regresjon. Hvis jeg projekterer \mathbf{x} på $S(\mathbf{1})$ så får jeg $\bar{x}\mathbf{1}$. Residualen er da den sentrerte vektoren der hver komponent er avvik fra gjennomsnitt. Hvis jeg regger \mathbf{y} på den sentrerte variabelen uten konstantledd får jeg samme helning som i den bivariate regresjonen. Hm.

FLW-theoremet gjør at vi kan isolere enkeltkomponentner i helningskoeffisienten og betrakte det som analog bivariat regresjon.

$$\beta_1 = \frac{cov(y, \tilde{\mathbf{x}}_k)}{var(\tilde{\mathbf{x}}_k)} \quad (7.10)$$

der $\tilde{\mathbf{x}}_k$ er residualen fra regresjonen av \mathbf{x}_k på $S(\mathbf{x}_{-k})$.

²Det betyr at vi kan skrive $\hat{b} = g(\{(y_n, \mathbf{x}_n) : n = 1, \dots, N\})$ der vi kjenner g . I praksis er det bedre å organisere utvalget i matrise som er helt analogt til representasjon i tabulær form som vi er vant til å se.

³Lurer på om jeg vil ha eget begrep for tap på hele utvalget i stedet for enkeltobservasjon... kostnad?

Tror det er enklest å tenke på dette som en to stegs prosess

$$\mathbf{M}_1 \mathbf{y} = \mathbf{M}_1 [\mathbf{X}_1 \hat{\beta}_1 + \mathbf{X}_2 \hat{\beta}_2 + \mathbf{M} \mathbf{y}] \quad (7.11)$$

$$= \mathbf{M}_1 \mathbf{X}_2 \hat{\beta}_2 + \mathbf{M} \mathbf{y} \quad (7.12)$$

Merk at $S(\mathbf{M}) \subset S(\mathbf{M}_1)$. Hjelper dette..?

Uansett, tror mye av poenget er at

$$\hat{\beta}_k = \frac{\text{cov}(\tilde{\mathbf{x}}_k, \tilde{\mathbf{y}}_k)}{\text{var}(\tilde{\mathbf{x}}_k)} \quad (7.13)$$

der $\tilde{\mathbf{x}}_k$ og $\tilde{\mathbf{y}}_k$ er residualen av projeksjon av henholdsvis \mathbf{x}_k og \mathbf{y} på span av underrommet til de andre forklaringsvariablene, $S(\mathbf{X}_{-k})$.

7.2.3 In-sample fit

Vi dekomponerer \mathbf{y} i komponent i $\text{span}(\mathbf{X})$ og dets ortogonale komplement. Vi vil si noe om hvor god vår tilnærmede løsning er. Det avhenger den relative størrelsen på komponentene; hvor stor avviket er i forhold til *størrelsen* på \mathbf{y} . Alt dette er jo vektorer så bedre å snakke om lengde enn størrelse...

- $TSS = \|\mathbf{y}\|^2$
- $RSS = \|\mathbf{M} \mathbf{y}\|^2$
- $ESS = \|\mathbf{P} \mathbf{y}\|^2$

fra pythagoras følger det at $TSS = RSS + ESS$. Vi kan definere enkel R^2 som andelen av den totale lengden som går i retning av kolonnerommet til \mathbf{X} ...

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS} = 1 - N \frac{R_{emp}(\hat{b})}{TSS} \quad (7.14)$$

Det følger at $R^2 \in [0, 1]$ og at det gir mål på in-sample fit.⁴ Men målet vårt er å generalisere til nye data; ikke lære mest mulig om det gitte utvalget vi besitter. Maksimering av R^2 er dårlig kriterie i modellselskjon siden vi alltid kan få den høyere ved å legge til nye variabler enten de er relevante eller ikke,

$$\text{span}(\mathbf{X}_a) \subset \text{span}(\mathbf{X}_b) \implies R_a^2 \leq R_b^2 \quad (7.15)$$

Det er heller ikke problem i seg selv at R^2 er lav... hvis vi estimerer kausal sammenheng så kan det være at eksponering for behandling forklarer liten andel av variasjon i utfall. Kan medføre problem med presisjon til koeffisientestimatene, men ser på dette under statistisk egenskaper.

⁴Hvis vi bruker annen estimering enn OLS og velger helt arbitrære R^2 så kan vi få negative verdier.

7.3 Statistiske egenskaper

De statistiske egenskapene sier oss noe om fordelingen til estimerte egenskaper dersom vi observerte (uendelig) mange realiserte datasett \mathcal{D} fra samme fordeling. For å utlede slike egenskaper må vi gjøre antagelser om prosessen som genererer data. Disse antagelsene kan til dels sannsynliggjøres fra vårt enkle utvalg, men det kan ikke bevises. Med antagelsene kan vi utlede egenskaper.

Antagelsene vi gjør avgrenser hvilke datagenereringsprosesser (DGP) vi er villig til å betrakte. Det definerer implisitt en mengde \mathbb{M} av DGPer. Hvis den sanne prosessen $DGP_0 \in \mathbb{M}$ er modellen riktig spesifisert og egenskapene vi utleder er sanne for prosessen vi betrakter. En (historisk) mye brukt modell er den klassiske lineære modellen som lar oss utlede statistiske egenskaper for vilkårlig utvalgsstørrelse

$$\mathbb{M} = \{y = \mathbf{x}\beta + u : u \sim NID(0, \sigma^2)\} \quad (7.16)$$

$$DGP_0 = \{y = \mathbf{x}\beta_0 + u : u \sim NID(0, \sigma_0^2)\}. \quad (7.17)$$

Merk at hvis vi inkluderer en irrellevante variabel z i modellen som ikke er inkludert i DGP_0 så vil det likevel være element i \mathbb{M} med koeffisient lik 0 for z . Egenskapene holder likevel. Hvis det derimot er utelatt variabel er $DGP_0 \notin \mathbb{M}$ og modellen er feilspesifisert.

7.3.1 Små utvalg

$$E[\hat{\beta}|\mathbf{X}] = E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'y|\mathbf{X}] \quad (7.18)$$

$$= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \mathbf{u})|\mathbf{X}] \quad (7.19)$$

$$= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E[\mathbf{u}|\mathbf{X}] \quad (7.20)$$

$$= \beta \quad (7.21)$$

$$V(\hat{\beta}) = E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'] \quad (7.22)$$

$$= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}'] \quad (7.23)$$

$$= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}] \quad (7.24)$$

$$= \sigma^2 E[(\mathbf{X}'\mathbf{X})^{-1}] \quad (7.25)$$

der jeg har brukt at $u_n \sim NID(0, \sigma^2) \implies V(\mathbf{u}) = E[\mathbf{u}\mathbf{u}'] = \sigma^2\mathbf{I}$. Kan omskrive

$$V(\hat{\beta}) = \frac{\sigma^2}{N} E[(\frac{1}{N}\mathbf{X}'\mathbf{X})^{-1}] \quad (7.26)$$

Variansen vokser proporsjonalt med varians til feilledd, omvendt proporsjonalt (?) med utvalgsstørrelse og avhenger også av kovariansstruktur mellom uavhengige variabler. Kan

bruke FWL til å finne uttrykk for variansen til enkelt koeffisientestimator,

$$V(\hat{\beta}_k) = \sigma^2(\mathbf{X}'\mathbf{M}_{-k}\mathbf{X})^{-1} \quad (7.27)$$

Litt usikker på den der. Avhenger av hvor mye variasjon i x_k som ikke er forklart med variasjonen i andre variabler...

7.3.2 Store utvalg

7.3.3 Presisjon til koeffisient

7.3.4 Presisjon til prediksjon

Antar at jeg har en respons som er betinget normalfordelt og at CEF er lineær.

$$\mathbf{y}|\mathbf{X} \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I}) \quad (7.28)$$

Hvis jeg observerer én realisering av denne simultanfordelingen får jeg én realisering av $\hat{\beta}$. Jeg kan beregne utvalgsfordelingen til $\hat{\beta}$ gitt antagelsen jeg har gjort over,

$$\hat{\beta} \sim N(\beta, \sigma^2\mathbf{I}) \quad (7.29)$$

Hvis jeg får en ny observasjon \mathbf{x}_{new} så vil min prediksjon fra den estimerte modellen være $\hat{y}_{new} := \mathbf{x}'_{new}\hat{\beta}$ siden dette er sentraltendensen i betinget fordeling av $y|\mathbf{x}_{new}$. På en annen side er det jo slik at dersom jeg hadde observert en annen realisering av simultanfordelingen ville jeg hatt en annen $\hat{\beta}$ og gitt annen prediksjon. Jeg vil forsøke å kvantifisere variasjonen i prediksjonen.⁵

$$V(\hat{y}_{new}|\mathbf{X}) = V(\mathbf{x}'_{new}\hat{\beta}|\mathbf{X}) \quad (7.30)$$

$$= E[(\mathbf{x}'_{new}\hat{\beta})^2|\mathbf{X}] - E[\mathbf{x}'_{new}\hat{\beta}|\mathbf{X}]^2 \quad (7.31)$$

$$= E[\mathbf{x}'_{new}\hat{\beta}\hat{\beta}'\mathbf{x}_{new}|\mathbf{X}] - (\mathbf{x}'_{new}\beta)^2 \quad (7.32)$$

$$= \mathbf{x}'_{new}E[\hat{\beta}\hat{\beta}'|\mathbf{X}]\mathbf{x}_{new} - \mathbf{x}'_{new}\beta\beta'\mathbf{x}_{new} \quad (7.33)$$

$$(7.34)$$

Bruker nå at

$$E[\hat{\beta}\hat{\beta}'|\mathbf{X}] = V(\hat{\beta}|\mathbf{X}) + E[\hat{\beta}|\mathbf{X}]E[\hat{\beta}|\mathbf{X}]' \quad (7.35)$$

$$= \sigma^2(\mathbf{X}'\mathbf{X})^{-1} + \beta\beta' \quad (7.36)$$

⁵Jeg tror målet mitt er å kunne angi et intervall der jeg med gitt sannsynlighet kan påstå at y_{new} vil ligge i. Det avhenger både av variasjon i \hat{y}_{new} og fordeling til avvik fra sentraltendens... Begynner i hvertfall med å se på variasjon til prediksjonen.

slik at

$$V(\hat{y}_{new}|\mathbf{X}) = \mathbf{x}'_{new}(\sigma^2(\mathbf{X}'\mathbf{X})^{-1} + \beta\beta')\mathbf{x}_{new} - \mathbf{x}'_{new}\beta\beta'\mathbf{x}_{new} \quad (7.37)$$

$$= \sigma^2\mathbf{x}'_{new}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_{new} + \mathbf{x}'_{new}\beta\beta'\mathbf{x}_{new} - \mathbf{x}'_{new}\beta\beta'\mathbf{x}_{new} \quad (7.38)$$

$$= \sigma^2\mathbf{x}'_{new}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_{new} \quad (7.39)$$

Tror det bør samsvare med variasjon jeg observerer når jeg sampler verdier av $\hat{\beta}$ og plotter linjer.. det vet jeg ikke om det gjør ..

7.3.5 Residual

Bruker residual til å estimere varians til feilledd,

$$\hat{\sigma}^2 = \frac{1}{N} \sum \hat{u}_n^2 \quad (7.40)$$

$$= \frac{1}{N} \sum (y_n - \mathbf{x}'_n \hat{\beta})^2 \quad (7.41)$$

$$= \frac{1}{N} (\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta}) \quad (7.42)$$

$$= \frac{1}{N} (\mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\hat{\beta} + \hat{\beta}'\mathbf{X}'\mathbf{X}\hat{\beta}) \quad (7.43)$$

der

$$\hat{\beta}'\mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (7.44)$$

$$= \mathbf{y}'\mathbf{X}\hat{\beta} \quad (7.45)$$

slik at

$$\hat{\sigma}^2 = \mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbf{X}\hat{\beta} \quad (7.46)$$

7.4 Hypotesetester

Knytte til test-prinsipp fra MLE...

7.4.1 Lineære restriksjoner av koeffisient

tror jeg vil ha t-test som special case av F-test

7.4.2 Diagnose

funksjonell form. vet ikke helt hva annet jeg vil teste..

7.5 Funksjonell form

Vi forsøker å tilnærme oss funksjonen $\mathbb{E}[y|\mathbf{x}]$. Vi har sett litt på i hvilke tilfeller forskjellene i observert utfall for ulike nivå av *behandling* kan ha kausal tolkning. Uansett er jeg interessert i å beskrive hvordan funksjonen endrer seg når variabler endres. I praksis bruker vi lineær regresjon. Det kan håndtere ikke-linearitet ved å først transformere variablene. Tolke den estimerte funksjonen; endring i forhold til opprinnelige størrelser.

Logaritmer

Vi kan transformere den avhengige variabelen. Mange utfall er strengt positive, så de kan ikke feilbedd være normalfordelt siden det er begrenset nedenfra. OLS er også veldig sensitiv for ekstremverdier, så greit å få skalert ned høye numeriske verdier av utfallet. Dessuten vil det ofte være slik at gjennomsnittlig utfall vokser omtrent proporsjonalt med variablene i stedet for lineært. Skal nå se hvordan vi tolker koeffisienter der vi har tatt logaritmisk transformasjon av forklaringsvariabel og/eller utfall.

Vi bruker at logaritmer har omtrent prosentvis tolkning for små endringer. Dette følger av første-orders Taylor-ekspansjon av logaritmen evaluert i $x = 1$,

$$\log(x) \approx \log(1) + \frac{d}{dx}\log(x)|_{x=1}(x-1) = x-1, \quad (7.47)$$

slik at $\Delta \log(x) := \log(x_1) - \log(x_0) = \log(x_1/x_0) \approx x_1/x_0 - 1 := \Delta x/x_0$.

Log-level

Har modell $\log(y) = \beta_0 + \beta_1 x + u$, $E(u|x) = 0$. Det følger at $\beta_1 = \frac{d}{dx}E[\log(y)]$. Okay, men vi er interessert i hvordan det endrer forventningsverdi til y . Observerer at

$$\Delta E \log(y) = \beta_1 \Delta x \quad (7.48)$$

$$\implies 100 \cdot \beta \approx 100 \cdot E \frac{\Delta y}{y_0} := \%E \Delta y \quad (7.49)$$

Det kan tolkes som prosentvis endring i forventet utfall når x endres med én enhet. Hvis jeg vil finne eksakt prosentvis endring så kan jeg bruke

$$\Delta E \log(y) = \beta_1 \Delta x \quad (7.50)$$

$$\exp(\Delta E \log(y)) = \exp(\beta_1 \Delta x) \quad (7.51)$$

$$\Delta E \log(y) = \exp(\beta_1 \Delta x) \quad (7.52)$$

$$\% \Delta E \log(y) = 100 \cdot (\exp(\beta_1 \Delta x) - 1) \quad (7.53)$$

Det er nødvendig å være litt forsiktig når vi bruker denne spesifikasjonen til å predikere verdi av y . Det kan være fristende å bruke

$$\widehat{\log y} = \mathbf{x}'\hat{\beta} \quad (7.54)$$

$$\hat{y} = \exp(\mathbf{x}'\hat{\beta}) \quad (7.55)$$

men dette er dårlig estimat på $E[y|x]$. For å se dette, observer at

$$E[y|x] = E[\exp(\mathbf{x}'\beta + u)|\mathbf{x}] = \exp(\mathbf{x}'\beta)E[\exp(u)] \quad (7.56)$$

der $E[e^u] = \alpha_0 \neq 0$. For å bruke spesifikasjonen til å predikere verdi av y må vi skalere opp $\hat{y} = \mathbf{x}'\hat{\beta}$ med $\hat{\alpha}_0 = E_{\hat{P}_N}[e^u] = \frac{1}{N} \sum_n e^{\hat{u}_n}$

Log-log

Kan toles som elastisitet. Følger av at

$$\Delta E \log(y) = \beta_1 \Delta \log(x) \quad (7.57)$$

$$\implies \beta_1 \approx \% \Delta E y / \% \Delta x \quad (7.58)$$

7.6 Fordeling til feilledd

Vi trenger bare momentbetingelsene for å konsistent estimere helningskoeffisientene og de er asymptotisk normalfordelte fra CLT. I utgangspunktet trenger vi ikke bry oss så mye om fordelingen til $\mathbf{u} := [u_1, u_2, \dots, u_N]'$ dersom vi kun vil ha mål sentraltendens i betingede fordelinger. Når vi estimerer med MLE antar vi gjerne at $\mathbb{V}[u|x] = \sigma^2$ og at observasjon er iid slik at $\mathbb{V}[\mathbf{u}|X] = \sigma^2 I$. Hvis det er heteroskedastisitet så blir dette målet på standardfeilen ikke riktig. I praksis har det ikke så mye å si og jeg skal vise at vi kan finne en mer generell formel for standardfeilen som ikke avhenger av den antagelsen. Videre kan vi også finne alternativ estimator som utnytter at noen observasjoner er mer informative om verdien til β for de feilleddet til observasjonen har mindre varians. Mer generelt så kan det være ønskelig å transformere variablene slik at de oppfyller $G - M$ antagelser og vi får mer effektiv estimering... tror ikke dette er så veldig relevant i praksis, men jeg tar det litt raskt.

7.6.1 Generalisert minste kvadrat

Vi kan skrive $\mathbb{V}[\mathbf{u}|X] = \sigma^2 \psi$. Vi har altså en parameter som er skalert med en matrise som kan avhenge av X . Jeg vil transformere variabelen \mathbf{u} slik at skaleringsfaktoren reduserer

til I . Merk først hvordan vi kan gå fram for å standardisere i én dimensjon,

$$\mathbb{V}[u|x] = a\sigma^2 \quad (7.59)$$

$$\implies \mathbb{V}\left[\frac{u}{\sqrt{a}}|x\right] = \frac{1}{a}\mathbb{V}[u|x] = \sigma^2 \quad (7.60)$$

Jeg bare skalerer variabelen med den inverse av kvadraturen av skaleringsfaktoren i uttrykket for variansen. Kan greit generalisere dette ved å finne A slik at $\psi^{-1} = A'A$.⁶ Kan da transformere modellen

$$A\mathbf{y} = A[X\beta + \mathbf{u}] \quad (7.61)$$

$$\mathbf{y}^* = X^*\beta + \mathbf{u}^* \quad (7.62)$$

og observere at

$$\mathbb{V}[Au|X] = A\mathbb{V}[u|X]A' \quad (7.63)$$

$$= A\sigma^2\psi A' \quad (7.64)$$

$$= \sigma^2 A(A'A)^{-1}A' \quad (7.65)$$

$$= \sigma^2 AA^{-1}(A')^{-1}A' \quad (7.66)$$

$$= \sigma^2 I \quad (7.67)$$

7.6.2 Vektet minste kvadrat

Hvis vi utelukker seriekorrelasjon er ψ en diagonalmatrise. I mange tilfeller er det rimelig at varians til feilledd⁷ avhenger av de observerte variablene. For eksempel er spredningen til mange variabler større for menn enn for kvinner. Det kan også være slik at størrelse avhenger av kontinuerlig variabel (mer variasjon i forbruk for rikinger). Vi kan generelt modellere dette som

$$\mathbb{V}[u|X] = \sigma^2\psi = \sigma^2 \text{diag}(h_n)^2 \quad (7.68)$$

der $h_n^2 := h(x_n)$. Tror det vil være en eller annen deterministisk funksjon av variablene.. Uansett, bare skalerer alle observasjoner med kvadratet av den inverse for å redusere ψ til I ,

$$\frac{y_n}{h_n} = \frac{x_n}{h_n}\beta + \frac{u_n}{h_n} \quad (7.69)$$

⁶Dette kan f.eks. gjøres med cholesky dekomponering. Eksisterer alltid fordi ψ er positiv semi definit (analog til positiv skalar), men er ikke unik.

⁷Siden feilledd kan tolkes som avvik fra sentraltendens i betinget fordeling så tilsvarer det varians i betinget fordeling.

I praksis så er det stor utfordring at h_n må estimeres. Hvis vi ikke påfører mer struktur så blir det like mange parametre som observasjoner. Litt usikker på hvordan jeg går frem i praksis. Et enkelt eksempel er grupperte observasjoner det jeg kun ser gjennomsnitt i gruppen, men antar at $u \sim IID(0, \sigma^2)$. Da vil $u_i := \frac{1}{N(i)} \sum u_i$ og $h_i^2 = \frac{\sigma}{N}$. Legger med vekt på grupper med flere observasjoner.

Feasible generalisert minste kvadrat

7.7 Robust estimering (sandwich)

Vi brukte i utgangspunktet antagelse om at feilleddene var *iid* for å beregne covariansmatrisen til estimatoren av koeffisientvektoren. Vi kan utvide til $E[\mathbf{u}\mathbf{u}'] = \mathbf{\Omega}$ på ukjent form. Hvis matrisen er diagonal er det kun heteroskedastisitet. Hvis det også er kovarians mellom feilleddene er det såkalt autokorrelasjon. Vi har et uttrykk for variansen til koeffisientestimatoren,

$$V(\hat{\beta} - \beta) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{\Omega}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \quad (7.70)$$

som ligner på en sandwich der $\mathbf{X}'\mathbf{\Omega}\mathbf{X}$ er fyllet i midten. Det er en utfordring at vi ikke kan konsistent estimere $\mathbf{\Omega}$ siden den har mer enn N parametre og vokser med utvalgsstørrelsen. Selve fyllet derimot er en $K \times K$ -matrise som kan estimeres konsistent likevel. Finnes ulike kandidater for $\hat{\mathbf{\Omega}}$.

7.8 Annet

Regresjonsanatomi

$$\text{cov}(y_i, \tilde{x}_{ki}) = \text{cov}(\mathbf{x}'\beta + u, \tilde{x}_{ki}) = \beta_k \text{var}(\tilde{x}_{ki}) \quad (7.71)$$

$$\implies \beta_k = \frac{\text{cov}(y_i, \tilde{x}_{ki})}{\text{var}(\tilde{x}_{ki})} \quad (7.72)$$

hmmm. frischhh

Anova (flytt til betinget fordeling/L2)

$$V(y) = V(E[y|x] + u) = V(E[y|x]) + V(u) \quad (7.73)$$

der $V(u) = E[u^2] = E[E[u^2|x]] = E[V[u^2|x]]$.

Generalisert lineær modell, flytte hvor?

Utvide... fordeling til u (eg betinget fordeling av y); ikke bare normal, andre medlem av eksponentialfordeling. Dessuten utvide til linkfunksjon... mer utvidet måte å modellere $E[y|x]$. Vet ikke hvor relevant og spennende dette egentlig er.. vet faen ingen ting.

Saturated model

Benchmark, start punkt. Så forenkle. Utvide til eksponering av behandling som kan ta mange verdier. Se på heterogen behandlingseffekt..

Kapittel 8

Tidsserier

I tidsserier gjør vi gjentagende observasjoner av samme enhet. Det kan for eksempel være utvikling i brutto nasjonalprodukt (BNP) til et land, størrelsen på populasjonen, pris på verdipapir over tid eller andre størrelser som endrer seg. Det vil være korrelasjon mellom observasjoner som er nær hverandre i tid slik at vi ikke kan bruke vanlige statistiske metoder der vi antar at de ulike observasjonene er uavhengige realiseringer fra samme fordeling.¹ I likhet med statistisk analyse i krysseksjon tar vi utgangspunkt i tabell med tall som vi kan betrakte som realiserte verdier fra en datagenereringsprosess. Det er vanskelig å lese tallene direkte, så vi må lage alternative representasjoner som gjør at vi kan lære om DGP og bruke dette til å svare på spørsmål.

Vi betrakter DGP som en følge av stokastiske variabler $(x_t)_{t \in T}$ der T er indeks til tidspunktene. Dette er en stokastisk prosess og følgen av realiserte verdier utgjør én realisering av prosessen. Vi kan visualisere den observerte realiseringen i utvalget ved å plote verdien opp mot tid, der vi gjerne bruker rette linjer mellom punktene som en tilnærming på den konseptuelt kontinuerlige underliggende tidsserien.² Videre bruker vi modell til å beskrive DGP og bruker de realiserte verdiene til å estimere parametre i denne modellen. Det finnes ulike modeller. De bestemmer hvor glatt tidsserien er.. hm.

1. Hvit støy: Ingen korrelasjon mellom ulike ledd av følgen og konstant spredning.
2. Moving average: Kan tenke at hver observasjon egentlig er løpende (vektet) gjennomsnitt av underliggende (uobserverte) verdier, eks: $w_t = \frac{1}{3}(v_{t-1}, v_t, v_{t+1})$.
3. Random walk med drift: $x_t = \mu + x_{t-1} + w_t$ der (w_t) er hvit støy.
4. Signal og støy: $x_t = f(t) + w_t$. Et eller annet gjentagende mønster pluss tilfeldig avvik.

¹Litt usikker på i hvilken grad det er overlapp mellom metodologi i tidsserier og i panel/hierarkisk data der det er korrelasjon mellom observasjon innad i grupper uten at det nødvendigvis er tidsdimensjon.

²Det at tidsserien er diskret og størrelsen på tidssavstand mellom observasjoner er gjerne bare en noe arbitrær egenskap ved hvordan data er samlet inn. Merk at de estimerte og observerte egenskapene til tidsserien kan være sensitive for størrelsen på dette gapet.

8.1 Deskriptivt

Vi begynner alltid analysen ved å plote tidsserien. Fra denne grafiske representasjonen kan vi få inntrykk av egenskaper ved tidsserien.

1. Langsiktig trend (endring i gjennomsnittsverdi).
2. Systematisk variasjon om den langsiktige trenden (over eller under i lengre tid). Kan ofte være sykklisk etter fast mønster på grunn av årstider eller andre gjentakende fenomen.
3. Ytterligere tilfeldig variasjon.
4. Kan se om det er noen outliars (verdier som avviker kraftig fra resten av tidsserien). Kan skyldes målefeil eller spesielle hendelser. Kan være aktuelt å fjerne disse for at de ikke skal påvirke videre analyse (forecasting), men må være litt forsiktig med dette.
5. Endringer i egenskapene over. Kan det for eksempel være at trenden endret seg på gitt tidspunkt? Kan indikere at det er aktuelt med separat analyse av ulike segment av tidsserie eller hm..

Dette er et greit første steg, men hvordan skal vi gå videre? Jeg vil tallfeste noen egenskaper i stedet for å bare *eyeballe* de fra figuren. Dessuten vil jeg dekomponere tidsserien slik at jeg kan rendyrke de ulike aspektene i ulike figurer i stedet for å ha alt i samme figur. Jeg kan for eksempel både se på den langsiktige trenden, de systematiske eller sykliske avvikene og deretter se på gjenstående tilfeldige avvikene (residualene) etter å ha forsøkt å modellere den systematiske delen. Med andre ord finner vi $y_t = f(t) + u_t$. I motsetning til i krysseksjon vil vi da gjerne observere at det er struktur i residualene ved at de er korrelert over tid og det virker som at mye i tidsserier handler om å håndtere dette. Dessuten overlapper dekomponering med modellering siden vi det ikke er en entydig måte å dekomponere tidsserien slik at vi må gjøre litt ulike valg. Har generelt tre fremgangsmåter for å fjerne trenden i tidsserie:

1. Parametrisk modell der $f(\cdot)$ er f.eks. et polynom
2. Filtring der vi beregner $f(\cdot)$ er beregnet ut fra vektet gjennomsnitt av nærliggende observasjoner. Må velge grad av smoothing.
3. Kan bruke differanser mellom observasjoner til å finne avvik fra trend uten å finne $f(\cdot)$ eksplisitt.

Deretter gjenstår det å modellere u_t ...

8.1.1 Stasjonaritet

Det er et problem for vår inferens at vi bare har én observasjon for hvert tidspunkt t . For å komme noen vei må vi anta at egenskapene til tidsserien er stabile over tid slik at vi får flere observasjoner til å lære egenskapene fra. Stasjonaritet er en påstand om sammenhengen mellom en delmende av observasjonene og en forskjøvet delmengde, altså om

$$(x_t, x_{t+1}, \dots, x_{t+k}) \text{ og } (x_{t+h}, x_{t+h+1}, \dots, x_{t+h+k}). \quad (8.1)$$

I sterk form er hele simultanfordelingen på to delmengdene den samme. Dette er sterkere enn vi trenger, lite realistisk og vanskelig å vurdere om det holder i praksis. Vi trenger bare svak stasjonaritet som kun legger begrensninger på første to moment. Med svak stasjonaritet er

1. $\mu_t = \mu$ for alle t . Forventningsverdien endrer seg ikke slik at vi i prinsippet får T observasjoner til å lære μ .³
2. $\gamma(s, t) = \text{gamma}(s + h, t + h)$. Størrelsen til autokovarians avhenger bare i avstand i tid mellom to observasjoner, $h = |t - s|$, ikke absolutt posisjon i tid. Det medfører at vi kan betrakte det som bare en funksjon av h og at vi får flere realiseringer.

8.1.2 Empirisk

Hvis tidsserien er stasjonær kan vi da konsistent estimere gjennomsnitt med \bar{x} og autokorrelasjon med lag j ved

$$\hat{\rho}(j) = \frac{\widehat{cov}(y_t, y_{t-j})}{\widehat{var}(y_t)} \quad (8.2)$$

der

$$\widehat{cov}(y_t, y_{t-j}) = \frac{1}{T-1} \sum_{t=1}^{T-j} (y_t - \bar{y})(y_{t-j} - \bar{y}) \quad (8.3)$$

og

$$\widehat{var}(y_t) = \frac{1}{T-1} \sum_{t=1}^{T-j} (y_t - \bar{y})^2. \quad (8.4)$$

Disse empiriske (eller utvalgs-) størrelsene kan vi alltid beregne for en tidsserie og det kan være et godt sted å begynne før vi gjør noe modellering. Vi kan plote autokorrelasjon

³Merk at hvis de ikke er uavhengige så lærer vi ikke like mye fra nye observasjon slik at større standardfeil til estimatet, ref. clustering i økonometri.

for laggene $h = 1, \dots, H$ for å undersøke strukturen. Hvis det kun er korrelasjon med perioden foran så kan det for eksempel være rimelig å bruke MA(1).

8.2 Modelling

Sammenhengen mellom fordelingen til de ulike leddene i den stokastiske prosessen er fullt ut beskrevet av simultanfordelingen som er gitt ved $F : (x_1, \dots, x_T) \mapsto \mathbb{R}$. Dette er et veldig komplisert objekt og umulig å få noe godt estimat fra én observasjon. Vi kan i stedet ta utgangspunkt i de marginale fordelingene til hvert av leddene F_t og betrakte forventningsverdien $\mu_t := E[x_t] = \int f_t(x_t)x_t dx_t$. Videre kan vi se på kovarians mellom variabler på ulike tidspunkt, $\gamma(s, t) := \text{cov}(x_s, x_t) = E((x_s - \mu_s)(x_t - \mu_t))$. Eller korrelasjonen.

I praksis er de fleste tidsserier ikke-stasjonære fordi gjennomsnittet endrer seg over tid ut fra en eller annen trend μ_t eller fordi spredningen endres. Vi kan dekomponere tidsserien slik at den består av trend og avvik fra trend, der avviket kan være en stasjonær tidsserie med $\mu = 0$ og autokovarians som vi konsistent kan estimere. Vi kan spesifisere en parametrisk form for trenden, for eksempel

$$\mu_t = \beta_0 + \beta_1 t \tag{8.5}$$

og estimere funksjonen med minste kvadrats metode slik at vi får estimerte residualer $\hat{u}_t := y - \hat{\mu}_t$ som utgjør en realisering av tidsserie med egenskaper vi kan estimere. Det finnes selvsagt alternative måter å dekomponere i trend og avvik; vi kan bruke ikke-parametriske (*smoothing*) estimatorer for å håndtere vilkårlig ikke linearitet i trend som funksjon av tid, eller vi kan også bruke tidligere realiserter verdier av egen og andre tidsserier som argument i funksjonen.⁴ Vi har sett på korrelasjon mellom variabler med ulike lag h . Dette er en deskriptiv egenskap ved verdiene i utvalget og vi kan betrakte det som et estimat på egenskap til simultanfordelingen som genererte data. Jeg ønsker å gå videre fra denne egenskapen til å modellere prosessen som genererer data. Vi kan bruke parametrisk spesifisering og estimere parametre slik at avhengighetsstrukturen samsvarer til dels med den observerte avhengigheten i utvalget. Det er fordel at fordeling til gitt variabel y_t avhenger av tidligere realiserter variabler slik at vi kan bruke modellen til å forutsi fremtidig variabler. For å forenkle notasjon med laggede variabler introduserer vi en såkalt *backshift operator* B der $B^k x_t = x_{t-k}$.

⁴Litt usikker på first differencing... får ikke trend, men får dekomponert tror jeg.

8.2.1 AR(p)

En mulig fremgangsmåte for å skape avhengighet er å ta regresjon på egne laggede verdier.⁵

$$y_t = \theta_1 y_{t-1} + \theta_2 y_{t-2} + \cdots + \theta_p y_{t-p} + u_t \quad (8.6)$$

8.2.2 Autoregressiv, AR(k)

I autoregressiv modell blir noe av utfall i forrige periode dratt med over til neste. Hvis det er spesielt høy verdi i én periode (høy ϵ_t) så blir det propagert videre i kommende perioder. Vi modellerer det med lagged utfall,

$$Y_t = \alpha + \rho Y_{t-1} + \epsilon_t \quad (8.7)$$

Hvor stor del av verdi som blir propagert videre avhenger av parameter som vi estimerer fra observert data. I praksis er det ofte enklere å jobbe med avvik fra gjennomsnitt, $y_t := Y_t - E[Y_t]$, og kan vises at

$$y_t = \rho y_{t-1} + \epsilon_t \quad (8.8)$$

Vil beskrive egenskap ved simultanfordeling til prosessen. Den er karakterisert ved såkalt autokovarians; korrelasjon mellom utfall på ulike tidspunkt.. Kan finne $var(y_t)$, $cov(y_t, y_{t-1})$ og $cov(y_t, y_{t-k})$..

8.2.3 MA(q)

Dette er en alternativ måte å beskrive sammenhengen mellom realisering på ulike tidspunkt. I stedet for at hele utfallet blir propagert videre er det nå bare selve sjokket ϵ_{t-1} som blir med å bestemme verdi i neste periode. Dette medfører at sjokket ikke påvirker verdi inn i evigheten,

$$Y_t = \mu + \alpha \epsilon_{t-1} + \epsilon_t \quad (8.9)$$

8.2.4 ARMA(p,q)

8.3 Tidsserier

Vi skal nå se på data der vi har gjentatt observasjon av samme enhet. Så langt har vi modellert sammenheng mellom variabler som blir realisert samtidig, $y = f(x) + u$. Hvis vi ønsker å predikere fremtidig verdi y_{t+k} når vi er i t så hjelper det oss ikke så mye å ha

⁵Merk at auto betyr self, eg autofellatio

god f siden vi uansett må predikere x_{t+k} for å bruke den. En alternativ fremgangsmåte er å modellere stokastisk prosess $(y_t)_{t \in \mathbb{N}}$. I de fleste tilfeller vil fordeling til realisering t avhenge av realisering i tidligere perioder. Jeg skal begynne med å se på enkle måter å modellere denne avhengigheten. Deretter skal jeg inkludere andre forklaringsvariabler..

8.3.1 Lineær trend

Kanskje den enkleste måten å beskrive trenden er med en enkel lineær sammenheng der vi lar tidsperiode t være uavhengige variabel,

$$y_t = \alpha_0 + \alpha_1 t + \epsilon_t \quad (8.10)$$

der parametrene er definert slik at de konstruerer feilledd med egenskap $E[t\epsilon_t] = 0$ og $E[\epsilon_t] = 0$. Kan jo også påstå at $E[\epsilon_t|t] = 0$ men det er en ganske sterk påstand siden gjennomsnittlig utfall i hver periode sjeldent endrer seg helt linært. Dersom vi har få perioder så kunne vi gitt en indikator for hvert tidspunkt som gir fullstendig fleksibel beskrivelse av $E[y|t] = \alpha_t d(t)$,

$$y_t = \sum_t \alpha_t d_t + \epsilon_t \quad (8.11)$$

der $E[\epsilon_t|t] = 0$, men har jo bare én observasjon per tidspunkt så blir like mange parametre som observasjoner. Denne fremgangsmåten blir mer hensiktsmessig med paneldata der jeg har mange observasjoner i hver t . Vil da modellere trend for å *de-trende*, fjerne trend fra sammenheng mellom behandling og utfall jeg interessert i. Kan enkelt utvide til eksponentiell trend ved å ta logaritmisk transformasjon av sammenheng slik at den blir lineær,

$$y_t = e^{\beta_0 + \beta_1 t + \epsilon_t} \quad (8.12)$$

$$\log y_t = \beta_0 + \beta_1 t + \epsilon_t \quad (8.13)$$

8.4 Annet

Vi kan bruke stokastiske prosesser med avhengighet til å modellere feilleddet i en regresjonsmodell dersom vi av ulike grunner ikke kan anta at de er uavhengige. Den observerte autokorrelasjonen mellom residualer kan skyldes feilspesifikasjon av modellen, for eksempel at den ikke fanger opp all ikke-linearitet og at det er korrelasjon mellom uavhengige variabler som blir samlet på ulik tidspunkt. Vi har i utgangspunktet to fremgangsmåter for å håndtere dette:

1. Bruke generalisert minste kvadrat for å utnytte at ikke alle observasjonene er like

informative om sammenhengen vi er interessert i. Transformerer modell slik at den oppfyller *iid* antagelse slik at vi både får mer effektive estimat og riktige standardfeil. Disse gode egenskapene forutsetter at vi treffer riktig på strukturen. Tror det er noe kobling til (bayesianske) hierarkiske modeller som modellerer avhengighetsstruktur mer eksplisitt.

2. Vi kan observere at MKM-estimatene fortsatt er konsistente så lenge den stokastiske prosessen til feilleddene er stasjonær. Hver observasjon vil være mindre informativ enn ved uavhengige feilledd slik at standardfeilene til estimatene blir større, men vi kan ta hensyn til dette ved å bruke autokorrelasjons-robuste feilledd som gir asymptotisk konsistente estimat av standardfeilene.⁶

Alternativt kan vi bruke stokastiske prosesser til å modellere forventningsverdi av en univariat tidsserie. En enkel fremgangsmåte ville vært å modellere $E[y|t] = f(t)$, for eksempel ved en lineær trend. Deretter kunne vi predikert fremtidige verdier \hat{y}_{t+j} med $\hat{f}(t+j)$. Dette er sikkert en grei tilnærming hvis vi ser på tidspunkt langt inn i fremtiden, men hva hvis $t+j$ er i nær fremtid? Da vil det gjerne være sann at hvis y_t ligger over trenden så er det kanskje også mer sannsynlig at y_{t+j} også ligger over? Dette gjelder spesielt hvis det er store rigiditer i tidsserien. Eller kanskje vi også vil ta hensyn til momentum i tidsserien? Hvis den holder på å enten vokse eller synke så er det kanskje større sannsynlighet for at den vil fortsette med dette?

8.4.1 Mer annet

Tidsserier består av gjentagende observasjoner av en størrelse. Konseptuelt kan tidsserien være kontinuerlig, men vi observerer verdier på diskret tidspunkt. Den kan være verdi på gitte tidspunkt eller aggregering av verdi i intervallet mellom tidspunkt.

⁶Litt usikker på om det er noe kobling til clustering her.

Kapittel 9

Statistisk læring

I maskinglæring bruker vi data til å lære en maskin å utføre oppgaver. Dette er i kontrast til tradisjonelle algoritmer der vi eksplisitt spesifiserer regler for hva maskinen skal gjøre. Maskinen får noe *input*, transformerer det til noe *output* og får tilbakemelding på i hvilken grad den klarte å utføre oppgaven. Algoritmen kan deretter tilpasse transformasjonen ut fra tilbakemeldingene den får. På denne måten kan maskiner blant annet lære å kjøre biler og spille sjakk. Jeg vil i hovedsak bruke det til å lære maskinen å predikere en såkalt *utfallsvariabel* til en observasjon med utgangspunkt i informasjon om andre egenskaper til observasjonen.

Statistisk læring kan betegnes som tilnærmingen til denne tematikken fra statistikk i stedet for informatikk. Maskinlæringsalgoritmer lærer fra data og dette kan vi formalisere innenfor et statistisk rammeverk. Mer spesifikt kan vi betrakte målet som å lære egenskaper til en fordeling \mathbb{P} som har generert data vi observerer. Dette perspektivet gjør det blant annet mulig å håndtere usikkerheten til en prediksjon på en systematisk måte. Rammeverket gir oss også verktøy for å lage transparente modeller som representerer egenskaper til populasjonen på en oversiktlig måte, selv om det nødvendigvis innebærer forenklinger og antagelser. I kontrast kan rene maskinlæringsalgoritmer være en såkalte *black-box* funksjoner $h : \mathbf{x} \mapsto \hat{y}$ som kun gir oss predikert utfall \hat{y} for hver input \mathbf{x} . I praksis er det stort overlapp mellom tilnærmingene så forskjellene bør ikke overdrives, men kapitlet har fått denne tittelen fordi jeg interessert i inferens og tolkning av modeller; ikke såkalt nevralt nettverk og andre maskinglæringsalgoritmer som er velegnet for ikke-tabulære data.

9.1 Bakgrunn og oversikt

Maskinlæringsalgoritmer bruker erfaring E til å bli bedre til å utføre oppgave T som målt ved kriterie P . Supervises så er erfaring data med både input \mathbf{x} og utfall y . Kan både oppnå bedre resultat og program som er enklere å utvikle og vedlikeholde, og som ikke trenger like mye *domain knowledge*.

Algoritme består av tre deler

1. Representasjon av egenskaper den lærer. Hypoteserom: kandidater av funksjoner h den søker over
2. Evaluerer: mapper hver kandidat til et mål på fit. F.eks: $g : \beta \mapsto RSS(\beta)$ i lineær regresjon
3. Optimering: må ha måte å effektivt søke over hypoteserommet for å finne gode kandidatfunksjoner

Kategorisere algoritmer: Supervised (reg og klassifikasjon) vs unsupervised (clustering, dimensjonsreduksjon, anamolie). Batch vs incremental learning. Instance vs model based. Finnes mange ulike... hvilken som er best avhegner av data. (no free lunch..) Mer struktur er bra hvis riktig struktur. I praksis: ensemble hvis maksimere prediksjon.

9.1.1 Generativ modell

Målet er å modellere assosiasjon¹ mellom input \mathbf{x} og output y .

I praksis kan vi ofte anta at \mathbf{x} er kjent eller at vi driter litt i fordeling, så trenger bare betinget fordeling $y|\mathbf{x} = \mathbf{s}$ til å svare på spørsmål.

Kan beskrive all info $f_{y|\mathbf{x}=\mathbf{s}}(y)$ der for alle $A \subset \mathbb{R}$ så er $\mathbb{P}(y \in A|\mathbf{x} = \mathbf{s}) = \int_A f_{y|\mathbf{x}=\mathbf{s}}(y)$

Vanskelig å jobbe med selv om vi kjente fordeling og enda vanskeligere å lære fra begrenset data! Vil ha sammendragsmål som egenskap. $E[y|\mathbf{x} = \mathbf{s}] = g(\mathbf{s}; \theta)$.

Tenker at jeg har lyst til å knytte statistisk læring i større grad opp mot Bayes siden det er vesentlig å kvantifisere usikkerhet.

9.1.2 Utfordringer

Generalisere fra utvalg

Utfordringer: For lite data. Signal og støy... støy jevner seg ut (per konstruksjon/definisjon), avvik fra sentraltendens. Avhenger av hvor sterkt signal er i forhold til støy. Problem at komplekse/fleksible algoritmer er veldig flinke til å finne mønster. Finner selv om det bare skyldes tilfeldigheter ved utvalget (alle med navn som begynner på 's' og slutter på 'e' er veldig smarte.. i utvalget. Generaliserer ikke). Problem som kan løses", men legger begrensninger på løsning.

Virkeligheten er komplisert og modellene er ofte veldig enkle. Det kan derfor være fristende å bruke mer fleksibel struktur som lar dataene snakke". Det er flere avveininger knyttet til dette. For det første kan det bli vanskeligere å beskrive og tolke den estimerte modellen. Dette er mindre problematisk dersom modellen skal brukes til prediksjon, men

¹Bruker det bekrepet siden det ikke er en sammenheng eller tilknytning, men den er ikke eksakt

det kan uansett være interessant å bruke modellen til å lære om hvordan input er relatert til output. Selv om vi kun er interessert i best mulig prediksjon er ikke alltid mer fleksibilitet bedre. Problemet er at vi lærer *for mye* om utvalget, mens vi egentlig er interessert i egenskaper ved prosessen som genererte det.

Det signalet vi er interessert i å lære er ofte $E[Y|X] = f(X)$. For hver y_n i utvalget er $y_n = f(x_n) + u_n$. Vi vil lære $f(\cdot)$, men med men for mye fleksibilitet vil det i for stor grad også fange opp støyen u_n i det gitte utvalget. Dette er et eksempel på bias-variance-tradeoff.

Et eksempel på at det ikke alltid er tilstrekkelig med høy test accuracy på data vi har tilgjengelig er klassifikasjon mellom ulv og hund (eg. husky). Hvis bilde av ulv er i område med snø, så vil snø være input som gjør at modell kan oppnå gode prediksjon. Men da bygger vi en modell som er god til å oppdage snø og det er kanskje ikke så nyttig for oss.

Ved å se på sammenheng mellom input og output kan vi avdekke om det er såkalt irrelevante features (partikulære egenskaper ved de gitte dataene) som har stor forklaringskraft. Dette vil indikere at modellen ikke kommer til å generalisere så bra. Ved å analysere de uriktige prediksjonene kan vi også lære mer om hvilke nye variabler / features som kan være nødvendig for å oppnå bedre prediksjon.

Dimensjonalitetens forbannelse

Det kan være en utfordring å jobbe med høydimensjonal data. Problemet er størrelsen på rommet vokser veldig raskt når vi øker antall dimensjoner slik den gjennomsnittlige avstanden mellom observasjonene også øker. Dette medfører at vi må ekstrapolere kurver til områder av rommet der vi har lite informasjon slik at det blir fort gjort å overfitte. Det medfører også at nabolaget ikke er så veldig lokalt og metoder som bygger på avstand mellom observasjon ikke fungerer så bra.²

Kvantifisere usikkerhet

Hvis vi observerer flere variabler kan vi redusere usikkerheten, men det kan også være variabler som er fundamentalt uobserverbare, spesielt når vi analyserer menneskelig atferd. En annen kilde til usikkerhet er målefeil i variablene, slik at observerasjoner med samme observerte x kan ha ulik verdi av de reelle størrelsene som påvirker utfallet. Eksistens av usikkerhet og kvantifisering av denne.

Usikkerhet til sentraltendens + gjennomsnittlig avvik fra sentraltendens.

²Tror litt av grunnen til at neurale nettverk fungerer så bra på høydimensjonal data er at det klarer å lære meningsfull representasjon i lavere dimensjon...

Ikke-representative data

Annen utfordring er verre: Ikke representative data. Kan skyldes tilfeldigheter ved utvalg (denne risikoen kan i prinsippet kvantifiseres... men kan gi dårlig performance... mer data er bra). Skjevhet som ikke løses med mer data (seleksjonsproblem, respons bias). Generaliserer ikke til populasjon. Dårlig data (målefeil). Viktig med data cleaning og god feature engineering... hvor mye modell klarer å lære avhenge av representasjon av data... teori + kryssvalidering.

9.2 Empirisk risikominimering

Vi kan formalisere læringsproblemet som et risikominimeringsproblem. Vi observerer $(\mathbf{z}_1, \dots, \mathbf{z}_N)$ der $\mathbf{z}_n = (y_n, \mathbf{x}_n)$ er realiseringer fra $\mathcal{L}(\mathbf{z}) = P$. Målet vårt er å finne en funksjon som tar verdi $f(\mathbf{x})$ og predikerer *output* y gitt at vi observerer *input* \mathbf{x} . I praksis vil det være avvik mellom predikert og sann verdi. Vi kan definere en *tapsfunksjon* L som avhenger av størrelsen på dette avviket.³ Noen vanlige tapsfunksjoner er

- Absolutt tap: $L(y, f(\mathbf{x})) = |y - f(\mathbf{x})|$
- Kvadratisk tap: $L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$
- Diskret tap: $L(y, f(\mathbf{x})) = I\{y \neq f(\mathbf{x})\}$

Merk at tapsfunksjonen er en tilfeldig variabel som tar verdi for hver realisering av \mathbf{z} og at fordelingen dermed avhenger av den ukjente simultanfordelingen P . Målet vårt er å minimere forventet tap som kan betegnes som *prediksjonsrisikoen* $R(f) = \mathbb{E}L(y, f(\mathbf{x}))$. Utfordringen er at vi ikke kjenner P slik at vi ikke kan evaluere prediksjonsrisikoen direkte. En mulig løsning er å bruke utvalgsanalogprinsippet og betrakte den *empiriske risikoen*

$$R_{emp}(f) \equiv \mathbb{E}_{\hat{P}_N} L(y, f(\mathbf{x})) = \frac{1}{N} \sum L(y_n, f(\mathbf{x}_n)) \quad (9.1)$$

En naiv tilnærming vil nå være å finne f som minimerer $R_{emp}(f)$, men dette vil ofte være en dårlig løsning. Dersom alle inputvektorene \mathbf{x}_n tar ulike verdier vil det alltid være mulig å finne funksjoner f slik at $f(\mathbf{x}_n) = y_n, n = 1, \dots, N$ og $R_{emp}(f) = 0$. Hvis vi minimerer empirisk risiko vil vi få en funksjon som har lært for mye om det gitte utvalget vårt og generaliserer dårlig til nye observasjoner. Dette problemet kalles *overfitting*.

For å unngå overfitting må vi begrense hvor mye algoritmen lærer fra det gitte realiserte utvalget. Vi har i hovedsak to måter å gjøre dette på. For det første kan vi modifisere tapsfunksjonen slik at det påføres ekstra kostnad dersom funksjonen tilpasser seg data i

³Dette rammeverket gir oss større fleksibilitet enn om vi kun ser på størrelsen av avviket $u := y - f(\mathbf{x})$. Det kan for eksempel være slik at kostnad med prediksjonsfeil er asymmetrisk slik at større kostnad ved å enten over- eller underpredikere. Litt usikker på hvor relevant dette er og om jeg kan gjøre det operativt...

utvalget. Dette kalles regularisering og vi skal se på det senere. Den andre muligheten er å avgrense oss apriori til å kun betrakte kandidatfunksjoner i et hypoteserom \mathcal{H} . Løsningen på det empiriske risikominimeringsproblemet kan da uttrykkes som

$$\hat{f} = \arg \min_{f \in \mathcal{H}} R_{emp}(f). \quad (9.2)$$

Hvis vi antar at $x \in \mathbb{R}$, så kan et mulig hypoteserom være mengden av alle polynomial av grad p ,

$$\mathcal{H}_p = \{f : f(x) = \beta_0 + \beta_1 x + \dots + \beta_p x^p\}. \quad (9.3)$$

Det empiriske risikominimeringsproblemet med kvadratisk tapsfunksjon er da

$$\arg \min_{\beta} \frac{1}{N} \sum ((y_n - (\beta_0 + \beta_1 x + \dots + \beta_p x^p))^2 \quad (9.4)$$

der vi kan finne løsningen analytisk med minste kvadrats metode. Med ulik grad p av polynomfunksjonen får vi ulike hypotesefunksjoner. Hvordan velger vi hvilken \mathcal{H}_p som er best? Vi kan ikke bruke empirisk risiko som mål fordi

$$\mathcal{H}_1 \subset \mathcal{H}_2 \implies R_{emp}(\hat{f}_1) \geq R_{emp}(\hat{f}_2) \quad (9.5)$$

Mer fleksibilitet vil alltid medføre at funksjonen kan lære mer fra data og få bedre *in-sample fit* og dermed lavere empirisk risiko. Målet vårt er derimot å generalisere til nye data og oppnå best mulig *out-of-sample fit*, altså prediksjonsrisiko. For å velge optimal struktur må vi bruke kryssvalidering.

9.2.1 Kryssvalidering

For å estimere prediksjonsrisiko til en gitt hypotesefunksjon \hat{h} trenger vi usette data som ikke ble brukt i opplæringen av funksjonen. Hvis vi setter av J observasjonen kan vi bruke

$$\widehat{R(\hat{h})} = \frac{1}{J} \sum L(y_j, \hat{h}(\mathbf{x}_j)). \quad (9.6)$$

som ikke systematisk favoriserer med fleksible funksjoner. Testing på usette data er det vesentlige, men vi forbedre testingen og utnytte data mer effektivt gjennom såkalt K-fold-kryssvalidering. I stedet for å kun oppdele i trenings- og testdata, kan vi partisjonere datasettet \mathcal{D} i K like store deler som vi angir med \mathcal{D}_k ($k = 1, \dots, K$). I hvert steg holder vi av én del som testdata og trener algoritmen på de restenrende delene. Målet på prediksjonsrisiko blir da gjennomsnittet av estimert risiko på hvert av testdelene. Med andre ord blir algoritmen da:

1. for k in $1, \dots, K$:
2. fit \hat{h} på \mathcal{D}_{-k}
3. finn $\hat{R}_k = \frac{1}{|\mathcal{D}_k|} \sum_{n:n \in \mathcal{D}_k} L(y_n, \hat{h}(\mathbf{x}_n))$
4. end for, finn $\widehat{R(\hat{h})} = \frac{1}{K} \sum \hat{R}_k$

Til slutt velger vi $\hat{h}_{opt} = \arg \min \widehat{R(\hat{h})}$. Merk at for å få forventningsrett estimat på hvor godt algoritmen fungerer på usette data må vi holde av enda et testsett som ikke har vært brukt i kryssvalideringen.

9.2.2 Dekomponering av risiko med kvadratisk tap

Det beste vi kan gjøre er å finne en funksjon h slik at

$$y = h(x) + \epsilon \quad (9.7)$$

der ϵ er uavhengig av x . Det medfører at funksjonen h fanger opp all informasjon om verdi av y slik at det resterende feilledet er uavhengig av x . Vi bruker da $\hat{y} = h(x)$ som predikert verdi av y . Vi bruker forventet kvadrert avvik som mål på prediksjonsfeil, og siden dette er det beste vi kan oppnå er den såkalte *irreducible error*

$$E[(y - h(x))^2] = \text{Var}[\epsilon] \quad (9.8)$$

Jamført med diskusjon om projeksjon i L_2 er $h(x)$ projeksjonen av y på underrommet som består av alle tilfeldige variabler som kan skrives som en deterministisk funksjon av x . Dette tilsvarer den betingede forventningsfunksjonen. I praksis så må vi estimere h fra realiserte verdier i utvalg. Vi finner da en annen \hat{h} i underrommet. Den kvadrerte avstanden fra \hat{h} til y er

$$E[(h(x) + \epsilon - \hat{h}(x))^2] = E[(h(x) - \hat{h}(x))^2] + \text{Var}[\epsilon] \quad (9.9)$$

der første ledd er prediksjonsfeilen vi kan ha håp om å redusere gitt x . Det er flere måter å vise denne dekomponeringen, men det følger av ortogonal projeksjon og pythagoras. Mye av statistisk læring handler om å finne best mulig \hat{h} . I praksis avgrenser vi oss ofte til å se på et underrom av mengden av tilfeldige variabler som kan skrives som funksjon av x ; for eksempel alle lineære funksjoner. Biasen vil være avstand mellom $h(x)$ og $h^*(x)$ som er beste variabel i den delmengden. Det er i tillegg varians i estimeringen.

Utvalgsanalogen til MSE er

$$E_{P_N}[(y - \hat{h}(x))^2] \quad (9.10)$$

Den kan virke rimelig å minimere dette for å finne \hat{h} . Problemet er at dette er en forventningsskjev estimator av MSE og alltid vil foretrekke mer fleksible funksjoner som kan lære mønster i utvalget. Men målet vårt er ikke å memorisere utvalget! Målet er å predikere fremtidige data. Vi bruker derfor kryssvalidering til å estimere MSE: se hvor god jobb hypotesefunksjonen gjør på usette data. Vi vil da se at forholdet mellom MSE og fleksibilitet har en U-form. Bias reduseres og varians øker.

9.3 Lineær regresjon

I regresjonsmodeller kan utfallsvariabelen ta verdier på et intervall av tallmengden. Vi vil finne en funksjon som predikerer utfallet med utgangspunkt i en input. Denne funksjonen kan være vilkårlig komplisert, men i praksis vil vi bruke såkalt *lineær regresjon*. Det medfører at vi avgrenser oss til å betrakte et hypoteserommet $\mathcal{H}_l := \{h : h(\mathbf{x}) = \mathbf{x}'\mathbf{b} \text{ for noen } \mathbf{x} \in \mathbb{R}^K\}$. Denne lineære modellen har flere attraktive egenskaper. For det første er funksjonen representert med et endelig antall helningsparametre β som betegner de marginale effektene, $\frac{\partial}{\partial x_k} f(\mathbf{x}) = \beta_k$. Hvor meningsfull parametrene er avhenger riktig nok av i hvilken grad den lineære funksjonen tilnærmer *CEF*, men den attraktive egenskapen er at rammeverket er fleksibelt nok til å fange opp ikke-lineariteter gjennom såkalte *basistransformasjoner*. Selv om vi transformerer variablene kan det fortsatt være mulig å tolke koeffisientverdiene i henhold til de opprinnelige variablene. Den tredje gode egenskapen er at med endelig mengde data vil den lineære strukturen begrense variansen slik at prediksjonene ikke nødvendigvis blir bedre med mer fleksible modeller.

Målet vårt er å nærme oss y ved å finne en funksjon $h(\mathbf{x})$ som minimerer $\|u\|_{L_2}$. Minste kvadrats metode er konsistent estimator av den lineære regresjonsfunksjonen $h^* := \arg \min_{h \in \mathcal{H}_l} \|u(h)\|_{L_2}$, der $\mathcal{H}_l := \{h(\cdot) : h(\mathbf{x} = \mathbf{x}'\beta)\}$. Vet å forlenge \mathbf{x} ved å legge til nye variabler kan vi utvide hypoteserommet \mathcal{H}_l slik at vi i teorien kan komme nærmere y . Det kan dermed være fristende å tenke at flere input-variabler alltid er et gode og kaste hele kjøkkenskapet inn i algoritmen. Problemet er at variansen øker og at det generalisere dårlig til nye data.

Vi skal se på valg av *feature space* som er viktig i lineære modeller siden det per konstruksjon er begrenset hvor mye algoritmene kan lære om mønster i data på egenhånd.

9.3.1 Feature space

Vi kan gjøre en transformasjon

$$\Phi : \mathbf{x} \mapsto \Phi(\mathbf{x}) = \begin{bmatrix} \Phi_1(\mathbf{x}) \\ \vdots \\ \Phi_J(\mathbf{x}) \end{bmatrix} \quad (9.11)$$

og behandle $\Phi(\mathbf{x})$ som om det var inputvektoren.⁴ De individuelle transformasjonene $\Phi_j(\cdot)$ betegnes som basistransformasjoner og verdimensjonen til $\Phi(\cdot)$ betegnes som feature space.⁵ Dette gir oss et nytt hypoteserom

$$\mathcal{H}_\Phi = \{l \circ \Phi : \Phi : \mathbb{R}^K \rightarrow \mathbb{R}^J \text{ og } l \text{ er lineær funksjon } l : \mathbb{R}^J \rightarrow \mathbb{R}\}. \quad (9.12)$$

Det empiriske risikominimeringsproblemet kan da skrives som

$$\hat{\gamma} = \arg \min_{\gamma \in \mathbb{R}^J} \frac{1}{N} \sum (y_n - \gamma' \Phi(\mathbf{x}))^2 \quad (9.13)$$

som vi løser for gitt transformasjon $\Phi(\cdot)$. I praksis er valg av $\Phi(\cdot)$ og dermed form på feature space vi søker over en viktig del av risikominimeringsproblemet. Det er bias-varians tradeoff og vi finner beste kandidat med kryssvalidering.

En mye brukt transformasjon er polynom for å modellere ikke-lineær sammenheng. Fra Weierstrass' theorem vet vi at kan oppnå vilkårlig god tilnærming av kontinuerlig funksjon med polynom av tilstrekkelig høy orden. Det kan derfor være tilforlatelig å gjøre en transformasjon

$$\Phi(x_n) = \begin{bmatrix} x_n^0 \\ x_n^1 \\ \vdots \\ x_n^J \end{bmatrix}, \quad \gamma' \Phi(x_n) = \sum_{j=0}^J \gamma_j x_n^j \quad (9.14)$$

men skal se at det finnes bedre måter å modellere ikke-lineær sammenheng dersom andre orden ikke er tilstrekkelig til å fange mønster.

Valg av featurespace

I statistisk læring er målet vårt enklere å måle fra data og vi kan bruke algoritme til å finne hvilken kombinasjon som generaliserer best.⁶

1. Kan teste alle mulige kombinasjoner, men det blir fort ganske mange... 2^K kombinasjoner
2. Greedy algoritme (forward/backward selection). I hvert steg legger til variabel som

⁴Hvis vi tolker de estimerte koeffisientene så vil vi ofte se på endring i forhold til opprinnelig input. Skal se på tolkning senere.

⁵I praksis er det ofte slik at basistransformasjoner bare avhenger av verdi til én av komponentene i inputvektoren, men kan avhenge av flere verdier eller ingen. Et grunnleggende eksempel er konstantledd, $\Phi(\mathbf{x}) = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$.

⁶Som målt ved kryssvalidering i henhold til metric eller andre kriterier som pålegger straff for fleksibilitet (justert R^2 , AIC, BIC,...). Vet ikke hvordan jeg utleder disse kriteriene eller hvorfor jeg skulle ønske å bruke det over kryssvalidering.

fører til størst reduksjon i RSS , deretter ta kryssvalidering.

Jeg synes fremgangsmåten over virker ganske slitsom. Bedre å kaste alt inn og la regulariseringsparameter ta seg av problemet.

9.3.2 Regularisering

Vi vil ta hensyn til at det er sannsynlighet for utfall som vi ikke observerer i det partikulære utvalget som vi trener modellen på, og at utfallet til verdi vi ikke observerer sannsynligvis er ganske like de nærmeste naboene vi observerer. Vi oppnår dette ved å straffe høye koeffisientverdier.

Vi kan se litt intuisjon for dette ved å betrakte to features som er høyt korrelert. Ettersom prediksjon er vektet gjennomsnitt av features så kan vi oppnå samme \hat{y} ved å skalere opp den ene koeffisienten og ned den andre. Dette kan gi bedre føyning i utvalget, men gjort at både \hat{y} og koeffisient blir ustabile. Vi kan gjøre det mer stabil ved å presse begge mot null. Dette oppnår vi ved å legge til et straffeledd i tapsfunksjonen. Har tre ulike måter i implementere det på i lineær regresjon.

Ridge

Kostnadsfunksjonen er

$$C(\theta) = MSE(\theta) + \alpha \sum_{k=1}^K \theta_k^2 \quad (9.15)$$

$$= MSE(\theta) + \alpha \mathbf{w}'\mathbf{w} \quad (9.16)$$

Merk at vi ikke straffer konstantleddet. Må standardisere features før vi regularisere slik at det ikke avhenger av måleenhet. Finnes en closed form løsning, men vet ikke hvor interessant det er.⁷

Lasso

Kostnadsfunksjoen er

$$C(\theta) = MSE(\theta) + \alpha \sum_{k=1}^K |\theta_k| \quad (9.17)$$

Bruker L_1 -norm i stedet. Vet ikke hvordan jeg kan skrive det på matrisiform. Fordelen med denne regulariseringen er at det setter koeffisienter lik 0 slik at den velger ut de relevante featurene. Her er marginalgevinsten ved å redusere koeffisient konstant. Kan også illustrere forskjell mellom ridge og lasso ved å sette det opp som betinget optimeringsproblem.

⁷Kan motivere at det har gode numeriske egenskaper som gjør inverse mer stabil eller noe sånt.

Elastic net

Kostnadsfunksjon har vektet gjennomsnitt av de to ulike regulariseringene,

$$C(\theta) = MSE(\theta) + \alpha \sum_{k=1}^K \theta_k^2 + (1 - r)\alpha \sum_{k=1}^K |\theta_k| \quad (9.18)$$

hmhm.

9.4 Andre regresjonsmetoder

Skal nå se på andre regresjonsmetoder. Tror splines er mest aktuell av disse.

9.4.1 Splines

9.4.2 Ikke-parametrisk regresjon

I lineær regresjon antar vi at $E[y|\cdot] := f(\cdot)$ er kjent opp til ukjent parameter. Vet ikke hvordan den ser ut på forhånd, men kan tilnærme arbitrære kontinuerlige funksjoner med basistransformasjoner og bruke kryssvalidering til å vurdere ekstern validitet. Dessuten har jeg sett av vi kan bruke splines til å partisjonere inputrommet slik at vi får mer fleksibilitet til å fange opp lokale sammenhenger. Sånn sett er parametriske metoder rimelig fleksible samtidig som de i prinsippet er mulige å tolke gjennom den parametriske representasjonen.

Har noen ikke-parametriske metoder som også kan være relevant å bruke i økonometri, men tror relevansen er avgrenset til regresjonsdiskontinuitet der det viktig å fange eksakte funksjonelle relasjonen. Metodene generaliserer veldig dårlig til input i flere dimensjoner; både fordi det er vanskelig å kommunisere den funksjonelle formen dersom den ikke kan visualiseres grafisk og fordi dimensjonalitetens forbannelse gjør at lokale metoder fungerer dårligere.

K-nærmeste naboer

Dette er en ganske direkte, naturlig og intuitiv estimator for $f(\mathbf{x}) = E[y|\mathbf{x}]$. Estimerer gjennomsnitt lokalt. Nedside: Ikke parametrisk representasjon av funksjonen, må ha alt data i minnet for å gjøre nye prediksjoner, fungerer dårlig i høy dimensjon.

Kernelmetoder

Sieve?

9.4.3 Kvantilregresjon

Den τ 'te kvantilen til en variabel y med cdf F er gitt ved $Q(\tau)$ der

$$\tau = F(Q(\tau)) \implies Q(\tau) = F^{-1}(\tau) \quad (9.19)$$

Vi må utvide definisjonen til å håndtere at ikke alle F er monotont voksende slik at den inverse ikke er definert på hele verdimengden til F .

$$Q(\tau) = \inf\{t : F(t) \geq \tau\} \quad (9.20)$$

Dette er den vanlige definisjonen, men vi kan også definere det som løsningen på et minimerings problem som involverer forventningsverdi av parametrisert funksjon av y . Skal da se at vi kan få det inn i ERM-rammeverket og kan estimere kvantiler fra utvalgsanalog.

$$L_\tau(y, \xi) = |(y - \xi)(\tau - I\{y < \xi\})| \quad (9.21)$$

$$= \textit{piecewise} \quad (9.22)$$

Kan vise at

$$Q(\tau) = \min_{\xi} \mathbb{E}[L_\tau(y, \xi)] \quad (9.23)$$

hmm... hm.

9.5 Klassifikasjon

I klassifikasjon angir utfallvariabelen hva slags kategori observasjonen tilhører. I binær klassifikasjon er det to kategorier. Det er vanlig å betegne den éne som *positiv* kategori og la den ta verdi 1, og den andre kategorien er *negativ* med verdi 0. Med flere kategorier kan utfallsvariabelen ta verdi $y \in G = \{1, \dots, K\}$. Merk at siden tallene bare er kode for kategori kan vi ikke bruke de numeriske verdiene til å si noe om avstanden mellom ulike kategorier eller deres rangering.

Fremgangsmåten i klassifikasjonsproblemer har likevel mange fellestrekk med regresjon. I praksis vil vi gjerne lære funksjoner som tar verdi $p_k(\mathbf{x}) := P(y = k|\mathbf{x})$. Disse betingede sannsynlighetsfunksjonene angir sannsynligheten for at en observasjon tilhører de ulike kategoriene gitt andre observerte egenskaper. De må oppfylle egenskapene til sannsynlighetsfunksjoner slik at $p_k(\mathbf{x}) \in [0, 1]$ for alle k og $\sum_k p_k(\mathbf{x}) = 1$. Merk at i binær

klassifikasjon er det tilstrekkelig å lære p_1 siden $p_0(\mathbf{x}) := 1 - p_1(\mathbf{x})$. Med flere kategorier er det vanlig å estimere ut fra *one-versus-rest* og eventuelt skalere slik at de summerer til én.

Vi vil også ha en hypotesefunksjon som tar verdi $h(\mathbf{x}) \in G$ og angir predikert kategori. Det kan vi lage ved å plassere observasjonen i kategorien med høyest predikert sannsynlighet,

$$h(\mathbf{x}) = \arg \max_{k \in G} p_k(\mathbf{x}) \quad (9.24)$$

Spesielt i binær klassifikasjon kan det være aktuelt å velge en *terskelverdi* k og plassere i positiv kategori dersom predikert sannsynlighet overstiger denne verdien,

$$h(\mathbf{x}) = I\{p(\mathbf{x}) > k\} \quad (9.25)$$

der valg av k avhenger av kostnad ved ulike typer feilprediksjon. Dette impliserer en partisjonering av inputmengden i delmengder som predikerer ulike kategorier og grensene mellom delmengdene er decision boundary

$$D(h) = D(p, k) = \{x : p(x) = k\} \quad (9.26)$$

som både avhenger av den estimerte betingede sannsynligheten og valg av threshold.

Empirisk risikominimering

Litt usikker på om jeg kan få klassifikasjon inn i rammeverket med empirisk risikominimering. Skal se at jeg kan bruke såkalt log-loss tap til å lære parameter i logistisk regresjon og betrakte det som empirisk risikominimering. Men jeg bruker jo annet kriterium til å vurdere risiko til h . Så er ikke like direkte kobling som i regresjon ...

9.5.1 Logistisk regresjon

Logistisk regresjon har en del fellestrekk med lineær regresjon og deler gode egenskaper, blant annet at det gir parametre som kan tolkes.

Utfallsvariabelen kan ta to verdier i $\{0, 1\}$ og den er da nødvendigvis bernoulli-fordelt. Fordelingen kan være betinget av \mathbf{x} slik at parameteren p er en funksjon av \mathbf{x} og det kan være ulike bernoulli-fordelinger for de ulike \mathbf{x} -verdiene, $y|\mathbf{x} \sim \text{bernoulli}(g(\mathbf{x}))$. Denne funksjonen g må tilfredstille $g(\mathbf{x}) \in [0, 1], \forall \mathbf{x}$. I praksis vil vi parametrisere funksjonen med $g(\mathbf{x}'\beta)$ slik at vi kan si noe om hvordan betinget sannsynlighet endrer seg når vi endrer input. Vi trenger en funksjon g som transformerer tallinjen til $[0, 1]$. Kumulative

fordelingsfunksjoner har denne egenskapen, og de to vanlige valgene er

$$g(\mathbf{x}) = \begin{cases} \Phi(z) = \int_{-\infty}^x (2\pi)^{0.5} \exp\{-s^2/2\} ds \\ \Lambda(z) = \frac{e^z}{1+e^z} \end{cases} \quad (9.27)$$

der første kalles probit og andre logit. Kan utvide til flere kategorier ved å lage egen parametervektor for hver kategori, $\theta^{(k)}$, $k = 1, \dots, K$. Den predikerte sannsynligheten for at input \mathbf{x}_n tilhører kategori j er da

$$\hat{p}_j = \frac{\exp(\theta^{(j)} \mathbf{x}_n)}{\sum_{k=1}^K \exp(\theta^{(k)} \mathbf{x}_n)} \quad (9.28)$$

Dette har visstnok noe med *softmax* og *cross entropy* å gjøre, men det må bli annen dag.⁸

Tolke koeffisienter

Merk nå at β er fra den underliggende latente modellen og ikke har noen opplagt tolkning. Det vi er interessert i er hvordan sannsynligheten for $P[y = 1|\mathbf{x}]$ avhenger av \mathbf{x} . For kontinuerlige variabler kan vi bruke kjerneregel til å derivere uttrykket,

$$\frac{\partial}{\partial x_k} F(\mathbf{x}'\beta) = \frac{\partial F(u)}{\partial u} \beta_k \quad (9.29)$$

$$= f(\mathbf{x}'\beta) \beta_k \quad (9.30)$$

Vi kan merke at effekt partiell effekt på betinget sannsynlighet alltid har samme fortegn som β_k siden $f(\cdot)$ er sannsynlighetstetthet, men størrelsen avhenger av hvor vi evaluerer \mathbf{x} . Vi kan betrakte $f(\mathbf{x}'\beta)$ som en skaleringsfaktor. Tre vanlige valg av skaleringer er

1. Plugge inn noen verdier. Interessant dersom vi har noen få dummies, men i praksis vil vi ofte ha enklere sammendragsmål.
2. Partial effect at average (PEA): Plugger in $\bar{\mathbf{x}}$. Litt problem dersom har transformerte variabler, siden tar gjennomsnitt etter transformasjon.
3. Average partial effect (APE): Evaluerer i hver \mathbf{x}_n som jeg observerer i utvalg og tar gjennomsnitt, $\frac{1}{N} \sum_n f(\mathbf{x}'_n \beta) \beta_k$.⁹

For variabler som er diskret er ikke den deriverte en meningsfull størrelse. Vi tar da differanse i verdi, $(x_k + 1) - x_k$, men det avhenger fortsatt av verdi til andre variabler.

⁸Også et poeng at vi kan kjøre log-reg som one-versus-all, men det skal i teorien være mulig å tolke $\theta^{(k)}$ fra multinomial... må prøve å få dette operativt i økonometri-delen, tror Cameron og Trivedi er best på dette.

⁹Dette blir da hele uttrykket for partial effect og ikke bare skaleringsfaktoren.

Kan bruke samme fremgangsmåter som over, f.eks blir APE:

$$APE(x_k) = \frac{1}{N} \sum_n [f(\mathbf{x}_{n,-k}\beta_{-k} + \beta_k(x_k + 1)) - f(\mathbf{x}'_n\beta)] \quad (9.31)$$

der $\mathbf{x}_{n,-k}\beta_{-k}$ er vektorene med de resterende $K - 1$ variablene.

9.5.2 Bayesianske metoder

Bruker bayes regels til å estimere betinget sannsynlighet

Diskriminantanalyse

Bayes-regel gjør at vi kan estimere betinget sannsynlighet på en annen måte:

$$f(y = k|\mathbf{x}) = \frac{f(\mathbf{x}|y = k)\mathbb{P}\{y = k\}}{f(\mathbf{x})} \quad (9.32)$$

$$= \frac{f_k(\mathbf{x})\pi_k}{\sum_j f_j(\mathbf{x})\pi_j} \quad (9.33)$$

$$\propto f_k(\mathbf{x})\pi_k \quad (9.34)$$

denne metoden er enklere å bruke på flere kategorier og dette må jeg si litt om.. Uansett, må nå i stedet velge parametrisk klasse til $\mathbf{x}|y$ og får kvadratisk discriminant analysis hvis jeg antar at $\mathbf{x}|y = k \sim N(\mu_k, \Sigma_k)$ og lineær hvis jeg i tillegg antar at $\Sigma_k = \Sigma, \forall k$. Med utgangspunkt i dette kan jeg finne decision boundary som funksjon av enkle størrelser som jeg kan estimere. Skal si mer om dette senere.

Naiv bayes

For å implementere bayes-regel kan vi estimere

$$f_k(\mathbf{x})\pi_k. \quad (9.35)$$

Har sett at lda/qda gir en måte å gjøre dette på. Den metoden har ganske sterke parametriske antagelser. Vil bruke svakere antagelser, men problem å estimere betinget simultanfordeling $f_k(\mathbf{x})$. Blir mye enklere dersom vi antar at de er uavhengige slik at

$$f_k(\mathbf{x}) = \prod_j f_{kj}(x_j) \quad (9.36)$$

mer om dette senere.

9.5.3 KNN

En ikke-parametrisk metode for å estimere betinget sannsynlighet. Bruker relativ av kategorier i nabolag til \mathbf{x} som estimat på betinget sannsynlig, der nabolaget $N_K(\mathbf{x})$ består av K observasjoner med minst $\|\mathbf{x}_n - \mathbf{x}\|$.

$$\hat{P}(y = j|\mathbf{x}) = \frac{1}{k} \sum_{n \in N_K(\mathbf{x})} I\{y_n = j\} \quad (9.37)$$

Hvis målet er å minimere feilrate blir klassifiseringsregelen h å predikere kategori som er mode i nabolag.

9.5.4 Support vector machines

9.5.5 Beslutningstrær

Beslutningstrær er en fleksibel og transparent metode som kan brukes til både regresjon og klassifikasjon. Metoden går ut på å partisjonere inputmengden og bruke mode eller gjennomsnitt i hver delmengde som predikert kategori for observasjon med input der. Formelt finner vi R_j der

$$\cup_{j=1}^J R_j = \mathcal{X}, \quad R_j \cap R_k = \emptyset, j \neq k \quad (9.38)$$

og predikert verdi kan representeres parametrisk i lineær modell som

$$\hat{y}_n = \sum_j I\{\mathbf{x}_n \in R_j\} \hat{\beta}_j \quad (9.39)$$

der $\hat{\beta}_j = \text{avg}(\{y_n : \mathbf{x}_n \in R_j\})$. Partisjonering kan representeres med et tre som er en special case av en graf.¹⁰ Treet består av nodes og koblinger mellom nodes som vi kan betegne som greiner (eller *branches*). Det begynner i root-node og splitter i to eller flere child nodes ut i fra verdi av en variabel. Nodes som ikke har childs betegnes som blader (eller *leaves*). Bladene på treet utgjør den endelige partisjoneringen og det er dette vi bruker til å gjøre prediksjoner.

For nye inputs kan vi bevege oss gjennom treet. Dette gjør estimatoren transparent siden vi ser hvilke inputs som fører til hvilke inputs. Mer spesifikt kan vi både se hvilke inputs som er viktige for å forklare forskjeller i observert kategori (tidlig split) og hvilken retning det påvirker predikert sannsynlighet for kategori. Vi kan observere predikert sannsynlighet i hver node og se hvordan den endres mens vi beveger oss i treet gjennom å gi gradvis mer informasjon om input til observasjon.

¹⁰Må ha litt formell definisjon av graf som jeg får fra algoritme/datastrukturer. Poeng at det har nodes og vertices.

En nedside med beslutningstrær er at de er veldig sensitive for treningsdata. Små endringer i data den blir opplært på kan få store konvekvenser for partisjoneringen (som bestemmer decision boundary og prediksjoner). Dette har til dels sammenheng med at den lager rektangulære partisjoner som er ortogonalt på aksene. Vi skal senere se at vi kan glatte ut boundaries ved å kombinere mange trær i en såkalt tilfeldig skog. Først skal vi se litt kort på algoritme for å konstruere hvert enkelt tre.

Algoritme for å konstruere trær

I hver node så søker vi over alle mulige cut-points for å finne partisjonering som fører til størst mulig reduksjon i såkalt *impurity*. Vi vil at labels i hver delmengde skal være mest mulig homogene. De to vanligste målene på impurity er *gini* og *entropy*. I praksis bruker vi en greedy algoritme som søker over grid og kalkulerer reduksjon i impurity for hver punkt i grid og bruker dette til å partisjonere. Deretter gjøres dette rekursivt helt til det ikke lenger er mulig å redusere impurity (alle inputs har enten samme features eller samme labels) eller til treet har nådd en spesifisert grense for dybde.

Jeg tror det er best å vokse ut hele treet og deretter trimme (*prune*) det ex-post for å fjerne oppdelinger som ikke fører til bedre fit out-of-sample. I praksis tror jeg vi bruker maks-dybde som regularisering selv om dette ikke er optimalt...

9.6 Ensemble

Vi kan oppnå bedre prediksjoner ved å kombinere flere estimatorer. Litt av intuisjonen bak dette er at idiosynkratiske feil jevner seg ut når vi tar gjennomsnitt.¹¹ Dette argumentet bygger på at det er variasjon i prediksjonene til de ulike estimatorene. Den enkleste måten å oppnå dette er å bruke algoritmer til å trene opp estimatorer på treningsdata og deretter bruke en avstemming til å predikere nye inputs. Vi kan da enten bruke simpel majoritet (såkalt hard voting) eller vi kan vekte ut i fra den predikerte sannsynligheten til de ulike estimatorene (såkalt soft voting). I praksis bruker vi to andre fremgangsmåter for å skape variasjon: vi sampler fra treningsdata for å skape variasjon i data eller vi trener estimatorer sekvensielt der det blir lagt større vekt på observasjonene som ble feilpredikert av forrige estimator.

9.6.1 Bagging og tilfeldig skog

Bagging er kort for bootstrap aggregering. Ved å sample med replacement fra treningsdata så sampler vi fra empirisk fordeling. Dette innfører litt mer bias, men ved å ta gjennom-

¹¹Kan koble til forsikring... Selv om hver enkelt estimator kun predikere riktig kategori i 51% av tilfellene vil andel som predikere riktig konvergere i sannsynlighet mot 51% slik at majoriteten tar riktig i 100% av tilfellene dersom de er uavhengige.

snitt av estimatorene så kan vi redusere varians. Fremgangsmåten er spesielt egnet for beslutningstrær siden de er sensitive for treningsdata. I praksis bruker vi derfor såkalte tilfeldig skog som er baggete beslutningstrær med litt ekstra triks for å oppnå mer varians. Det er for eksempel vanlig å avgrense mengden av variabler den kan bruke til å partisjonere til en tilfeldig delmengde.

Den tilfeldige skogen er litt mindre transparent enn et enkelt beslutningstre siden vi ikke kan følge hvordan input beveger seg langs greinene, men vi kan få et mål på feature importance ut fra gjennomsnittlig bidrag til reduksjon i impurity fra tra trærne i skogen. Det er også en fordel at trærne kan blir trent opp parallelt.

9.6.2 Boosting

Dette er en alternativ fremgangsmåte der estimatorene blir trent opp sekvensielt og som dermed ikke kan paralleliseres. I stedet for å bruke mange unbiased estimatorer med høy varians så forsøker vi å sekvensielt redusere bias ved å legge mer vekt på observasjonene som ble feilpredikert av forrige estimator. Det finnes i hovedsak to fremgangsmåter for å booste: Ada(ptive)Boost og gradient boosting.

Adaptive Boosting

Eksplisitt endringer av vekter i kostnadsfunksjon, må gjør algoritmen formell en annen gang.

Gradient Boosting

Fitter residual av forrige, vet ikke om det fungerer på klassifikasjon.

9.6.3 Stacking

Kan også forsøke å lære den beste mulige måten å kombinere de ulike estimatorene... Hard vs soft voting osv; alt kan læres og valideres...

9.7 Vurderingskriterier

9.7.1 Confusion matrix

I klassifikasjon er ikke vurderingskriteriet like entydig som i regresjon. Det mest intuitive kriteriet er *accuracy* som angir andelen av observasjoner som blir plassert i riktig kategori, men dette er ofte ikke et godt mål på hvor egnet modellen er. For det første kan vi med ubalanserte kategorier oppnå høy treffsikkerhet ved å alltid predikere majoritetskategorien som ikke er så nyttig. Dessuten er det ofte ulike kostnader assosiert med ulike *typer* feil.

Vi kan kategorisere ulike typer feil gjennom en såkalt *confusion matrix* som deler inn observasjoner ut fra faktisk kategori og predikert kategori. Med binær klassifikasjon gir dette fire muligheter og vi bruker dette til å lage vurderingskriterier

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} \quad (9.40)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (9.41)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (9.42)$$

Anta nå at vi bruker algoritmen til å diagnostisere om personer har en gitt sykdom (f.eks. covid). Presisjon til algoritmen angir andelen av de som tester positiv som faktisk er smittet. Denne kan vi få arbitrært høy gjennom å kun diagnostisere de som helt klart er syke. Recall (sensitivitet) angir derimot andelen av de som faktisk er syke som tester positivt. Hvis testen er lite sensitiv så er det mange syke som vil gå under radaren. Vi kan igjen få denne arbitrært høy ved å si at alle som tester seg er syke.

Isolert sett gir disse kriteriene ikke noe godt mål siden vi kan lage rimelig trivielle algoritmer som maksimerer kriterium uten å være nyttig. Et alternativ kan være å ta et gjennomsnitt av presisjon og sensitivitet. F1-score tar harmonisk gjennomsnitt.¹² Dette kan gi et greit sammendragsmål for å sammenligne ulike algoritmer, men i praksis er det bedre å undersøke tradeoff mellom presisjon og sensitivitet for å finne den beste balansen til vårt formål.

9.7.2 Presisjon vs Recall trade-off

Algoritme gir gjerne et mål på hvor sikker den er at en observasjon tilhører en kategori.

$$P(\widehat{y_n = 1} | \mathbf{x}_n) = \hat{p}_n \quad (9.43)$$

$$\hat{y} = I\{\hat{p}_n \geq k\} \quad (9.44)$$

ved å øke *threshold* k kan vi øke presisjon og redusere recall. Vi kan visualisere dette gjennom å tegne $(k, \text{Pres}(k))$ og $(k, \text{Recall}(k))$ i et diagram. Det er mer nyttig å tegne output fra $f : k \mapsto (\text{Pres}(k), \text{Recall}(k))$ som angier en såkalt *mulighetskurve* med recall vi kan oppnå for gitt presisjon. Hvis kurven er bratt så må vi gi opp masse recall for å oppnå litt mer presisjon. Vi kan bruke denne kurven til å finne k som korresponderer med balansen av presisjon og recall som vi foretrekker. Kan også tegne kurver fra ulike algoritmer. Ulike algoritmer kan ha ulik performance på ulike deler av kurven. Så dersom vi er veldig opptatt av å ha f.eks. over 90% recall, så kan vi se hvilken algo som oppnår

¹²litt usikker på hvorfor vi ikke tar aritmetisk snitt. Harmonisk venter slik at det blir større straff dersom én av de er lav..

høyest presisjon i det intervallet.

Et annet mye brukt mål er den såkalte *ROC-kurven* som plotter True Positive Rate (?) vs False Positive Rate (?). Igjen kan vi undersøke kurven eller bruke areal under kurven (*AUC*) som et sammendragsmål for valg av algoritme og treshold k .

Kapittel 10

Læring uten tilsyn

Data uten labels.

10.1 Dimensjonalitetsreduksjon

Vi kan ønske å finne en lavere dimensjonal representasjon av data som bevarer mest mulig informasjon.¹ Hva som utgjør informasjon har ikke en eksakt definisjon og tenker at det avhenger litt av kontekst.² En måte å redusere dimensjonene er å droppe variabler som vi anser som irrelevante. Vi skal nå se på nye fremgangsmåte som i stedet konstruerer nye variabler som fanger opp informasjon fra de eksisterende variablene, slik at vi kan finne lavere dimensjonal representasjon som bevarer mest mulig info.

10.1.1 Principal component analysis

Metode for å finne et d -dimensjonalt underrom for datasett med k variabler der $d \leq k$. Vi vil bevare mest mulig informasjon. For hver d finner vi derfor underrommet som bevarer mest mulig av variansen. Dette er ekvivalent med å at det minimerer MSE av projeksjon av data ned på underrommet.³ Det som er veldig fint er at greedy algoritme også gir optimal løsning: kan finne k principale vektorer $[\mathbf{v}_1, \dots, \mathbf{v}_k]$ og for $d < k$ så velger vi bare delmengde som består av d første.

For å finne disse principale komponentene kan vi bruke singulærverdidekomposisjon som er en måte å faktorisere en matrise,

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}' \quad (10.1)$$

¹Når vi klassifiserer siffer så er mange av pixlene hvite for alle bildene slik at de ikke er så informative. Vi kunne droppet disse og få færre variabler (dimensjoner) per observasjon uten at vi taper info.

²Ren unsupervised eller preprocessing i supervised... hvilke variabler er informative i en regresjon liksom.

³Tror jeg kan vise dette med dekomponering av varians... vil knytte til ting jeg kan om prosjektering fra før, men blir litt anderledes siden jeg nå projekterer en matrise i stedet for vektor...

der $\mathbf{V} = [\mathbf{v}_1 \dots \mathbf{v}_k]$. Vi finner projektering ned på underrommet som er utspent av komponentene med

$$\hat{\mathbf{X}}_d = \mathbf{X}\mathbf{V}_d \quad (10.2)$$

der $\mathbf{V}_d = [\mathbf{v}_1 \dots \mathbf{v}_d]$. Vi kan også forsøke å gjøre invers transformasjon for å forsøke å gjenskape det opprinnelige datasettet med

$$\hat{\mathbf{X}} = \hat{\mathbf{X}}_d \mathbf{V}_d' \quad (10.3)$$

Veldig kjekt at vi kan finne andel av forklart variasjon til hver komponent slik at vi kan plote dette og bruke til å bestemme hvor mange dimensjoner vi trenger til å representere informasjonen i datasettet.

10.1.2 Andre metoder

Kan bruke noe kernel eller manifold ... for noen representasjoner er det vesentlige greier som går tapt når vi projekterer på underrom.

10.2 Clustering

Metoder som forsøker å konstruerer clusters og plassere observasjoner i kategori. Vi vil at de skal være homogene innad og ha distinksjon mellom andre cluster. For gitte cluster blir det litt sånn som klassifisering bare at vi ikke kjenner tolkning til cluster-label.. Kan skille mellom hard cluster og soft cluster, der sistnevne beregner sannsynlighet for at observasjon tilhører de ulike clusterene.

10.2.1 K-means

Finner K *centroids* og angir label til observasjon ut fra nærmeste centroid. Algoritmen er enkel: bruk en tilfeldig initialisering, label data og oppdater plassering av centroids slik at det er i midten av observasjonene som ble angitt til den centroiden. Deretter angi labels ut fra oppdatert posisjon til centroids og fortsett slik til det konvergerer. Det er en utfordring at det kan konvergere mot lokal minimum som ikke er globalt optimalt, men dette kan vi håndtere ved å kjøre flere ganger og velger det som minimerer avstand innad i clusterene.

Den større utfordringen er valg av K . Vi har noen vektøy for å gjøre informert valg om dette: såkalt inertia og silhouette, men må se på dette en annen gang.

10.3 Tetthetsestimering

Vi har sett at inferens handler om å estimere egenskaper til fordelingen som genererte utvalget vi observerer. Nå skal vi forsøke å estimere selve fordelingen. Dette har vi vært inne på i MLE der vi avgrensner oss til å betrakte en enkel parametrisk familie slik at problemet reduseres til å estimere parametre som karakteriserer denne. Svakheten med denne fremgangsmåten er at det gir dårlig resultat dersom den gitte parametriske familien ikke er fleksibel nok til å tilnærme den sanne tettheten til fordelingen som genererte data. Vi skal derfor utvide til ikke-parametriske metoder som kan tilnærme tettheter med vilkårlig form (så lenge de kontinuerlige). Deretter skal jeg se på såkalte *mixture models* som kan betraktes som et semi-parametrisk kompromiss; har både struktur og fleksibilitet.

For å strukturere diskusjonen og vurdere de relative styrkene og svakhetene til de ulike tilnærmingene vil jeg innføre mål på avstand mellom tettheter.

10.3.1 Mål på avstand mellom tetthetsfunksjoner

Jeg tenker vi kan betrakte en mengde av funksjoner, for eksempel alle funksjoner f der $f : \mathbb{R}^d \rightarrow \mathbb{R}$. Jeg vil si noe om avstand mellom elementene i mengden. Hva vil det si at to funksjoner er nærme hverandre? Det blir vel dersom de tar omtrent like funksjonsverdier for de ulike inputene.

For å formalisere dette kan vi definere en p -norm som er analog til andre vektorrom,

$$\|f\|_p := (|f|^p)^{1/p} := \left(\int |f(\mathbf{s})|^p d\mathbf{s} \right)^{1/p} \quad (10.4)$$

Dette gir et mål på avvik mellom funksjoner

$$d_p(f, g) = \|f - g\|_p. \quad (10.5)$$

Når vi har definert normen kan vi betrakte mengden som et L_p rom bestående av funksjoner med definert p -norm, $L_p = \{f : \|f\|_p < \infty \text{ og der } f : \mathbb{R}^d \rightarrow \mathbb{R}\}$. De vanligste valgene er $p = 1$ eller $p = 2$. Det kan være enklere å jobbe med L_2 siden normen er analog til eukledisk avstand slik at resultatet fra vanlige vektorrom kan generaliseres, men L_1 gir nok et bedre og mer intuitivt mål på avstand mellom funksjoner.

Fremstillingen over gjelder for funksjoner generelt. Vi betrakter kun mengden av funksjoner som oppfyller egenskapene til tetthetsfunksjoner. For tetthetsfunksjoner har vi allerede sett at *total variation distance* som angir den maksimale forskjellen i sannsynlighet to tetthetsfunksjoner tillegger en hendelse,

$$TVD(f, g) := \sup_{A \in \mathcal{B}(\mathbb{R})} \left| \int_A f - \int_A g \right|, \quad (10.6)$$

og det kan vises at

$$TVD(f, g) = \frac{1}{2} \|f - g\|_1 \quad (10.7)$$

Forventet avstand og konsistens

Jeg tenker at vi kan ha et utfallsrom Σ som består av de ulike tetthetsfunksjonene som kan bli realisert avhengig av verdiene i utvalget, $\Sigma = \{\hat{f}(\omega) : \omega \in \Omega\}$. Videre tenker jeg at vi kan definere en tilfeldig variabel på megden, $\gamma(\hat{f}) = d_1(\hat{f}, f)$ og finne forventningsverdi til denne med hensyn på fordelingen f . Dette kan vi evaluere med simulering fra et utvalg på N uavhengige observasjoner fra f .

Alternativt kan vi bruke asymptotisk teori. En følge med tilfeldige tetthetsfunksjoner $(\hat{f}_N)_{N \in \mathbb{N}}$ er L_p konsistent for f hvis

$$\|\hat{f}_N - f\|_p \xrightarrow{p} 0 \quad (10.8)$$

når $n \rightarrow \infty$. Dersom modellen er feilspesifisert slik at $P_0 \notin \{P_\theta : \theta \in \Theta\}$ så har avviket en nedre begrensning

$$\delta(P_0, \hat{P}_\theta) = \inf_{\theta \in \Theta} \|P_0 - \hat{P}_\theta\|_p \quad (10.9)$$

10.3.2 Histogram

Vi observerer realiseringer X_1, \dots, X_N på utfallsrom $Z = [0, 1]$ fra en fordeling med tetthet f . En veldig naiv tilnærming er å bruke relativ andel av hver observasjon som estimat på fordelingen, $\hat{f}(x) = \frac{1}{n} \sum I\{X_n = x\}$. Det vil jo vanligvis være sannsynlighet for verdier vi ikke observerer i det gitte utvalget vårt, og dette gjelder spesielt siden vi antar at den sanne fordelingen er kontinuerlig. Histogram gir en litt bedre tilnærming. Vi finner en partisjonering av Z som er en mindre $\{B_1, \dots, B_K\}$ der $\cup B_k = Z$ og $B_k \cap B_j = \emptyset$ for $k \neq j$. Dette er bins med lengde $1/K$. Sannsynligheten for å få utfall i hver bin er $p_k = \int_{B_k} f(x) dx$ og estimator er $\hat{p}_k = \frac{1}{N} \sum I\{x \in B_k\}$. Dette gir en diskontinuerlig funksjon på Z som vi kan skrive som $\hat{f}(x) = \sum_k \hat{p}_k I\{x \in B_k\}$.

Hyperparameteren i histogrammet er antall bins som også bestemmer lengden på intervallene. Hvis det er for få bins klarer det ikke fange mønsteret i den sanne tettheten (høy bias), mens for mange bins gjør at estimat fra ulike utvalg blir veldig forskjellige (høy varians). Dette kan man til en viss grad ta på øyemål, men analogt til estimasjon av regresjonsfunksjon/betinget sannsynlighet kan vi definere en tapsfunksjon og forsøke å velge antall bins som minimiserer risiko. Tar det senere.

En nedside med histogram er at det gir en diskontinuerlig funksjon. Skal nå se en måte å utlede kontinuerlig tetthetsfunksjon.

10.3.3 Kernel density estimation

Intuisjonen er at vi vil spre litt tetthet rundt de observerte realiseringene, der tyngden avtar mer avstand. For å gjøre dette bruker vi en kernel funksjon $K(\cdot)$ med egenskaper $\int K(x)dx = 1$, $\int xK(x)dx = 0$, $K(x) = k(-x)$. For hver observasjon transformerer vi input til $\frac{x-x_n}{h}$ slik at tyngden er sentrert rundt x_n og der h er den såkalte *bandwidth* som justerer for raskt tyngden avtar... må omskrive litt senere. Som eksempel, la $K(\cdot)$ være standardnormalfordeling

$$K(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \quad (10.10)$$

og der

$$g_n(x) = \frac{1}{Nh} K\left(\frac{x - x_n}{h}\right) \quad (10.11)$$

$$) = \frac{1}{N} \frac{1}{\sqrt{2\pi}h^2} \exp\left(-\frac{(x - x_n)^2}{2h^2}\right) \quad (10.12)$$

og $\hat{f}(x) = \sum_n g_n(x)$. Ikke helt opplagt for meg hvorfor vi deler på h utenfor kernel, men ser at det gir mening for at det fortsatt skal gi normalfordeling..

gammelt

Kernel density estimators gir en alternativ fremgangsmåte for å estimere tetthetsfunksjonen under svakere antagelsen (ie. ikke avgrenset til spesifikk parametrisk klasse). Estimatoren er en den skalerte summen av N tetthetsfunksjoner som hver er sentrert på de observerte \mathbf{x}_n ($n = 1, \dots, N$). Det er en såkalt *bandwidth* som justerer spredningen på de individuelle tetthetsfunksjoner og dermed glattheten (*smoothness*) til summen. Med lavere bandwidth blir større del av tyngden konsentrert på rundt de observerte verdiene. Formelt kan vi skrive estimatoren som

$$\hat{f}_N(\mathbf{s}) = \frac{1}{Nh^p} \sum K\left(\frac{\mathbf{s} - \mathbf{x}_n}{h}\right) \quad (10.13)$$

10.3.4 Mixture models

Kapittel 11

Økonometri

For meg er økonometri synonymt med programevaluering.¹ Vi bruker data til å estimere effekt av behandling på utfall til individer.² For å kvantifisere dette vil vi ideelt sett observere utfallene til hvert av individene i en verden der de blir eksponert for behandling og i en verden uten behandling. Dessverre er dette umulig siden kun ett av tilfellene kan inntreffe, og vi kan dermed aldri kvantifisere individuelle behandlingseffekter. I praksis beregner vi gjennomsnittlige behandlingseffekter ved å sammeligne forskjell i gjennomsnittlig utfall til en behandlingsgruppe og en kontrollgruppe. For at denne observerte forskjellen skal gi et godt mål på behandlingseffekten må kontrollgruppen være en god *proxy* for det kontrafaktiske utfallet til behandlingsgruppen dersom de ikke ble eksponert for behandling.

11.1 Programevaluering

Det er rimelig å betrakte kontrollgruppen som en proxy dersom de to gruppene er omtrent like bortsett fra at den éne ble eksponert for behandling. I så fall kan observerte forskjeller i utfall tilskrives denne ene dimensjonen der gruppene er forskjellig.³ Det sentrale spørsmålet i programevaluering er hvorvidt dette er rimelig antagelse med de dataene som foreligger i analysen. Ettersom mange variabler som påvirker utfallet er uobserverte (og uobserverbare) kan det ikke testes med de gitte dataene, og må i stedet sannsynliggjø-

¹Økonometri omfatter også greier med (makroøkonomiske) tidsserier og estimering/kalibrering(?) av parametre i økonomiske modeller, men disse greiene vet jeg lite om. Det er forøvrig andre fagfelt som holder på med programevaluering og det er overlapp i problemstillinger og faglige tilnærminger. Det som kjennetegner økonometri er bruk av såkalte naturlige eksperiment og metoder for å analysere disse. I biostatistikk bruker de mer kontrollere eksperiment. Andre har mer naiv tilnærming for å isolere kausal effekt gjennom matching/justering for andre observerte egenskaper.

²Terminologi stammer fra medisinske eksperimenter. Det som betegnes som behandling kan være andre former for tiltak og reformer. Enhetene som blir eksponert for disse kan være aggregerte størrelser som foretak.

³Noe som også impliserer at det ikke ville vært observerte forskjeller dersom de ikke fikk ulik eksponering for behandling. Med begrensede utvalg vil det alltid være litt tilfeldig variasjon, men dette abstraherer vi stort sett vekk fra når vi diskuterer kausalitet.

res gjennom en beskrivelse av hvordan data er generert. Vi skal nå se på tre ulike typer beskrivelser av hvordan eksponering for behandling er bestemt.

Den første kategorien er tilfeldige eksperimenter der eksponering for behandling blir bestemt av forskere i henhold til en randomiseringsregel.⁴ Ettersom eksponeringen er tilfeldig vil det ikke være systematiske forskjeller mellom gruppene langs noen egenskaper, hverken observerte eller uobserverte. Tilfeldige eksperiment regnes som gullstandard i programevaluering siden antagelsen om at kontrollgruppen er god proxy er veldig troverdi. Det er likevel vesentlige begrensinger ved slike eksperiment. Det er mange interessante kausale spørsmål som ikke kan besvares på denne måten, enten fordi det er praktisk umulig, for dyrt eller uetisk. Det er også mange mange praktiske utfordringer, og det kan argumenteres for at de kan ha begrenset ekstern validitet.⁵ På tross av dette er det gjennomført en del eksperiment i stor skala og det har blitt gradvis mer fremtredende også i økonometri.⁶

En alternativ fremgangsmåte er å utnytte at ytre omstendigheter kan skape variasjon i behandling selv om det ikke er planlagt som et tilfeldig eksperiment.⁷ Det faktum at individene ikke selv velger egen eksponering for behandling gjør det ofte mer kredibelt at det ikke er systematiske forskjeller i uobserverte egenskaper. Et eksempel på en slik situasjon er at egenskaper ved institusjoner at eksponering for en behandling er bestemt ved om individ havner over eller under en noe abitrær *cut-off*.⁸ I den grad det er vanskelig for individ å strategisk velge side så blir behandling som om tilfeldig fordelt i populasjonen i nærheten av cut-off. Dersom hele grupper blir eksponert for ulike behandlinger avhengig av geografi (fordi ulike policy på ulik sted) eller fødselsalder (fordi endring i policy som rammer personer født etter gitt dato) har vi også verktøy for å sammenligne forskjeller og vurdere i hvilken grad det skyldes effekt av ulik eksponering for behandling. Vi kan også håndtere omstendigheter som skaper noe variasjon i behandling uten å bestemme det eksakt. Ved hjelp av såkalte instrumentelle variabler kan vi i store utvalg isolere variasjonen i behandling som skyldes den ytre omstendigheten.⁹

I fravær av slik eksogen variasjon kan det være fristende å stratifisere observasjonsdata og aggregere forskjell mellom behandling og kontroll innad i hvert strata. Selv om behandlings- og kontrollgruppen samlet sett er systematisk forskjellig kan vi konstruere

⁴Kan være tilfeldig på hele utvalget eller tilfeldig innad i strata definert av observert egenskap, for eksempel kjønn.

⁵For eksempel kan det være at folk oppfører seg annerledes hvis de vet at de er med på et eksperiment (den såkalte Hawthorne-effekten). Det er praktisk utfordring å få randomisert utvalg; for de som er med er det enkelt å dele inn i behandling og kontroll, men kan være systematisk skjevhet fra resten av populasjonen som vi vil generalisere til.

⁶Eksempler på store eksperiment i er STAR som undersøkte effekt av klassestørrelse på barns utfall og noe greier med effekt av helseforsikring på pasientenes utgifter i USA. Eksperimenter er viktig i atferdsøkonomi og har blitt viktig del av utviklingsøkonomi.

⁷Tror vi betegner det som eksogen variasjon.

⁸Noen eksempel er karakterkrav for å komme inn på skole og helsetiltak som avhenger av nyfødt barns vekt.

⁹Liker dårlig denne formuleringen siden IV i praksis bare er skalering av redusert form..

delutvalg der individer med ulik eksponering for behandling er omtrent like langs andre observerte egenskaper. Vi kan da estimere behandlingseffekter innad i hvert delutvalg og forsøke å aggregere dette til en gjennomsnittlig behandlingseffekt i populasjonen. Denne fremgangsmåten kan motiveres med at individer som er like langs observerte egenskaper forhåpentligvis også er ganske like langs uobserverte egenskaper som kan påvirke utfallet. Problemet er at det alltid er en grunn til at individene velger ulik eksponering for behandling innad i hvert strata og det er lite kredibelt at denne grunnen ikke også påvirker utfallet.¹⁰ For å publisere i gode tidskrift er det nødvendig å ha et forskningsdesign som isolerer eksogen variasjon i behandling. Ellers regnes det som lite troverdig at kontrollgruppen er proxy for det kontrafaktiske utfallet til behandlingsgruppen i fravær av behandling, slik at de observerte forskjellene i utfall ikke samsvarer med kausal effekt av behandling.¹¹

Jeg skal nå utlede et rammeverk som formaliserer idéen om kontrollgruppe som proxy.

11.1.1 Potensielle utfall

Vi har nå et rammeverk som lar oss beskrive relasjon mellom variabler og estimere dette fra data. Denne relasjonen består både av en eventuell kausal relasjon mellom variablene og spuriøs korrelasjon som følge av andre variabler som er korrelert med både utfall og forklaringsvariabler. Den kausale effekten kan defineres som differansen i de potensielle utfallene med og uten behandling. Det grunnleggende problemet er at kun én av tilstandene blir realisert for hver observasjon. Vi innfører notasjonen

$$y_i = \begin{cases} y_i^0, & D_i = 0 \\ y_i^1, & D_i = 1 \end{cases} \quad (11.1)$$

som også kan skrives som

$$y_i = y_i^0 + D_i(y_i^1 - y_i^0) \quad (11.2)$$

Det er ikke mulig å estimere individuell kausal effekt siden vi aldri kan observere de kontrafaktiske utfallene $y_i^1|D_i = 0$ og $y_i^0|D_i = 1$, men vi kan forsøke å estimere gjennomsnittlig effekt for en avgrenset populasjon ved å se på differansen i utfall til de som blir eksponert for behandlingen og kontrollgruppen som ikke blir eksponert. Intuisjonen bak denne sammenligningen er at utfallet til kontrollgruppen gir en proxy for det kontrafak-

¹⁰Individer er sånn omtrent rasjonelle og vi kan betrakte eksponering for behandling som løsning på et optimeringsproblem. Det er lite rimelig at forskjellene bare er tilfeldig. Det kan skyldes ulike preferanser: de som spiser vitaminer større preferanse for 'sunnhet' og vil gjerne dermed være sunnere uavhengig av eventuell behandlingseffekt. Eller kanskje de kompenserer for usunt kosthold. Uansett: vanskelig å isolere behandlingseffekt fra andre systematiske forskjeller.

¹¹Dette er til dels en konsekvens av den såkalte kredibilitetsrevolusjonen.

tiske utfallet til behandlingsgruppen slik at den observerte differansen tilsvarer differanse i potensielle utfall. Uten randomisering vil observert differanse bestå av både kausal effekt og seleksjonsskjevhet.

$$E[y_i|D_i = 1] - E[y_i|D_i = 0] = E[y_i^1|D_i = 1] - E[y_i^0|D_i = 1] \quad (11.3)$$

$$+ E[y_i^0|D_i = 1] - E[y_i^0|D_i = 0] \quad (11.4)$$

For at denne naive sammenligningen mellom behandling og kontroll skal isolere kausal effekt trenger vi randomisering av behandling. Dette sikrer at potensielle utfall er uavhengig av behandling. Sagt på en annen måte; observert behandling gir oss ikke noe informasjon om kontrafaktisk utfall.

$$(y_i^1, y_i^0) \perp\!\!\!\perp D_i \implies E[y_i^j|D_i] = E[y_i^j], j = 0, 1 \quad (11.5)$$

Vi kan utvide til setting der randomiseringen skjer betinget av observerbar egenskap¹²

$$(y_i^1, y_i^0) \perp\!\!\!\perp D_i | X \implies E[y_i^j|D_i, X] = E[y_i^j|X], j = 0, 1 \quad (11.6)$$

Vi kan da bruke matching eller regresjon til å estimere denne effekten, noe jeg skal se på senere. I praksis er dette som oftest lite troverdig siden det er en grunn til at observasjonene i hver kategori valgte ulik behandling og det er lite troverdig at dette ikke også påvirker potensielle utfall. Vi trenger derfor et forskningsdesign som skaper tilfeldig (eksogen) variasjon i behandling. Vi kan analysere variasjon i behandling D direkte dersom vi har randomisert eksperiment eller vi kan se på variasjonen som skyldes et instrument Z . En viktig kilde til eksogen variasjon er såkalte *naturlige eksperiment*. En viktig kilde til slike eksperiment er kunnskap om institusjonelle regler... som vi skal analysere med regression discontinuity. Senere skal jeg også se på paneldata som lar oss kontrollere for uobservert heterogenitet som er konstant over tid.

11.1.2 Matching

Matching er strategi for å estimere behandlingseffekt ved å konstruere undergruppe med samme covariates, finne forskjell i gjennomsnittlig utfall til behandling og kontroll innad i hver undergruppe og aggregere forskjellene ved å finne et vektet gjennomsnitt av forskjellene. La x være en diskret variabel og anta at CIA er oppfylt slik at $E[Y_i^j|D_i, X_i] = E[Y_i^j|X_i]$, $j = 0, 1$. Matching er fint siden vi kan knytte det direkte til CIA og finne

¹²For eksempel at tilfeldig utvalg av 40% av menn og 60% av kvinner får behandling. Gitt at vi vet kjønn til observasjon vil ikke informasjon om behandlingsstatus gi oss ny informasjon om forventet potensielle utfall med ulik eksponering for behandling

estimator som er enkel utvalgsanalog.

$$\delta_{ate} = E(y_i^1 - y_i^0) \quad (11.7)$$

$$= E[E(y_i^1 - y_i^0 | x)] \quad (11.8)$$

$$= E[E(y_i^1 | D_i = 1, x) - E(y_i^0 | D_i = 0, x)] \quad (11.9)$$

$$= E[E(y_i | D_i = 1, x) - E(y_i | D_i = 0, x)] \quad (11.10)$$

$$= E[\delta_x] \quad (11.11)$$

$$= \sum_x \delta_x P(X = x) \quad (11.12)$$

Åpner for heterogenitet i behandlingseffekt avhengig av covariates (hvilken undergruppe). Vi er interessert i gjennomsnittlig behandlingseffekt for hele populasjon så gir større vekt til større undergrupper. Kan tilsvarende finne gjennomsnittlig behandlingseffekt for de som blir behandlet ved å i stedet vekte på betinget fordeling i stedet for marginal,

$$\delta_{att} = \sum_x \delta_x P(X = x | D = 1) \quad (11.13)$$

Kan finne utvalgsanalog til forventningene for å evaluere. Matching er greit å gjøre operativt dersom vi har et fåtall veldefinerte undergrupper, men hva hvis covariate er kontinuerlig? Hva hvis inndeling blir for fin slik at mange grupper der vi ikke både oppserverer behandling og kontroll? Kan delvis håndteres ved å minimere avstand. For hvert individ kan vi finne kontroll individ(er) som er mest mulig lik bortsett fra behandlingsstatus og aggregere opp individuelle behandlingseffekter,

$$\hat{\delta}_{ATT} = \frac{1}{N_T} \sum_{i:d_i=1} \left(y_i - \sum_{j \in N(i)} w_{ij} y_j \right) \quad (11.14)$$

der $N(i)$ er indeksene til observasjon i et nabolag til observasjon i og det vektene w_{ij} avhenger av avstand til observasjon. Summerer til én slik at det blir et vektet gjennomsnitt.

11.1.3 Dårlig kontroll

Av og til kan det være fristende å kontrollere for variabler som er bestemt etter behandlingen. Dette er som oftest en dårlig idé siden behandling endrer sammensetning av undergruppene vi ser på slik at forskjellene i utfall ikke kan tilskrives en kausal effekt av behandlingen. Anta at myndighetene innfører et tiltak der et tilfeldig utvalg får gratis personlig trener ($D_i = 1$) og samtidig observerer utfallene til en kontrollgruppe ($D_i = 0$). Vi kan se på gjennomsnittlig effekt av tiltaket på ulike utfall y ved å beregne $E[y_i | D_i = 1] - E[y_i | D_i = 0] = E[y_i^1 - y_i^0]$. Det er ikke nødvendig å betinge for andre variabler, men det kan øke presisjon til estimat og vi kan også bruke stratifisering til å un-

dersøke heterogentitet i behandlingseffekt. Anta nå at vi vil undersøke effekt av tiltaket på undergruppen av observasjoner som trener etter at tiltaket blir iverksatt ($T_i = 1$). Denne beslutningen kan avhenge av D_i så vi kan skrive det opp i potensielt utfall rammeverk,

$$y_i = y_i^0 + D_i(y_i^1 - y_i^0) \quad (11.15)$$

$$T_i = T_i^0 + D_i(T_i^1 - T_i^0) \quad (11.16)$$

Finner differanse i gjennomsnitt i undergruppe,

$$E[y_i | D_i = 1, T_i = 1] - E[y_i | D_i = 0, T_i = 1] \quad (11.17)$$

$$= E[y_i^1 | T_i^1 = 1] - E[y_i^0 | T_i^0 = 1] \quad (11.18)$$

$$= E[y_i^1 | T_i^1 = 1] - E[y_i^0 | T_i^1 = 1] \quad (11.19)$$

$$+ E[y_i^0 | T_i^1 = 1] - E[y_i^0 | T_i^0 = 1] \quad (11.20)$$

der siste linje er seleksjonseffekt. I dette tilfelle vil seleksjonseffekten sannsynligvis være negativ siden gruppen som trener uavhengig av eksponering av tiltak gjerne har bedre utfall enn gruppen som trener på tiltak. Dette skaper seleksjonseffekt selv om behandling i utgangspunktet var tilfeldig fordelt. Analogt så bør man være forsiktig med å legge til utfallsvariabler som kontroll i regresjoner også på observasjonsdata.

Propensity score matching

I stedet for å matche på $\mathbf{x} \in \mathbb{R}^k$ kan vi modellere sannsynlighet for at observasjon mottar behandling, $E[D_i = 1 | \mathbf{x}] := p(\mathbf{x}) \in \mathbb{R}$, og matche på dette. Tror det kan forenkle problemet litt og det er i mange tilfeller enklere å modellere hvordan \mathbf{x} påvirker sannsynlighet for behandling enn utfallet. På en annen side er fremgangsmåten litt mer *non-standard*. Det er ulike valg av vekting, konstruering av feilledd mm. slik at konklusjon kan avhenge av valg til forsker. Noe av fordelene med regresjon er at alt er standardisert!

Regresjon som matching

Hvis vi har modell er saturert i diskret x og antar homogen behandlingseffekt,

$$y = \sum_x d_{xi} \alpha_x + \delta_R D + u_i, \quad (11.21)$$

så kan det vises at

$$\delta_R = \frac{E[\sigma_D^2(X_i) \delta_X]}{E[\sigma_D^2(X_i)]} \quad (11.22)$$

der $\sigma_D^2(X_i) := E[(D_i - E[D_i|X_i])^2|X_i]$, betinget varians av D_i gitt undergruppe X_i . Dette betyr at regresjon - siden det er type effektiv estimator som utnytter informasjon - legger mer vekt på grupper der det er større variasjon i behandling. Med dummy behandling vil den legge mer vekt på grupper der $P(D = 1, X = x) \approx 0.5..$ i stedet for å bare se på størrelsen av gruppen. Har ikke så mye å si dersom behandlingseffekt er omtrent homogen, men hvis regresjon og matching gir veldig ulike resultat kan det være grunn til å tenke litt på hvorfor.

11.1.4 Knytte potensielle utfall til regresjonsligning

For å knytte dette til kausalitet kan vi bruke potensielle utfall der parameter kan korrespondere med (gjennomsnittlig) kausal effekt. Anta først at $y_i^1 - y_i^0 = \delta$.

$$y_i = y_i^0 + \delta D_i \quad (11.23)$$

$$= E[y_i^0] + \delta D_i + y_i^0 - E[y_i^0] \quad (11.24)$$

$$= \alpha + \delta D_i + u_i \quad (11.25)$$

$$(11.26)$$

Ser nå at feilleddet har konkret innhold og hvis vi nå sier at $cov(D_i, u_i) = 0$ så er dette en veldig sterk antagelse som impliserer at $cov(y_i^0, D_i) = 0$. Hvis vi åpner for heterogen behandlingseffekt får vi en tilfeldig koeffisient $\delta_i := y_i^1 - y_i^0$ og ligningen kan skrives som

$$y_i = E[y_i^0] + (E[y_i^1] - E[y_i^0])D_i + y_i^0 - E[y_i^0] + D_i\{(y_i^1 - y_i^0) - (E[y_i^1] - E[y_i^0])\} \quad (11.27)$$

der feilledd kan være korrelert med behandling vis de som har større effekt av behandling i større grad velger å eksponerer seg for den. Da vil lineær regresjon gi skjevt estimat av gjennomsnittlig behandlingseffekt $\delta := E[y_i^1] - E[y_i^0]$. Kan utvide til behandling med flere nivåer. La den nå være gitt ved s der $s \in \{0, 1, \dots, \bar{s}\}$

$$y_i = \begin{cases} y_i^0, & s_i = 0 \\ y_i^1, & s_i = 1 \\ \vdots \\ y_i^{\bar{s}}, & s_i = \bar{s} \end{cases} \quad (11.28)$$

som mer kompakt kan skrives $y_i = f_i(s_i)$. Hvert individ har sin egen kurve med potensielle utfall $f_i(s)$, men vi observerer bare utfallet assosert med den realiserde eksponeringen $s = s_i$. Kan saturere modell og ha parameter for hvert nivå, men i praksis så modellerer vi med lineær funksjon som gir estimat av gjennomsnittlig behandlingseffekt ved å øke eksponering med ett nivå.

11.1.5 Målefeil

11.1.6 Utelatte variabler

Jeg har lyst til å undersøke sammenhengen mellom antall år med skolegang s og inntekt y . Det finnes ingen deterministisk funksjon som forklarer relasjonen siden personer med lik skolegang kan ha ulik lønn av andre grunner. Vi setter derfor opp en modell

$$y = \alpha + \beta s + \epsilon \quad (11.29)$$

der det stokastiske feilleddet ϵ blir et mål på vår uvitenhet. Koeffisientene i ligningen over er ikke entydig bestemt siden vi får enhver $f(\cdot)$ kan definere $\epsilon = y - f(\cdot)$. Generelt så vil vi finne $f(\cdot)$ som minimerer $\|\epsilon\|$ og jeg har vist over at dette er CEF. På en annen side vil vi ha en enkel funksjon med parametre som vi kan tolke. Jeg har derfor avgrenset til å se på lineære funksjoner som er parametrisert med β . Ved å påføre restiksjonen $cov(s, \epsilon) = 0$ korresponderer parameteren i ligningen over med PRF som er beste lineære tilnærming til CEF. Dette kan vi konsistent estimere med OLS og gir oss et greit sammendragsmål på *assosiasjonen* mellom skolegang og inntekt. Gitt at CEF er tilnærmet lineær vil β kunne tolkes som differansen i forventet inntekt mellom to individer med ett års differanse i skolegang. Denne differansen fanger både opp en eventuell kausal effekt av skolegang på inntekt og at individer med ulik skolegang er systematisk forskjellig på måter som påvirker inntekt. Det kan for eksempel tenkes at individ som velger lengre utdanning har høyere evner og motivasjon som vi kan benevne som a . Jeg skal nå undersøke mulighet til å isolere kausal effekt av skolegang og undersøke tolkningen av parametre i regresjonsmodeller.

Vi har formelt definert kausal effekt som differanse i potensielle utfall. For å muliggjøre estimering av kausale effekter fra observasjonsdata der behandling ikke er tilfeldig kan det være lurt å beskrive den kausale prosessen som genererer utfallet. I et laboratorium kan vi i noen sammenhenger beskrive eksakt hvordan et utfall avhenger av egenskaper ved eksperimentet, $y = f(\mathbf{x})$. På grunn av målefeil kan det være små avvik fra sammenhengen. Hvis dette er tilfeldig støy så er $y = f(\mathbf{x}) + \epsilon$ og $E[y|x] = f(\mathbf{x})$ slik at forventningsverdien gjenfanger den deterministiske, kausale sammenhengen. Det reduserer problemet til å estimere CEF. Virkeligheten er langt mer komplisert, men vi kan tenke på det som Guds laboratorium. I utgangspunktet kan vi tenke at fremtidige utfall er deterministisk bestemt av en fullstendig beskrivelse av verdens tilstand på et tidspunkt. La oss tenke på hva som bestemmer personers lønn.

$$y = f(\mathbf{z}) = \delta s + \gamma a + \beta' \mathbf{x} + \epsilon \quad (11.30)$$

Vi antar at den deterministiske $f(\cdot)$ eksisterer, men den kan være vilkårlig komplisert. Det er bare fantasien som setter grenser. Jeg har forenklet den ved å anta at inntekt avhenger

lineært av skole, evne og noen andre variabler \mathbf{x} . På et eller annet tidspunkt må vi nesten slutte å liste opp variabler. Samler sammen bidraget til resten av variablene i et tilfeldig støyled ϵ , der $E[\epsilon|s, a, \mathbf{x}] = 0$. Anta nå at vi observerer (s, a, \mathbf{x}) . Ved å kjøre regresjon

$$y = \alpha + \delta s + \gamma a + \beta' \mathbf{x} + \epsilon \quad (11.31)$$

kan vi gjenfinne hele den kausale sammenhengen. Regresjon er verktøy for å estimere assosiasjon mellom variabler, men i dette tilfellet samsvarer det med den kausale sammenhengen. Anta nå at vi ikke kan observere \mathbf{x} . Vi legger det kumulative bidraget fra disse variablene inn i feilleddet, slik at den kausale sammenhengen nå er

$$y = \delta s + \gamma a + u \quad (11.32)$$

der $u := \beta' \mathbf{x} + \epsilon$. Hvis vi kjører regresjon

$$y = \alpha + \delta s + \gamma a + u \quad (11.33)$$

så vil OLS bestemme parametrene ved å konstruere et feilledd $\hat{u} := y - \hat{\alpha} - \hat{\delta}s - \hat{\gamma}a$ som er ukorrelert med s og a . Dette vil kun samsvare (asymptotisk) med de kausale parametrene dersom feilleddet u fra den kausale prosessen faktisk er ukorrelert med disse variablene.

¹³ Anta nå at $cov(s, u) = 0$ og $cov(a, u) \neq 0$. Kan vi fortsatt finne den kausale effekten av skolegang (δ)? Dette skjer dersom feilleddet konstruert ved OLS er ukorrelert med s . OLS lager et nytt feilledd \hat{u} .. hm... må se litt mer på dette.

$$\hat{u} = \phi a + v, cov(a, v) = 0 \quad (11.34)$$

slik at

$$y = \alpha + \delta s + \phi a + v, \quad cov(s, v) = cov(a, v) = 0 \quad (11.35)$$

Ser at vi fortsatt finner δ , men ikke γ . Generelt kan vi bare håpe å finne en kausal parameter og er ikke noe problem om de andre variablene er korrelert med feilledd.

Hva skjer dersom vi ikke observerer a , men likevel prøver å estimere kausal effekt fra observerte (y, s) ? Vi kan vise at PRF i kort regresjon ikke samsvarer med kausale parameteren.

$$\frac{cov(s, y)}{var(s)} = \frac{cov(s, \alpha + \delta s + \phi a + v)}{var(s)} = \delta + \phi \frac{cov(s, a)}{var(s)} \quad (11.36)$$

Merk at OLS alltid konsistent estimerer PRF, men PRF i den korte regresjonen ikke samsvarer med den kausale parameteren og at OLS derfor er forventningsskjev estimator

¹³Merk at feilleddet u har en konkret tolkning utover å bare være det mekaniske avviket $y - f(\cdot)$.

for den parameteren vi egentlig er interessert i. Når det er selvseleksjon til behandling er det lite troverdig at vi kan obsevere alle *confounding variables*. Vi trenger derfor en kilde til eksogen variasjon i behandling som gjør den ukorrelert med alle disse andre variablene. Dette bringer oss til instrumentelle variabler.

11.2 Instrumentelle variabler

For å kvantifisere effekt av behandling trenger vi i utgangspunktet eksogen variasjon for å unngå seleksjonsskjevhet. Når dette ikke er mulig kan vi forsøke å kontrollere for andre variabler som påvirker utfallet og som er korrelert med behandling.¹⁴ Dette kan motiveres ved at det gjør CIA-antagelsen mer plausibel og det vil redusere bias fra enkle naive sammenligninger av utfall mellom behandling og kontroll, men estimatene vil ofte være lite kredible siden vi aldri kan utelukke at det er flere relevante utelatte variabler. Vi skal nå se på en alternativ fremgangsmåte som utnytter eksogen variasjon i en annen variabel enn selve behandlingen. Denne variabelen kan fungere som et instrument for behandlingen dersom den kun er assosiert med utfallet gjennom å påvirke eksponering for behandling. I den grad det eksisterer en assosiasjon mellom instrument og utfallet må dette skyldes behandlingen, og vi kan kvantifisere behandlingseffekten gjennom å skalere assosiasjonen med assosiasjonen mellom instrument og behandling.

Instrumentelle variabler er derfor en helt annen strategi enn å kontrollere for utelatte variabler. I stedet for å finne variabler som er korrelert med disse, så prøver vi å finne en variabel som er ukorrelert med disse og gjerne egentlig ikke har noe med den kausale prosessen som bestemmer utfallet. Samtidig må det både være korrelert med behandling og ha eksogen variasjon. Det finnes mange kreative instrument i litteraturen, men jeg tror i praksis jeg bare vil bruke det til å analysere eksperiment med delvis compliance samt fuzzy regresjonsdiskontinuitet.

11.2.1 Estimering

Formelt er det to kriterier for at en variabel Z skal være et gyldig instrument for D :

1. Relevans: $cov(D, Z) \neq 0$
2. Eksogenitet: $cov(D, \eta) = 0$

Kan da identifisere den strukturelle parameteren

$$cov(y, Z) = cov(\alpha + \delta D + \eta, Z) = \delta cov(D, Z) \implies \delta = \frac{cov(y, Z)}{cov(D, Z)} \quad (11.37)$$

¹⁴Kan også bruke proxy for slike variabler; variabler som ikke selv inngår i kausal prosess, men som er korrelert med slike variabler og derfor fanger opp effekt fra disse.

Dersom instrumentet er binært kan vi få et enklere uttrykk for estimatoren,

$$E[y_i|Z_i = 1] - E[y_i|Z_i = 0] = \quad (11.38)$$

$$\delta(E[D_i|Z_i = 1] - E[D_i|Z_i = 0]) + E[\eta_i|Z_i = 1] - E[\eta_i|Z_i = 0] \quad (11.39)$$

$$\implies \delta = \frac{E[y_i|Z_i = 1] - E[y_i|Z_i = 0]}{E[D_i|Z_i = 1] - E[D_i|Z_i = 0]} \quad (11.40)$$

gitt at $E[\eta_i|Z_i = 1] - E[\eta_i|Z_i = 0] = 0$ og $E[D_i|Z_i = 1] - E[D_i|Z_i = 0] \neq 0$. Kan også motivere det som en momentestimatorer på matriseform, men synes egentlig at disse matrisegreiene ikke er så relevant for programevaluering siden jeg bare er interessert i én parameter...

$$\mathbb{E}[\mathbf{z}\eta(\beta)] = \mathbb{E}[\mathbf{z}(y - \mathbf{x}'\beta)] = \mathbf{0} \implies \beta = \mathbb{E}[\mathbf{z}\mathbf{z}']^{-1}\mathbb{E}[\mathbf{z}y] \quad (11.41)$$

11.2.2 Heterogen behandlingseffekt

Vi kan tenke at instrumentet setter i gang en kausal kjedereaksjon. Instrumentet påvirker eksponering for en behandling som igjen påvirker utfallet. I utgangspunktet kan utfallet avhenge av både verdi til instrument og behandling slik at det må skrives $y_i(d, z)$. Den første antagelsen vi gjør er at instrumentet kun påvirker utfallet gjennom behandlingen, slik at

$$y_i(d, 1) = y_i(d, 0) \quad (11.42)$$

$$\implies y_i(1, 1) = y_i(1, 0) := y_i^1 \text{ og } y_i(0, 1) = y_i(0, 0) := y_i^0 \quad (11.43)$$

Denne antagelsen er ikke nødvendig for å finne kausal redusert form, altså kausal effekt av Z på y , men er nødvendig for å isolere effekt av behandling D . Vi kan også skrive opp en såkalt *first stage* som er analog til vanlig kausal effekt, men der *behandling* er instrumentet og utfallet er behandling,

$$D_i = D_i^0 + Z_i(D_i^1 - D_i^0) = \pi_0 + \pi_1 Z_i + v_i \quad (11.44)$$

Den første antagelsen er at instrumentet er så godt som tilfeldig fordelt. Det betyr at den observerte eksponeringen ikke sier oss noen ting om de potensielle utfallene.

$$[Y_i^1, Y_i^0, D_i^1, D_i^0] \perp\!\!\!\perp Z_i \quad (11.45)$$

Dette er analog til tilfeldig fordeling av behandling og impliserer at $E[D_i|Z = j] = E[D_i^j]$ og $E[Y_i|Z = j] = E[Y_i^j]$ slik at vi kan lære kausal first stage og redusert form. For at dette skal være oppfylt i praksis må instrumentet ikke være informativt om uobserverte variabler som påvirker utfallet. Med andre ord så må observasjon med ulik eksponering for

instrument ikke være systematisk forskjellige langs andre egenskaper som påvirker utfall. Vi kan utvide denne antagelsen til CIA ved å legge til covariates noen som kan gjøre det mer kredibelt hvis eksponering ikke kan betraktes som et rent eksperiment.

Videre må instrumentet være relevant for behandling, altså at eksponering for instrument endrer behandling for noen av observasjonene

$$E[D_i|Z_i = 1] - E[D_i|Z_i = 0] = E[D_i^1 - D_i^0] = E[\pi_{1i}] \neq 0 \quad (11.46)$$

Til slutt vil vi anta at first stage er monoton slik at $\pi_{1i} \geq 0$ eller $\pi_{1i} \leq 0$ for alle observasjonene. Det medfører at instrument enten øker eksponering for behandling eller reduserer det, men ikke begge deler. Vi kan dele populasjonen inn i fire kategorier ut fra hvordan instrument påvirker deres behandlingsstatus,

1. *Always takers* der $D_i^1 = D_i^0$ og $\pi_{1i} = 0$.
2. *Never takers* der $D_i^1 = D_i^0$ og $\pi_{1i} = 0$.
3. *Compliers* der $D_i^1 = 1$ og $D_i^0 = 0$ og $\pi_{1i} = \dots$ hm.
4. *Defiers* der $D_i^1 = 1$ og $D_i^0 = 0$.

Når vi bruker instrument isolerer vi den variasjonen i eksponering i behandling som er skapt av instrumentet. Siden vi ekskluderer *defiers* per antagelse, så vil all all variasjon skyldes compliers og vi estimerer $E[y_i^1 - y_i^0 | complier]$. Med antagelse om konstant behandlingseffekt kan dette generaliseres som gjennomsnittlig behandlingseffekt for populasjonen, men i praksis kan compliers være systematisk forskjellig fra de andre gruppene. Med heterogen behandlingseffekt isolerer vi bare den lokale gjennomsnittlige behandlingseffekt for compliers. Gitt de fire antagelsene:

1. Eksklusjonskriteriet, eksklusiv kausal kanal gjennom behandling D_i
2. (Betinget) uavhengighetskriteriet, instrument er så godt som tilfeldig fordelt
3. Relevans, kausal first stage
4. Monotoniet

kan jeg vise at wald er LATE. Med andre ord:

$$\frac{E[y_i|Z_i = 1] - E[y_i|Z_i = 0]}{E[D_i|Z_i = 1] - E[D_i|Z_i = 0]} = E[y_i^1 - y_i^0 | D_i^1 > D_i^0] \quad (11.47)$$

For å ta beviset begynne vi med nevneren,

$$E[y_i|Z_i = 1] - E[y_i|Z_i = 0] \quad (11.48)$$

$$= E[y_i^0 + D_i^1(y_i^1 - y_i^0)|Z_i = 1] - E[y_i^0 + D_i^0(y_i^1 - y_i^0)|Z_i = 0] \quad (11.49)$$

$$= E[y_i^0 + D_i^1(y_i^1 - y_i^0)] - E[y_i^0 + D_i^0(y_i^1 - y_i^0)] \quad (11.50)$$

$$= E[(D_i^1 - D_i^0)(y_i^1 - y_i^0)] \quad (11.51)$$

$$= E[y_i^1 - y_i^0|(D_i^1 > D_i^0)P(D_i^1 > D_i^0)] \quad (11.52)$$

Ser at den reduserte formen fanger opp effekt på *compliers* skalert opp med andelen *compliers*. Skal deretter se at *first stage* i nevneren gir andel compliers,

$$E[D_i|Z_i = 1] - E[D_i|Z_i = 0] \quad (11.53)$$

$$= E[D_i^0 + Z_i(D_i^1 - D_i^0)|Z_i = 1] - E[D_i^0 + Z_i(D_i^1 - D_i^0)|Z_i = 0] \quad (11.54)$$

$$= E[D_i^1 - D_i^0] \quad (11.55)$$

$$= E[D_i^1 - D_i^0|D_i^1 > D_i^0]P(D_i^1 > D_i^0) \quad (11.56)$$

$$= P(D_i^1 > D_i^0) \quad (11.57)$$

Vi trenger antagelse om monotonitet fordi ...

Vi trenger ikke antagelsen dersom vi antar konstant behandlingseffekt fordi ...

Praktiske hensyn

Må vurdere antagelse om monotonitet som ikke kan testes. Må også vurdere antagelsen om eksogen variasjon i instrument. Dette kan ikke testes fra data. Det vi derimot kan teste er relevansen til instrumentet. Det gjør vi med hypotesetest i first stage. Vanlig huskeregel er at $F > 10$, men ny studie viser at dette er for lite. Med lite relevant instrument så blir estimatet usikkert. Dette er ikke problem i seg selv, men problem at standardfeil og fordeling bygger på asymptotisk teori og tar veldig lang tid å konvergere med svak instrument slik at rapporterte størrelser er feil i små utvalg slik at inferens blir feil.

Har tester for overidentifikasjon og om IV/OLS konvergerer mot ulike parametre, men er ikke så interessant.

Karakterisere compliant subpopulasjon

Vi finner en lokal behandlingseffekt for delpopulasjonen som faktisk responderer på instrumentet. Heldigvis for oss er dette gjerne effekten som er mest interessant fordi dette er personer som er på marginen i valg om å ta behandling og dermed er mest sensitiv for policy som gjør det enklere tilgjengelig. Vi vil uansett ønske å si noe om andelen *compliers* og hvordan de skiller seg fra resten av populasjonen langs andre observerte egenskaper.

Ettersom vi ikke kan observere både D_i^1 og D_i^0 kan vi aldri observere om gitt enhet er *compliant* eller *always-taker*. Vi kan aldri vite om en person uansett ville tatt behandlingen selv uten eksponering for instrumentet. Likevel kan vi forsøke å beskrive egenskaper til *compliers*. Vi kan finne andel fra first stage og vi kan visstnok også beskrive den betingede fordelingen av andre covariates.

11.2.3 Eksperiment med delvis compliance

I mange eksperiment er det tilfeldig utvalg som blir plassert i gruppen som får tilbud om behandling, men det er ikke alltid mulig å tvinge observasjonene til å eksponere seg for behandling. Dette medfører selv-seleksjon og det er dermed ikke mulig å finne den kausale behandlingseffekten ved å sammenligne utfall til de som blir behandlet og de som ikke. Vi kan finne kausal effekt av å bli *tilbudt* behandlings som betegnes som *Intention to treat*, men dette er ofte ikke like interessant. For å finne kausal effekt kan vi bruke instrumentvariabel der tildeling til behandlingsgruppe er instrument. Dette vil da være helt analogt med den mer generelle diskusjonen over, bortsett fra at vi nå kan ekskludere *always-takers* dersom kun behandlingsgruppen har tilgang på eksponering for behandling. Dette forenkler formelen til

$$\frac{E[y_i|Z_i = 1] - E[y_i|Z_i = 0]}{P[D_i|Z_i = 1]} = E[y_i^1 - y_i^0|D_i = 1] \quad (11.58)$$

Instrumentelle variabler ble først benyttet til å estimere parametre i simultane ligningssystem. Jeg begynner med å beskrive dette fordi mye av terminologien stammer derfra. I praksis brukes det som oftest til å håndtere problemet med utelatte variabler. Det er enklest å motivere det i tilfeldige eksperimenter med delvis *compliance*. Vi kan deretter bruke det til å analysere naturlige eksperimenter.

11.2.4 Simultane ligningssystem

Det klassiske eksempelet på simultant ligningssystem er tilbuds- og etterspørselskurven. Pris og kvantum i marked blir bestemt i samspill av tilbudskurve og etterspørselskurve,

$$Q^s = \alpha_0 + \beta_0 P + \gamma_0 z + u_0 \quad (11.59)$$

$$Q^d = \alpha_1 + \beta_1 P + u_1 \quad (11.60)$$

$$Q = Q^s = Q^d \quad (11.61)$$

der z er observerbar variabel som skifter kurven og vi antar at helning er konstant over tid. Eventuelt kan vi betrakte det som en gjennomsnittlig helning og det blir litt analogt til heterogen behandlingseffekt... Vi sier at tilbuds- og etterspørselskurven er strukturelle

ligninger der parameterne at en kausal tolkning.¹⁵ Konkret så sier det oss endring i henholdsvis tilbudt og etterspurt kvantum dersom vi endrer pris med én enhet og holder alt annet likt. Slike strukturelle sammenhenger er ofte ikke mulig å observere fra tilgjengelig data, så da må de nødvendigvis komme fra teori. Jeg skal nå se på muligheten til å lære (β_0, β_1) fra data. Utfordringen vår er at P er *endogen*.¹⁶

$$Q = \alpha_0 + \beta_0 P + \gamma_0 z + u_0 \quad (11.62)$$

$$Q = \alpha_0 + \beta_0 [(Q - \alpha_1 - u_1)/\beta_1] + \gamma_0 z + u_0 \quad (11.63)$$

$$Q = \frac{1}{1 - \frac{\beta_0}{\beta_1}} \left[\alpha_0 - \frac{\beta_0}{\beta_1} \alpha_1 + \gamma_0 z + u_0 - \frac{\beta_0}{\beta_1} u_1 \right] \quad (11.64)$$

$$Q = \pi_0 + \pi_1 z + v \quad (11.65)$$

... det er intuitivt at vi ikke kan ha eksogen variasjon i P siden det også påvirker etterspørsel, men ser ikke med én gang korrelasjon i feilledd. Merk at siste ligning er såkalt *reduisert form* fordi vi skriver en *endogen* variabel som funksjon av *eksogene* variabler og parametre. De ulike parametrene i redusert form er ikke-lineære funksjoner av underliggende strukturelle parametre og har ikke interessant tolkning i seg selv. Utfordringen er å lære atferdsparameterne fra enkeltligninger og det er her IV kommer inn.

Grafisk så gir det mening av vi kan bruke z til å lære β_1 fordi den skifter tilbudskurven opp og ned. Vil knytte dette til IV estimator. Til slutt kan vi utlede estimatoren med utgangspunkt i simultant ligningssystem som beskrevet over. Vi har en strukturell ligning med en kausal parameter som vi vil estimere

$$y = \alpha + \delta D + \eta \quad (11.66)$$

I første omgang tar jeg ikke med kontrollvariabler, men jeg kan utvide senere. Jeg tenker at motivasjonen for å inkludere disse er litt analogt til kontrollvariabler i vanlig regresjon. For det første kan det øke presisjonen til estimatoren ved å redusere feilleddet. Videre kan det gjør antagelsen om eksogenitet mer kredibel ved at instrumentet er så godt som tilfeldig fordelt innenfor hvert strata (delgruppe). Betingelsen om relevans blir da et spørsmål om *partiell kovarians*; det vil si at det korrelasjon mellom instrument og behandling innenfor strata. Dette kan vi undersøke med helningskoeffisient i multivariat regresjon. Uansett, vi kan ikke lære parameter ved å konstruere residual som er ukorrelert med D i utvalget fordi D er *endogen*. Vi modellerer eksponering for behandling i såkalt *first stage*,

$$D = \pi_{00} + \pi_{01} z + v_0 \quad (11.67)$$

¹⁵Vi kan også si at de inneholder såkalte *atferdsparametre* som sier noe om endring i atferd dersom vi endrer egenskap ved systemet

¹⁶Litt usikker på om endogen har en presis definisjon. Kan tenke at det er korrelert med feilledd, men endogenitet kan jo også henspille på at variabelen blir bestemt innenfor systemet...

som ikke er en strukturligning. Dette beskriver bare deskriptiv sammenheng i data og det er derfor slik at $cov(z, v_0) = 0$ per konstruksjon. Det er et verktøy for å finne kausal parameter. Vi finner redusert form ved å plugge first stage inn i strukturligningen,

$$y = \alpha + \delta[p_{i00} + \pi_{01}z + v_0] + \eta \quad (11.68)$$

$$y = \alpha + \delta p_{i00} + \delta \pi_{01}z + \eta + \delta v_0 \quad (11.69)$$

$$y = \pi_{10} + \pi_{11}z + v_1 \quad (11.70)$$

der $\pi_{11} := \delta \pi_{01}$. Merk at $cov(z, v_0) = 0$ per konstruksjon, men trenger også at $cov(z, \eta) = 0$ for at $cov(z, v_1) = 0$ slik at vi kan lære π_{11} .¹⁷ Kan nå se to nye strategier for å lære δ ,

1. Kjøre first stage og redusert form separat, $\delta = \pi_{11}/\pi_{01}$
2. Plugge $\hat{D} = p_{i00} + \pi_{01}z$ inn for D i strukturligning og kjøre den.

11.2.5 Generalisering av wald

Jeg kan utvide dette til å se på flere instrument der vi finner den lineære kombinasjonen som er mest mulig korrelert med D . Dette finner vi uansett i first stage, så ikke noe problem å legge til flere. Det er ikke nødvendigvis så interessant siden hvert instrument estimerer en instrument-spesifikk lokal behandlingseffekt, det vil si behandlingseffekt fra delpopulasjon som complier til det gitte instrument. Når vi kombinerer får vi en saus; et vektet gjennomsnitt. Jeg starter heller enkelt med å se på wald-estimatorer som er spesialtilfelle av IV der både instrument og behandling er dummyvariabler.

$$\delta = \frac{E[Y|Z = 1] - E[Y|Z = 0]}{E[Y|D = 1] - E[Y|D = 0]} \quad (11.71)$$

11.3 Regresjonsdiskontinuitet

Regresjonsdiskontinuitet utnytter at regler som bestemmer eksponering for behandling har vilkårlige cutoffs og at det er begrenset mulighet for observasjonene til tilpasse seg eksakt hvilken side den havner på.¹⁸ Mer formelt er sannsynlighet for å bli eksponert for behandling, er diskontinuerlig i en verdi s_0 av en såkalt *running variable* s . Dette kan for eksempel være en reform som ble innført på en gitt dato s , institusjonelle regler som gjør at enheter som ikke oppnår en gitt verdi av s må innføre tiltak eller at man trenger en viss score s for å få en gevinst. I skarp rdd er behandlingen da for eksempel $D_i = I\{s_i \geq s_0\}$. Hvis det er tilfeldig hvilken side av cutoff s_0 observasjonene havner vil ikke dette si oss noe om potensielle utfall i fravær av behandling slik at differanse i utfall

¹⁷vet også at jeg trenger $\pi_{01} \neq 0$ men klarer ikke se hvordan det kommer inn her.

¹⁸Det finnes vist to tilnærminger. Continuity based og basert på lokal randomisering. Vet ikke helt hva som er forskjell i praksis

identifiserer kausal effekt. Dersom det ikke er mulig for enhetene å eksakt bestemme sin verdi av s er det kredibelt at det er tilfeldig for et intervall rundt s_0 .

Det er en utfordring at forventet utfall uten behandling avhenger av s . Det kan være både fordi det eksisterer en direkte sammenheng mellom s og utfallet y , men også fordi andre variabler som påvirker utfall er korrelert med s . Vi modellerer derfor sammenhengen $\mathbb{E}[y|s]$ og bruker eventuell diskontinuitet i s_0 som estimat på kausal effekt. Det er også mulig å knytte til potensielle utfall. Anta at $\mathbb{E}[y_i^0|s_i] = f(s_i)$ og at det er konstant effekt slik at $y_i^1 = y_i^0 + \rho$. Da er

$$\mathbb{E}[y_i|s_i] = \begin{cases} \mathbb{E}[y_i^0|s_i] = f(s_i), & s_i < s_0 \\ \mathbb{E}[y_i^1|s_i] = f(s_i) + \rho, & s_i \geq s_0 \end{cases} \quad (11.72)$$

$$= f(s_i) + \rho D_i \quad (11.73)$$

der $D_i := I\{s_i \geq s_0\}$. Dette impliserer at

$$y_i = \mathbb{E}[y_i|s_i] + (y_i - \mathbb{E}[y_i|s_i]) \quad (11.74)$$

$$= f(s_i) + \rho D_i + u_i \quad (11.75)$$

Vi kan spesifisere en parametrisk form på $f(\cdot)$ som kan estimeres med OLS, for eksempel en polynom av orden p . Da blir ligningen

$$y_i = \alpha + \beta_1 s_{1i} + \dots + \beta_p s_i^p + \rho D_i + u_i \quad (11.76)$$

Det er også mulig å generalisere til ulike $f(s)$ på hver side av s_0 . Dette kan gjøres parametrisk ved å ha ulike parametre i polynom og interkasjon med indikator for om det er over s_0 . Blir litt sånn som splines. Ellers kan vi også bruke ikke-parametrisk lokal regresjon. I alle tilfeller er behandlingseffekten identifisert av

$$\rho = \lim_{s \rightarrow s_0^+} E[y|s] - \lim_{s \rightarrow s_0^-} E[y|s] \quad (11.77)$$

Det er nødvendig å spesifisere et intervall rundt cutoff, $[s_{min}, s_{max}]$, som avgrenser hvilke observasjoner vi betrakter.¹⁹ Lengden på intervallet omtales som *bandwidth*. Et lengre bandwidth medfører flere observasjoner som kan mer presise estimat, men medfører også at observasjonene blir mer ulike. Jo lenger vekk fra s_0 , jo mer sannsynlig at observasjoner er ulike langs andre dimensjoner. Dette gjelder spesielt hvis observasjonene kan tilpasse seg for oppnå payoff hvis $s > s_0$.

Det er viktig å treffe med funksjonell form på $f(\cdot)$ for å få riktig estimat på diskontinuitet. Hvis den modelleres som med en lineær likning vil ikke-linearitet kunne bli fanget

¹⁹Det finnes data-driven måter å gjøre dette på slik at vi slipper å velge det ad-hoc.

opp som en diskontinuitet. Det er vanlig å bruke polynom, splines eller ikke-parametrisk lokal regresjon som er vektet med en kernel. Som alltid er det ikke én metode som dominerer; valg av struktur avhenger av antall observasjoner og hvor fleksibel funksjonen må være for å fange funksjonell form. Det er en klassisk bias-varians tradeoff.

For at strategien skal identifisere kausal effekt kan det ikke være slik at andre ting endres brått i s_0 . For eksempel kan det være flere reformer som innføres samtidig, flere ting som endres når en person blir pensjonist. Da vil vi ikke diskontinuiteten identifisere akkurat den endringen vi er interessert i. Det må også være tilfeldig hvilken side observasjonen havner på. Hvis det er strategisk selv-selektering ved at observasjon bevisst velger $s \leq s_0$ kan de være systematisk forskjellig langs uobserverte variabler.

Strategi for å underbygge kredibilitet til identifiserende antagelse:

1. Density test: Vil at sannsynlighetstetthet til s skal være kontinuerlig rundt s_0 . Hvis de klumper seg sammen på éne siden gir det indikasjon på strategisk selvseleksjon.
2. Covariate balance: Vil at andre variabler ikke skal endres brått i s_0 . Kan ikke observere alt, men kan kjøre opplegget med andre covariates vi observerer som utfall og se om det er diskontinuitet. Alternativt kan vi bare se om de har like gjennomsnitt.. men føler at vi kan håndtere at det varierer med s ..
3. Placebo-tester: Vil ikke se hopp der vi ikke forventer hopp. Kan kjøre opplegget med andre verdier av s som cut-off.

11.3.1 Fuzzy rdd

Diskontinuerlig hopp i sannsynlighet for eksponering for behandling i $s = s_0$,

$$P(D = 1|x) = \begin{cases} g_0(x), & x < x_0 \\ g_1(x), & x \geq x_0 \end{cases} \quad (11.78)$$

$$= g_0(x) + I\{x \geq x_0\}[g_1(x) - g_0(x)] \quad (11.79)$$

bruker det som first stage men er litt usikker på praktisk gjennomføring..

11.4 Paneldata

Paneldata er datastruktur der vi har flere observasjoner fra ulike grupper eller gjentatte observasjon av individ.²⁰ Det har gode egenskaper i programevaluering siden vi kan se

²⁰Det er litt glidende overgangs fra krysseksjonsdata siden vi ikke trenger tidsdimensjon. Poeng at vi har to (eller flere) indekser ut fra gruppetilhørighet, men kan kjøre dummy regresjons i stedet... Det finnes stor litteratur om multi-level modellering og såkalte random effects som jeg føler at ikke er så relevant for programevaluering dersom behandling ikke er tilfeldig fordelt.

på forskjeller i forskjeller til behandling og kontrollgruppe i stedet for absolutte nivåer.²¹ Dette gjør at vi kan dekomponere feilledd og se bort fra den komponenten som er konstant for gruppen. Det gjør det mer kredibelt at behandling er ukorrelert med den gjenstående komponenten av feilleddet og kan redusere bias, men det er fortsatt en sterk antagelse. Vi kan jo ikke kvantifisere effekt av en gitt behandling dersom den er korrelert med bruk av annen behandling som vi ikke observerer.

Det er også en utfordring at det kan være lite variasjon i behandling innad i gruppen slik at estimatene blir usikre. Det medfører også at de kan være mer sensitive for målefeil en kryssseksjonsdata..

11.4.1 Dekomponering av feilledd

Anta at marginal kausal effekt av behandling x på utfall y er konstant $\partial y / \partial x = \beta$, som verken avhenger av mengden behandling eller andre egenskaper til enheten. Dette medfører at

$$y = \beta x + (y - \beta x) = \beta x + u \quad (11.80)$$

der $u := (y - \beta x)$ representerer det kumulative bidraget til utfallet av alle andre variabler. Jeg vil at feilleddet skal ha forventning lik null, så jeg sentrerer variabelen og omdefinerer,

$$y = \beta x + \mathbb{E}[u] + (u - \mathbb{E}[u]) = \beta x + \alpha + u_c \quad (11.81)$$

der jeg heretter lar $u_c := u$. La $d(i)$ være dummy som indikerer om observasjonen tilhører gruppe i . Vi kan da dekomponere feilleddet,

$$u = \mathbb{E}[u|d(i)] + (u - \mathbb{E}[u|d(i)]) = \alpha_i + \epsilon \quad (11.82)$$

slik at ligningen kan skrives som

$$y_{it} = \alpha + \beta x_{it} + \alpha_i + \epsilon_{it} \quad (11.83)$$

der α_i er en parameter som fanger opp uobservert heterogenitet. Det er bidraget av uobserverte variabler som er felles innad i gruppe. Tror jeg kan tenke på det som konkrete variabler eller bare dekomponering og projektering av tilfeldig variabel... Betegnes som fast-effekt fordi den er konstant innad i gruppe, f.eks. alle egenskaper innad i familie som ikke varierer (foreldres utdanningsnivå?). Jeg kan estimere β konsistent med OLS dersom $cov(x_{it}, \epsilon_{it}) = 0$. Denne antagelsen er mer kredibel nå som jeg fått den fasteffekten ut av feilleddet, men det er fortsatt en ganske sterk antagelse på observasjonsdata. Hvis det er uobserverte ting som påvirker utfallet og samvarierer med behandlingen så kan vi ikke

²¹Ekvivalent kan vi isolere behandlingseffekt ved å se på variasjon innad i gruppe..

isolere effekten av denne. Tenke for eksempel om vi kjører to behandlingen samtidig og kun observerer eksponering for ene. Tror vi kan bruke instrument på samme måte som i kryssseksjon ved å sette det opp på first-difference form, men vet lite om dette.

11.4.2 Identifikasjon

Skal forsøke å formalisere dette med rammeverket for potensielle utfall. Anta heretter at i er individ og t er tid. Begynner med å anta at eksponering for behandling er så godt som tilfeldig fordelt betinget av egenskaper ved individet som er konstant over tid A_i , tidspunkt t og noen andre observerte covariates \mathbf{x}_{it} som kan variere over tid for hvert individ,

$$y_{it}^0 \perp\!\!\!\perp D_{it} | \mathbf{x}_{it}, A_i, t \quad (11.84)$$

$$\implies E[y_{it}^0 | \mathbf{x}_{it}, A_i, t, D_{it}] = E[y_{it}^0 | \mathbf{x}_{it}, A_i, t] \quad (11.85)$$

Antar en lineær en lineær modell,

$$E[y_{it}^0 | \mathbf{x}_{it}, A_i, t] = \alpha + \delta_t + A_i' \gamma + \mathbf{x}_{it}' \beta \quad (11.86)$$

Tror vi kun kan estimere en additativ kausal effekt,

$$E[y_{it}^1 | \mathbf{x}_{it}, A_i, t] = E[y_{it}^0 | \mathbf{x}_{it}, A_i, t] + \delta \quad (11.87)$$

$$\implies E[y_{it} | \mathbf{x}_{it}, A_i, t, D_{it}] = \alpha + \delta_t + A_i' \gamma + \mathbf{x}_{it}' \beta + \delta D_{it} \quad (11.88)$$

$$\implies y_{it} = \alpha_i + \delta_t + \mathbf{x}_{it}' \beta + \delta D_{it} + \epsilon_{it} \quad (11.89)$$

der $\alpha_i := \alpha + A_i' \gamma$ og $\epsilon_{it} := y_{it}^0 - E[y_{it}^0 | \mathbf{x}_{it}, A_i, t]$.²²

11.4.3 Estimering

Si litt high level om dette. I praksis knyttes til estimering av ligningsystem.. staker matriser og sånn, blir heavy notasjon i lineær algebra. Kan også si noe om forskjell på fixed og random effects.

Fixed effects (within)

Vi vil se på variasjon innad i gruppe. En naturlig måte er å demean alle variabler innad i hver gruppe. Finner først gjennomsnitt i hver gruppe, der jeg åpner for at antallet

²²Vet ikke hvorfor det ikke er $\epsilon_{it} := y_{it} - E[y_{it} | \cdot]$ Må se på dette senere

medlemmer kan variere.²³

$$\frac{1}{T(i)} \sum_t y_{it} = \frac{1}{T(i)} \sum_t [\alpha + \beta x_{it} + \alpha_i + \epsilon_{it}] \quad (11.90)$$

$$\bar{y}_i = \alpha + \beta \bar{x}_i + \alpha_i + \bar{\epsilon}_i \quad (11.91)$$

Deretter trekker jeg dette fra hver observasjon

$$\ddot{y}_{it} := y_{it} - \bar{y}_i = \beta \ddot{x}_{it} + \ddot{\epsilon}_{it} \quad (11.92)$$

slik at vi kun trenger at $cov(x_{it}, \epsilon_{it}) = 0$.²⁴ Alternativt kan jeg bruke dummies til eksplisitt å estimere fast-effektene. Dette illustrerer litt koblingen mellom panel og dummies generelt. Ved å introdusere dummies så definerer man eksplisitt undergrupper og identifiserer effekt av behandling ved å se på variasjon innad i gruppene. Således kan jo regresjon med kjønnsdummy betraktes som panel der hvert kjønn utgjør gruppe.

Pooled OLS

Random effects

Betrakte α_i som tilfeldig variabel. For at det estimator skal være konsistent må vi anta at $cov(x_{it}, \alpha_i) = 0$ slik at det egentlig ikke er noe stort poeng i å bruke panel struktur, bortsett fra at vi kan bruke strukturen til å få mer effektiv estimator enn vanlig OLS. Denne strukturen utnyttes gjennom feasible generalized least square (FGLS). Tror ikke det er veldig stort poeng, bortsett fra at det ikke er kosher.

Det kan derimot være relevant å modellere på denne måten dersom man ikke er interessert i kausal parameter, men kun endring i forvetningsverdi... bruker da MLE til å estimere.

Kan bruke såkalt hausmannstest for å vurdere om $cov(x_{it}, \alpha_i) = 0$. Hvis den er oppfylt vil både F.E og R.E estimatorene være konsistent, men hvis den ikke er oppfylt vil de konvergere mot ulike størrelser. Huasmann gir oss test som sier om differansen er stor nok til at vi kan forkaste nullhypotesen om at $cov(x_{it}, \alpha_i) = 0$.

Between estimator

En mulig strategi er å kollapse variablene innad i hver gruppe slik at vi bare har gjennomsnittene. Mister variasjon innad i gruppe (tidsdimensjon hvis gjentatt observasjon av individ) og ser bare på forskjeller mellom grupper. Kan identifisere kausaleffekt hvis gjennomsnittlig behandling ikke er korrelert med uobservert heterogenitet α_i , men det vil

²³Såkalt ubalansert panel. Ikke stort problem dersom ikke skyldes systematisk skjevt frafall som er brudd på forutsetning om tilfeldig utvalg.

²⁴Tror det er ekvivalent med $cov(\ddot{x}_{it}, \ddot{\epsilon}_{it}) = 0$. kunne kanskje vist.

være en lite effektiv estimator siden vi mister mye informasjon,

$$\bar{y}_i = \bar{\mathbf{x}}_i\beta + \bar{u}_i \quad (11.93)$$

GLS som vektet gjennomsnitt av within og between

hm. kan ta resultat og intuisjon, men utledning er ganske fucked. kunne brukt det som anledning til å si noe om GLS generelt, men gjør det et annet sted. Angrist sier at GLS er teit uansett, men kan gi litt intuisjon om 2SLS som GLS på wald estimatorer... eller noe sånt.

11.4.4 Dynamisk panel

Kan ønske å utnytte tidsseriedimensjonen i panel, altså det faktum at vi observerer samme individer på flere tidspunkt, til å modellere dynamikk. For det første kan det være slik at en behandling x_{it} påvirker utfall ikke bare i tidspunkt t men også i fremtidige perioder $t+1, t+2, \dots$. Dette kan vi ta hensyn til ved å inkludere laggede uavhengige variabler. Vi kan også ha lyst til å ta med lagged avhengig variabel. Tror at motivasjonen for dette er at vi vil sammenligne likt med likt og derfor må ta med historien til individene. Selv om de ser like ut på et gitt tidspunkt t , så kan informasjon om utfallene deres på tidligere tidspunkt gi informasjon om hva slags type person de er...

11.5 Dynamiske modeller

11.5.1 Lagged uavhengig variabel

Så langt har jeg sett på univariate tidsserier. Dette er ofte greiest dersom jeg vil bruke modellen til forecasting. Men jeg kan jo også være interessert i relasjon mellom variabler. Tenk for eksempel at jeg vil modellere hvor sulten jeg er og at jeg for hver time t rapporterer sult-nivå (y_t) samt egenskaper ved ting jeg har gjort i den aktuelle perioden \mathbf{x}_t . La for eksempel x_t være antall kalorier jeg har spist. Dette vil ikke bare påvirke sultnivå i t , men også i fremtidige perioder $t+1, t+2, \dots$ med gradvis mindre effekt. Modellen kan beskrives med

$$y_t = \alpha_0 + \beta_0 x_t + \beta_1 + x_{t-1} + \beta_2 x_{t-2} + \epsilon_t \quad (11.94)$$

Anta at jeg spiser 100 kalorier i $t = 1$ og ingen i de andre. Funksjonene blir da

$$y_1 = \alpha_0 + \beta_0 \cdot 100 + \beta_1 \cdot 0 + \beta_2 \cdot 0 + \epsilon_1 \quad (11.95)$$

$$y_2 = \alpha_0 + \beta_0 \cdot 0 + \beta_1 \cdot 100 + \beta_2 \cdot 0 + \epsilon_2 \quad (11.96)$$

$$y_3 = \alpha_0 + \beta_0 \cdot 0 + \beta_1 \cdot 0 + \beta_2 \cdot 100 + \epsilon_3 \quad (11.97)$$

$$(11.98)$$

Hvis effekten er avtagende så vil $\beta_0 > \beta_1 > \beta_2$ slik at $E[y_1|(x_1, x_2, x_3) = (100, 0, 0)] > E[y_2|(x_1, x_2, x_3) = (100, 0, 0)]$ osv. Gir sånn passe mening dette her.

11.5.2 Lagged avhengig variabel

Vi kan motivere dette med evaluering av arbeidsmarkedstiltak. Personer på tiltak har ofte hatt tap av inntekt. Hvis de har høyere lønnsvekst etter tiltak så kan dette fange opp at de har høyere *potensiell lønn* enn sammenlignbare personer i kontrollgruppe? Vet ikke helt, men tar litt utledning uansett. Enkleste modell er

$$y_{it} = \gamma y_{i,t-1} + u_{it}, \quad \text{der } u_{it} = \alpha_i + \epsilon_{it} \quad (11.99)$$

der $(\epsilon_{it})_{t \in \mathbb{N}}$ er hvit støy (altså ingen seriekorrelasjon). Skal nå se på mulighet til å estimere den kausale γ . Det forutsetter at den uavhengige variabelen ikke er endogen.²⁵ Jeg har et enkelt oppsett for å vurdere konsistens til OLS-estimator i enkel univariat regresjon,

$$\hat{\beta} = \frac{\sum x_n y_n}{\sum x_n^2} = \frac{\sum x_n (\beta x_n + u_n)}{\sum x_n^2} = \beta + \frac{\sum x_n u_n}{\sum x_n^2} \quad (11.100)$$

Vi kan da vurdere konsistens ved å skalere med $1/N$ og slik at teller konvergerer til $cov(x_n, u_n)$ når $n \rightarrow \infty$. Mange estimatorer kan betraktes som OLS på transformerte variabler, så her er det bare å plugge inn. Prøver først med fixed effect der $x_n := \sum_{t=1} (y_{i,t-1} - \bar{y}_{i,-1})$ og $y_n = \sum_{t=1} (y_{i,t} - \bar{y})$

$$\hat{\beta} = \frac{\sum_n \sum_{t=1} (y_{i,t-1} - \bar{y}_{i,-1})(y_{i,t} - \bar{y})}{\sum_n \sum_{t=1} (y_{i,t-1} - \bar{y}_{i,-1})^2} \quad (11.101)$$

$$= \gamma + \frac{\sum_n \sum_{t=1} (y_{i,t-1} - \bar{y}_{i,-1})(u_{i,t} - \bar{u}_i)}{\sum_n \sum_{t=1} (y_{i,t-1} - \bar{y}_{i,-1})^2} \quad (11.102)$$

$$(11.103)$$

Det kan sikkert vises at det stemmer, men jeg bare plugger inn de analoge størrelsene fra transformert OLS. Litt usikker på hvordan jeg skal vise hva som er problemet her.

²⁵Jeg er ikke 100% komfortabel med begrepet endogen. Brukes vel bare short-hand for å være korrelert med feilledd, men tror det er greit å se litt på ligningssystem for at terminologi skal gi litt mening.

Tror problemet er at $y_{i,t-1}$ er positivt korrelert med α_i i feilleddet. Personer som har (uobserverte) konstante egenskaper som gir de høy lønn i forrige periode vil i gjennomsnitt ha høyere feilledd i denne perioden ikke sant.. så hvis høyt utfall i forrige periode er behandling", så vil behandlingseffekten fange opp dette og vi får skjevt estimat av kausal effekt.

Løsningen på dette er å bruke momentbetingelser implisert av modellen. For at et instrument z skal være gyldig må

1. Relevans, må forklare noe variasjon i behandling: $cov(y_{i,t-1}, z) \neq 0$
2. Eksogenitet, må ikke være korrelert med uobserverte variabler som påvirker utfall: $cov(u_{i,t}, z) = 0$.

Siden prosess er autoregressiv vil utfall i tidligere periode propagere gjennom prosessen. Laggede avhengige variabler er relevant og ikke korrelert med idiosynkratiske delen av feilleddet i hvertfall. Vil være korrelert med uobservert heterogenitet, men dette kan vi bli kvitt ved å ta first difference. Les om Arellano-Bond hvis dette blir relevant i fremtiden...

11.5.3 Instrument

Selv i med fixed effect må vi anta at $cov(x_{it}, \epsilon_{it}) = 0$, altså at behandling er ukorrelert med idiosynkratisk (?) uobserverte variabler. Kan forsøke å legge inn forklaringsvariabler slik at komponent av ϵ_{it} som ikke er forklart av kontroll-variablene ikke er korrelert med behandling, men i likhet med i krysseksjon kan de fortsatt være relevante utelatte variabler som varierer over tid. En alternativ strategi er å bruke instrument til å identifisere behandlingseffekt med utgangspunkt i variasjonen av behandling som kan forklares med instrumentet. Tror jeg trenger litt recap på instrument i krysseksjon først.

11.6 Forskjeller i forskjeller

Fixed effect er vel og bra hvis vi har variasjon i eksponering for behandling innad i grupper og føles oss komfortabel med å si at behandling er omtrent betinget uavhengig av potensielle utfall. I praksis er ofte behandling bestemt overfra og ned i betydningen at noen grupper blir eksponert for tiltak/reform og andre ikke. Hvis vi kun har krysseksjon med utfallet til gruppene etter behandling er det vanskelig å trekke noen konklusjoner siden gruppene nok var forskjellige i utgangspunktet. Hvis vi derimot har informasjon om gruppene over tid kan vi se på forskjell i *endring* for hver gruppe i stedet for nivå. Dette gjør det mulig å identifisere behandlingseffekten under antagelse om at de ville hatt felles trend i fravær av behandling. Dette skal jeg nå vise formelt med to grupper, to tidsperioder og to behandlingsnivåer. Deretter skal jeg utvide til flere tidsperioder, flere behandlingsnivåer og kontroll for individuelle covariates.

11.6.1 Identifikasjon

Potensielt utfall i periode t :

$$Y_{it} = Y_{it}^0 + D_i(Y_{it}^1 - Y_{it}^0) \quad (11.104)$$

Vi vil estimere average treatment effect on treated (ATT):

$$\delta_1 = E[Y_{i1}^1 - Y_{i1}^0 | D_i = 1] \quad (11.105)$$

Identifiserende antagelse er felles trend i fravær av behandling:

$$E[Y_{i1}^0 | D_i = 1] - E[Y_{i0}^0 | D_i = 1] \quad (11.106)$$

$$= E[Y_{i1}^0 | D_i = 0] - E[Y_{i0}^0 | D_i = 0] \quad (11.107)$$

Antar altså at $(Y_{i1}^0 - Y_{i0}^0) \perp\!\!\!\perp D_i$. Bevis for identifikasjon:

$$\delta_1 = E[Y_{i1}^1 | D_i = 1] - E[Y_{i1}^0 | D_i = 1] \quad (11.108)$$

$$= E[Y_{i1} | D_i = 1] - E[Y_{i1}^0 | D_i = 1] \quad (11.109)$$

Vi observerer $E[Y_{i1} | D_i = 1]$, men $E[Y_{i1}^0 | D_i = 1]$ er kontrafaktisk. Kan finne proxy gitt identifiserende antagelse:

$$E[Y_{i1}^0 | D_i = 1] = E[Y_{i0}^0 | D_i = 1] + (E[Y_{i1}^0 | D_i = 1] - E[Y_{i0}^0 | D_i = 1]) \quad (11.110)$$

$$= E[Y_{i0} | D_i = 1] + (E[Y_{i1}^0 | D_i = 0] - E[Y_{i0}^0 | D_i = 0]) \quad (11.111)$$

$$= E[Y_{i0} | D_i = 1] + (E[Y_{i1} | D_i = 0] - E[Y_{i0} | D_i = 0]) \quad (11.112)$$

Alle størrelsene er observerbare og den kausale effekten er identifisert:

$$\delta_1 = E[Y_{i1} | D_i = 1] - E[Y_{i0} | D_i = 1] - (E[Y_{i1} | D_i = 0] - E[Y_{i0} | D_i = 0]) \quad (11.113)$$

Merk at antagelsen er ekvivalent med stabil seleksjonsskjevhets:

$$E[Y_{i1}^0 | D_i = 1] - E[Y_{i1}^0 | D_i = 0] \quad (11.114)$$

$$= E[Y_{i0}^0 | D_i = 1] - E[Y_{i0}^0 | D_i = 0] \quad (11.115)$$

11.6.2 Flere grupper og flere tidsperioder

I det enkleste eksempelet trenger vi egentlig bare fire størrelser: gjennomsnitt i hver av gruppene før og etter eksponeringen for behandling. Dersom vi har individdata er det flere

fordeler ved å sette opp en regresjonsmodell på form

$$y_{ist} = \gamma_s + \delta_t + \delta D_{st} + \epsilon_{ist} \quad (11.116)$$

11.7 Limited Dependent Variable

Skal nå begynne å se på utfallsvariabler med begrensede utfallsrom. Det kan være fordi utfallet er kategorisk slik at utfallsverdi bare er kode for den realiserte kategorien, fordi utfallet bare kan ta heltallsverdier, kun kan være positivt eller er avgrenset på andre måter. Felles for disse modellene er at vi estimerer dem med maximum likelihood. Et annet fellestrekk er at vi modellerer en underliggende latent variabel med ubegrenset utfallsrom som gjerne har en deterministisk komponent som avhenger lineært av inputvariabler samt et feilledd som vi gjør litt antagelser om. Vi definerer en regel som knytter den latente variabelen til det observerte utfallet og bruker dette til å lære om parametre i den underliggende latente modellen. Vi bruker deretter den latente modellen til å svare på spørsmål vi er interessert i.

Den latente variabelen kan i noen tilfeller ha en fysisk tolkning der den kun er uobserverbar fordi vi ikke ser realiseringen i datasettet vårt. I andre tilfeller er det en ren konstruksjon som brukes til å modellere datagenereringsprosessen. Vi begynner med sistnevnte når vi modellerer klassifikasjon gjennom stokastisk nytte..

11.7.1 Stokastisk nytte

Vi antar at nytte ved valg av gode 1 (i forhold til gode 0) har en systematisk komponent som avhenger av observerte egenskaper ved valgsituasjonen \mathbf{x}_i . Kan anta at den er felles for ulike individer eller så på det som et slags gjennomsnitt. Dessuten er det en tilfeldig komponent som er avviket av nytten til gitt individ fra gjennomsnitt til andre individ med samme observerte egenskap. Kan ha ulike preferanser og være ulike på andre måter.

$$u_i = v(\mathbf{x}_i) + \epsilon_i \quad (11.117)$$

$$y_i = I\{u_i > 0\} \quad (11.118)$$

Der vi setter terskelverdi lik 0 uten tap av generalisering siden nytte ikke har noe skala. Vi modellerer systematisk komponent med $v(\mathbf{x}_i) = \mathbf{x}_i' \beta$. Vi kan lære om parameter i denne funksjonen ved å observere at

$$P[y = 1|\mathbf{x}] = P[\mathbf{x}'\beta + \epsilon > 0|\mathbf{x}] \quad (11.119)$$

$$= \mathbb{P}[\epsilon > -\mathbf{x}'\beta|\mathbf{x}] \quad (11.120)$$

$$= F(\mathbf{x}'\beta) \quad (11.121)$$

der $F(\cdot)$ er kumulativ fordeling til $-\epsilon$, som vi antar er symmetrisk slik at det også er cdf til ϵ . Det er tre vanlige valg av $F(\cdot)$ som gir probit, logit og LPM. Kan definere LPM slik at det blir gyldig cdf, men i praksis bruker jeg bare helningen og ikke predikert verdi.

Flere kategorier

Med flere kategorier blir det

$$u_{ij} = v_{ij} + \epsilon_{ij} \quad (11.122)$$

$$y_i = \arg \max_{j \in \{1, \dots, J\}} u_{ij} \quad (11.123)$$

dersom vi antar at feilleddet er fra type 1 ekstremverdifordeling kan vi vise at

$$P(Y = k | v_{ij}) = \frac{\exp(v_{ik})}{\sum_j \exp(v_{ij})} \quad (11.124)$$

Det følger fra ligningen at relativ sannsynlighet til to alternativer ikke blir påvirket av endring i egenskap til tredje alternativ. Dette er en egenskap, og det er ikke nødvendigvis ønskelig.. Har i hovedsak to typer modeller:

1. Conditional logit der egenskaper varierer mellom ulike kategori (og for ulike individ): $v_{ij} = x_{ij}\beta$. Poeng at individer har preferanse for ulike egenskaper ved alternativ (for eksempel pris) og vi beregner vekt for disse ulike egenskapene alternativene kan ha. Vi kan da bruke vektene til å analysere sannsynlighet for at individ vil bruke et nytt tilbud $J + 1$ dersom vi kjenner x_{iJ+1} .
2. Multinomial logit der vi kun bruker egenskap til individ $v_{ij} = x_{ij}\beta_j$. Vi finner i hvilken grad de ulike alternativene vekt egenskapene, for eksempel om kvinner foretrekker gitt alternativ. Denne analysen kan ikke brukes på nye alternativ eksempel siden vi ikke kjenner deres vekter.

kan også kombineres. Det er noe greier om at man kan estimere priseffekt/krysspriseffekt og analysere verdsetting av tid. Relevant for transportøkonomi, men jeg vet ingenting.

Ordnet logit

Hvis alternativene kan rangeres kan vi modellerer de fra underliggende latent variabel, bare at vi finner en partisjonering slik at $y = j$ hvis $y^*(\gamma_{j-1}, \gamma_j]$. Tror grenseverdiene blir parametere. Bestemmer at $\gamma_1 = 1$ for de skal være entydig siden $\mathbf{x}\beta$ ikke har noen naturlig skala.

Merk at det er kun for kategoriene i endene der fortegn på koeffisient har entydig effekt på sannsynlighet for at observasjon tilhører kategori.

11.7.2 Sensurert regresjon (tobit)

En annen situasjon - som egentlig er ganske forskjellig! - som analyseres på tilsvarende måte er valg med hjørneløsning. I mange situasjoner er det ikke mulig å velge negativ kvantum slik at det blir en opphopning av verdier i $y = 0$. Jeg tenker at vi da kunne ha modellert sannsynlighet for positivt kvantum med probit og modellert $E[y|x, y > 0]$ med OLS hver for seg. Med Tobit gjør vi begge deler samtidig ved å anta at begge avhenger latent variabel.

Vi begynner som alltid ved å beskrive modell for latent variabel og hvordan den er relatert til observert utfall:

$$y^* = \mathbf{x}'\beta_0 + u, \quad \text{der } u|\mathbf{x} \sim N(0, \sigma_0^2) \quad (11.125)$$

$$y = \max\{y^*, 0\} \quad (11.126)$$

Estimering

Jeg vil beskrive hvordan betinget sannsynlighet til y gitt \mathbf{x} avhenger av parametre. Vil har et uttrykk for tetthet til sannsynlighet av y for en gitt \mathbf{x} ikke sant, men dette er en mixed fordeling som punktsannsynlighet i $y = 0$. For å håndtere dette deler jeg opp den betingede fordelingen og behandler hvert tilfelle separat.

$$f(y|\mathbf{x}) = \begin{cases} 0, & y < 0 \\ P(y^* \leq 0|\mathbf{x}), & y = 0 \\ f(y^*|\mathbf{x}), & y > 0 \end{cases} \quad (11.127)$$

Jeg kan nå bruke antagelsene fra den latente modellen til å gi et eksplisitt uttrykk for hver størrelse, noe som samtidig gir meg likelihood-funksjonen når jeg betrakter det som funksjon av parameter slik at jeg kan estimere parametre fra observerte data.

$$P(y = 0|\mathbf{x}) = P(y^* \leq 0|\mathbf{x}) \quad (11.128)$$

$$= P(\mathbf{x}'\beta_0 + u \leq 0|\mathbf{x}) \quad (11.129)$$

$$= P(u < -\mathbf{x}'\beta_0|\mathbf{x}) \quad (11.130)$$

$$= P\left(\frac{u}{\sigma_0} < -\frac{\mathbf{x}'\beta_0}{\sigma_0}|\mathbf{x}\right) \quad (11.131)$$

$$= \Phi\left(-\frac{\mathbf{x}'\beta_0}{\sigma_0}\right) \quad (11.132)$$

$$= 1 - \Phi\left(\frac{\mathbf{x}'\beta_0}{\sigma_0}\right) \quad (11.133)$$

For de positive verdiene kan vi utlede på samme måte som i vanlig regresjon. Begynner med å ta sannsynlighet for intervall siden punktsannsynlighet til tetthet=0..

$$P(y < y|\mathbf{x}, y > 0) = P(y^* < y|\mathbf{x}, y > 0) \quad (11.134)$$

$$= P(\mathbf{x}'\beta_0 + u < y|\mathbf{x}, y > 0) \quad (11.135)$$

$$= P\left(\frac{u}{\sigma_0} < \frac{y - \mathbf{x}'\beta_0}{\sigma_0}|\mathbf{x}, y > 0\right) \quad (11.136)$$

$$= \Phi\left(\frac{y - \mathbf{x}'\beta_0}{\sigma_0}\right) \quad (11.137)$$

Deriverer for å finne uttrykk for tetthet

$$f(y|\mathbf{x}, y > 0) = \frac{1}{\sigma_0} \phi\left(\frac{y - \mathbf{x}'\beta}{\sigma_0}\right) \quad (11.138)$$

Trenger ikke gjøre det relativt til standardisert fordeling.. skal se om jeg kan omskrive dette senere. Uansett, det gir en representasjon av loglikelihoodfunksjon,

$$\log L(\beta, \sigma) = \sum_n I\{y_n = 0\} \log\left(1 - \Phi\left(\frac{\mathbf{x}'_n \beta_0}{\sigma_0}\right)\right) \quad (11.139)$$

$$+ I\{y_n > 0\} \log\left(\frac{1}{\sigma} \phi\left(\frac{y - \mathbf{x}'_n \beta}{\sigma}\right)\right) \quad (11.140)$$

Når vi har estimert parametrene kan vi gi estimerte mål på størrelser vi er interessert i og se hvordan de er relatert til parameter β fra underliggende latent modell.

Tolke koeffisient

Når vi har fått estimert β så ligner output litt på vanlig regresjon og det kan være fristende å tolke det på vanlig måte. Men β er parameter i $\mathbb{E}[y^*|\mathbf{x}] = \mathbf{x}'\beta$. Tror dette kan ha direkte tolkning hvis data er sensuert, men ikke dersom vi modellerer hjørneløsning. Da vil vi istedet bruke latent modell til å beskrive sentraltendens i det faktiske utfallet.

$$\mathbb{E}[y|\mathbf{x}] = \mathbb{E}[y|\mathbf{x}, y > 0]P(y > 0|\mathbf{x}) + 0 \quad (11.141)$$

$$= \mathbb{E}[y|\mathbf{x}, y > 0]\Phi\left(\frac{\mathbf{x}'\beta_0}{\sigma_0}\right) \quad (11.142)$$

der jeg har brukt at $P(y > 0|\mathbf{x}) = 1 - P(y = 0|\mathbf{x})$. Det er ganske rimelig at vi kan evaluere forventingsverdi ved å ta vektet sannsynlighet av betingede forventninger på partisjone-ring, men litt usikker på hvordan jeg viser det formelt. Se lov om iterated expectations.

Vil finne et uttrykk for betinget sannsynlighet av utfall gitt at det er positivt.

$$\mathbb{E}[y|\mathbf{x}, y > 0] = \mathbb{E}[y^*|\mathbf{x}, y > 0] \quad (11.143)$$

$$= \mathbb{E}[\mathbf{x}'\beta_0 + u|\mathbf{x}, y > 0] \quad (11.144)$$

$$= \mathbf{x}'\beta_0 + \sigma_0 \mathbb{E}\left[\frac{u}{\sigma_0} \mid \mathbf{x}'\beta_0 + u > 0\right] \quad (11.145)$$

$$= \mathbf{x}'\beta_0 + \sigma_0 \mathbb{E}\left[\frac{u}{\sigma_0} \mid \frac{u}{\sigma_0} < -\frac{\mathbf{x}'\beta_0}{\sigma_0}\right] \quad (11.146)$$

$$= \mathbf{x}'\beta_0 + \sigma_0 \frac{\phi(-c)}{1 - \Phi(-c)} \quad (11.147)$$

$$= \mathbf{x}'\beta_0 + \sigma_0 \frac{\phi(c)}{\Phi(c)} \quad (11.148)$$

$$= \mathbf{x}'\beta_0 + \sigma_0 \lambda(c) \quad (11.149)$$

der $c := \frac{\mathbf{x}'\beta_0}{\sigma_0}$ og $\lambda(c) := \frac{\phi(c)}{\Phi(c)}$ betegnes som den inverse mills ratioen. Dette gir oss et uttrykk for $\mathbb{E}[y|\mathbf{x}, y > 0]$ som er en størrelse vi kunne forsøkt å estimere ved å avgrense til kun observasjoner med positivt utfall. Kan se at det består både av helning til forventningsverdi av latentverdi og et ekstra ledd. Litt usikker på hvordan vi tolker, men det har noe sammenheng med at observasjoner med positivt utfall har større verdi av uobservert variabel u i latent modell. Dersom vi er interessert i parameter β_0 fra latent modell vil vi derfor få skjevt estimat dersom vi avgrenser til observasjon med positivt utfall.. Uansett, kan nå også plugge inn størrelsen og få et uttrykk

$$\mathbb{E}[y|\mathbf{x}] = \mathbb{E}[y|\mathbf{x}, y > 0]P(y > 0|x) \quad (11.150)$$

$$= [\mathbf{x}'\beta_0 + \sigma_0 \lambda(c)] \Phi(c) \quad (11.151)$$

$$= \mathbf{x}'\beta_0 \Phi(c) + \sigma_0 \phi(c) \quad (11.152)$$

der $c := \frac{\mathbf{x}'\beta_0}{\sigma_0}$. Når jeg har estimert parametrene er estimatene av $\mathbb{E}[y|\mathbf{x}]$ og $\mathbb{E}[y|\mathbf{x}, y > 0]$ bare to deterministiske funksjoner av \mathbf{x} . Disse funksjonene er en forenklet representasjon av egenskaper ved virkeligheten og er dessuten upresise fordi utvalget gir begrenset informasjon om den *sanne* prosessen²⁶ som genererer data. Uansett, det er ihvertfall objekter jeg kan jobbe med og bruke til å svare på spørsmål. Vi begynner med å se på hvordan $\mathbb{E}[y|\mathbf{x}, y > 0]$ endres når vi endrer en variabel x_j .

$$\frac{\partial}{\partial x_j} \mathbb{E}[y|\mathbf{x}, y > 0] = \frac{\partial}{\partial x_j} (\mathbf{x}'\beta_0 + \sigma_0 \lambda(c)) \quad (11.153)$$

$$= \beta_j \sigma \frac{\partial}{\partial c} \lambda(c) \frac{\beta_j}{\sigma} \quad (11.154)$$

²⁶Eller vår representasjon av den..

Må finne $\frac{\partial}{\partial c}\lambda(c)$. Bruker at $\partial\phi(c)/\partial c = -c\phi(c)$.²⁷

$$\frac{\partial}{\partial c} \frac{\phi(c)}{\Phi(c)} = \frac{-\mathbf{x}'\beta\phi(c)\Phi(c) - \phi(c)^2}{\Phi(c)^2} \quad (11.155)$$

$$= \frac{-\phi(c)[c\Phi(c) + \phi(c)]}{\Phi(c)^2} \quad (11.156)$$

$$= -\lambda(c) \frac{c\Phi(c) + \phi(c)}{\Phi(c)} \quad (11.157)$$

$$= -\lambda(c)[c + \phi(c)] \quad (11.158)$$

slik at

$$\frac{\partial}{\partial x_j} \mathbb{E}[y|\mathbf{x}, y > 0] = \beta_j \{1 - \lambda(c)[c + \phi(c)]\} \quad (11.159)$$

Siden $\mathbb{E}[y|\mathbf{x}] = \mathbb{E}[y|\mathbf{x}, y > 0]P(y > 0|x)$ og er $P(y > 0|x) = \Phi(c)$, kan vi bruke produktregel og plugge inn

$$\frac{\partial}{\partial x_j} \mathbb{E}[y|\mathbf{x}] = \beta_j \{1 - \lambda(c)[c + \phi(c)]\} \Phi(c) + [c + \sigma\lambda(c)] \frac{\beta_j}{\sigma} \phi(c) \quad (11.160)$$

$$= \beta_j \{\Phi(c) - \phi(c)[c + \phi(c)]\} + [c + \sigma\lambda(c)] \frac{\beta_j}{\sigma} \phi(c) \quad (11.161)$$

$$= \beta_j \Phi(c) - \beta_j \phi(c)c - \beta_j \phi(c)^2 + c \frac{\beta_j}{\sigma} \phi(c) + \lambda(c) \beta_j \phi(c) \quad (11.162)$$

$$(11.163)$$

wtf. Bruker istedet

$$\mathbb{E}[y|\mathbf{x}] = \mathbf{x}'\beta\Phi(c) + \sigma\phi(c) \quad (11.164)$$

$$\frac{\partial}{\partial x_j} \mathbb{E}[y|\mathbf{x}] = \beta_j \Phi(c) + \mathbf{x}'\beta\phi(c) \frac{\beta_j}{\sigma} - c\sigma\phi(c) \frac{\beta_j}{\sigma} \quad (11.165)$$

$$= \beta_j \Phi(c) + c\phi(c)\beta_j - c\phi(c)\beta_j \quad (11.166)$$

$$= \beta_j \Phi(c) \quad (11.167)$$

der jeg igjen har brukt $\partial/\partial c\phi(c) = -c\phi(c)$.

Partial effect avhenger av hvor vi evaluerer \mathbf{x} . Kan bruke average partial effect (APE) på samme måte som i probit. Merk også at derivert gir en lokal lineær tilnærming. Kan få eksakt endring ved én enhets endring ved å plugge verdier inn i $g(\cdot) := \mathbb{E}[y|\cdot]$...

²⁷hvorfor?

Kausalitet i Tobit

Så langt har jeg sett på Tobit som fremgangsmåte for å modellere $E[y|x]$ og $E[y|x, y > 0]$ for utfall som er begrenset til å være positiv gjennom en lineær latent modell. Har sett at $E[y|x]$ består av to komponenter: sannsynlighet for å ha positivt utfall og verdi gitt positiv. Hvis vi vil undersøke kausal effekt av binær tilfeldig behandling D_i på et slikt begrenset utfall kan det være fristende å bruke Tobit for å beregne de ulike effektene.²⁸ Hvis vi plotter histogram av betinget fordeling av utfall for behandling og kontroll grupper så vil de ha tyngdepunkt i 0 og en eller annen fordeling for positive verdier. Vi kan føle at forskjellen i gjennomsnittet ikke fanger alle aspekter ved behandlingseffekten, og det er for såvidt sant. Vi kan bruke Tobit til å dekomponere behandlingseffekt,

$$E[y_i|D_i = 1] - E[y_i|D_i = 0] \quad (11.168)$$

$$= E[y_i|D_i = 1, y_i > 0]P[y_i > 0|D_i = 1] \quad (11.169)$$

$$- E[y_i|D_i = 0, y_i > 0]P[y_i > 0|D_i = 0] \quad (11.170)$$

$$= E[y_i|D_i = 1, y_i > 0]P[y_i > 0|D_i = 1] - E[y_i|D_i = 1]P[y_i > 0|D_i = 0] \quad (11.171)$$

$$+ E[y_i|D_i = 1]P[y_i > 0|D_i = 0] - E[y_i|D_i = 0, y_i > 0]P[y_i > 0|D_i = 0] \quad (11.172)$$

$$= \{P[y_i > 0|D_i = 1] - P[y_i > 0|D_i = 0]\}E[y_i|D_i = 1, y_i > 0] \quad (11.173)$$

$$+ \{E[y_i|D_i = 1, y_i > 0] - E[y_i|D_i = 0, y_i > 0]\}P[y_i > 0|D_i = 0] \quad (11.174)$$

Hmm... litt usikker på om jeg skal beholde det. Poeng at *conditional on positive* ($y > 0$) er å betinge på utfall av behandling. Gruppene med positivt utfall i behandling og kontroll er systematisk forskjellig på andre måter slik at forskjellene ikke fanger kausal effekt. Noe å tenke på dersom man bruker Tobit til å analysere kausal effekt...

11.7.3 Heltallsverdier (Poisson-regresjon)

I mange sammenhenger kan utfallsvariabelen kun ta ikke-negative heltallsverdier. Et eksempel er antall barn en kvinne føder eller antall dager før person havner tilbake i fengsel. Hvis vi bruker MLE må vi modellere hele den betingede fordelingen, men vi kan også avgrense oss til kun å modellere sentraltendensen $\mathbb{E}[y|\mathbf{x}] = g(\mathbf{x}, \beta)$. Siden y ikke kan ta negative verdier gir det lite mening om denne funksjonen gjør det. Et mulig valg er $g(\mathbf{x}, \beta) = \exp\{\mathbf{x}'\beta\}$. Denne funksjonen er ikke lineær i parametre, så vi kan ikke bruke vanlig closed form OLS. Vi kan derimot bruke ikke-lineær OLS som er enkel generalisering og løse det numerisk.

I praksis kan det være greit å påføre mer struktur ved å modellere den betingede

²⁸Merk at binær behandling så er CEF, $E[y|D]$ nødvendigvis lineær slik at det er god grunn til å bruke vanlig regresjon. Mer generelt er de gjerne ikke-lineære med LDV slik at det blir større motivasjon for å beregne form med ikke-lineær kurve... Det vil uansett gjerne være slik at det ikke er så stor forskjell i gjennomsnittlig marginal effekt.

fordelingen. Hvis vi velger en parametrisk fordeling så kan vi estimere den betingede sannsynligheten for de ulike utfallene i stedet for bare å ha sentraltendens og mål på spredning. Et vanlig valg er poisson-fordeling.

$$p(y; \lambda) = \frac{e^{-\lambda}}{y!} \lambda^y \quad (11.175)$$

$$p(y|\mathbf{x}; \beta) = \exp(-\exp(\mathbf{x}'\beta)) \exp(\mathbf{x}'\beta)^y / y! \quad (11.176)$$

$$\text{Log}L_n(\beta) = y_n \mathbf{x}'\beta - \exp(\mathbf{x}'\beta) \quad (11.177)$$

Poisson-fordelingen har egenskapen at $\mathbb{E}[y|\mathbf{x}] = \mathbb{V}[y|\mathbf{x}] = \lambda(\mathbf{x}, \beta) = \exp(\mathbf{x}'\beta)$. I praksis kan dette stemme dårlig med den gitte fordelingen vi observerer og dette bør vi ta hensyn til. MLE-estimatoren kan gi konsistent estimat på $\mathbb{E}[y|\mathbf{x}]$ i en større klasse av fordelinger, men standardfeilen som vi utleder fra informasjonsmatrisen vil være feil. Kan innføre ny parameter σ der $\mathbb{V}[y|\mathbf{x}] = \sigma \lambda(\mathbf{x}, \beta)$. Bruker som vanlig en sandwich til å skalere, men i dette tilfelle er estimator en skalar, slik at

$$\widehat{se}_{QMLE} = \hat{\sigma} \widehat{se}_{MLE} \quad (11.178)$$

må finne ut av dette senere.

11.8 Modellere seleksjon

Så langt har vi antatt at vi har et tilfeldig utvalg slik at det er representativt for populasjonen. Av ulike grunner så kan det være slik at vi ikke observerer alle egenskaper vi er interessert i for alle enheter. Noen kan la være å svare på alle spørsmål, noen kan la vær å svare overhodet, noen kan falle fra i forskningsopplegget slik at vi ikke ser utfallet. Selv med observasjonsdata kan det være systematiske skjevheter i hvilke variabler vi observerer for ulike personer. Man må for eksempel ha jobb for at vi skal observere lønn. Med slike ufullstendige utvalg kan ikke nødvendigvis konklusjon fra utvalg generaliseres til populasjon.

Jeg skal nå utvikle et rammeverk for å betrakte ufullstendige utvalg. Jeg vil se på egenskapene til estimatoren med ulike former for seleksjon og se på mulighet til å korrigere for eventuell skjevhet.

Poenget er at vi kan konstruere en binær tilfeldig variabel $s_n := I\{\text{observerer hele}(\mathbf{x}_n, y_n)\}$. Det er en tilfeldig variabel med betinget fordeling. Vi kan modellere denne og forsøke å undersøke hva som påvirker om vi observerer.

Vi kan anta at prosessen $(\mathbf{x}_n, y_n)_{n \in \mathbb{N}}$ oppfyller alle gode egenskaper, men istedet for tilfeldig utvalg $\{(\mathbf{x}_n, y_n) : n = 1, \dots, N\}$ observerer vi $\{s_n(\mathbf{x}_n, y_n) : n = 1, \dots, N\}$. Setter alle verdier lik null dersom noen mangler.

11.8.1 Avkortet (truncated) regression

I sensurert regression kan vi observere et tilfeldig utvalg av enheter i populasjon og deres karakteristikk, med unntak av utfall kan være sensuert i endepunktene. I avkortet regresjon blir observasjoner med gitte verdier av utfall ekskludert fra utvalget. Det er ikke lenger et tilfeldig (representativt) utvalg. Det er ingenting i veien for å estimere størrelser betinget av eksklusjonskriteriet, men dette kan avvike fra egenskaper ved underliggende latent modell som kan være det vi egentlig er interessert i. Begynner som alltid med å beskrive latent modell og relasjon mellom latente utfall og observerte utfall.

$$y^* = \mathbf{x}'\beta_0 + u, \quad u|\mathbf{x} \sim N(0, \sigma_0^2) \quad (11.179)$$

$$y = \begin{cases} y^*, & y^* > 0 \\ (y, x) \text{ ikke observert,} & y^* \leq 0 \end{cases} \quad (11.180)$$

11.8.2 Heckit

Vi kan tenke på seleksjon til utvalg som to-steg prosess. Først blir individer trukket og deretter blir en tilfeldig variabler s_n realisert og bestemmer om vi observerer realiseringen av resten av variablene.²⁹ Hvis s_n er helt tilfeldig, så fører ikke frafallet til at utvalget er systematisk forskjellig fra populasjonen og vi kan se bort i fra det. Hvis det derimot er korrelert med sammenhengen vi vil estimere så kan det føre til skjevheter og konklusjoner om utvalget kan ikke lenger generaliseres til populasjonen. Det er et poeng at vi alltid kan estimere størrelser betinget av seleksjon, men dette er ikke alltid det vi er interessert i. En mulig løsning er å modellere seleksjonsprosessen og forsøke å korrigere for seleksjon.

I heckit modellerer vi to latente variabler

$$y^* = \mathbf{x}'\beta + \epsilon, \quad \epsilon := y^* - E[y^*|\mathbf{x}] \quad (11.181)$$

$$s^* = \mathbf{z}'\gamma + v \quad (11.182)$$

og der de korresponderende utfallene vi observerer i data er

$$s = I\{s^* > 0\} \quad (11.183)$$

$$y = \begin{cases} y^*, & s = 1 \\ \text{uobservert,} & s = 0 \end{cases} \quad (11.184)$$

²⁹Vi kan være interessert i latent utfall. På samme måte som over kan latent utfall i prinsippet være fysisk observerbar, mens i det i andre tilfeller er litt mer konstruert størrelse.

Den betingede forventningsfunksjonen betinget av seleksjon er

$$E[y|\mathbf{x}, s = 1] = E[y^*|\mathbf{x}, s^* > 0] \quad (11.185)$$

$$= \mathbf{x}'\beta + E[\epsilon|v > -\mathbf{z}'\gamma] \quad (11.186)$$

Hvis jeg estimerer med OLS vil jeg asymptotisk finne beste lineære tilnærming til denne funksjonen. Det er skjevhet som er litt analog til utelatt variabel. Kan isolere β ved å modellere det andre leddet og ta det med i estimeringen. For å gjøre dette antar jeg at feilleddene i de latente ligningene er multivariat normalfordelt.

$$(\epsilon, v) \sim N(\mathbf{0}, \Sigma) \quad (11.187)$$

De vil ha positiv korrelasjon dersom individer med uobservert egenskap som gjør de over gjennomsnittlig stor sannsynlighet for å være selektert (høy v) også gir de over gjennomsnittlig høyt utfall (høy ϵ). Uansett, denne antagelsen medfører at

$$1. E[v|v > -c] = \frac{\phi(c)}{\Phi(c)} := \lambda(c)$$

$$2. E[\epsilon|v] = \sigma_{\epsilon,v}E[v]$$

slik at

$$E[\epsilon|v > -\mathbf{z}'\gamma] = \sigma_{\epsilon,v}\lambda(\mathbf{z}'\gamma) \quad (11.188)$$

Heckmans' to-steg estimator er da

1. Finne $\hat{\gamma}$ fra estimering av $P(s = 1|\mathbf{z}) = \Phi(\mathbf{z}'\beta)$. Bruk det til å lage $\hat{\lambda}_n = \lambda(\mathbf{z}'_n \hat{\gamma})$
2. Kjør OLS på $y_n = \mathbf{x}'_n\beta + \theta\hat{\lambda}_n + \text{feilledd}$, der $\hat{\theta}$ blir estimat på kovarians. Tror at hvis den er null så er det ikke sammenheng i seleksjon når vi kontrollerer for \mathbf{x} (og \mathbf{z} ?). I så fall får vi lignende estimat hvis vi kjører vanlig OLS, men greit å ha sjekket uansett..

Til slutt kan jeg nevne at poeng å ha identifiserende variabel som påvirker seleksjon men ikke utfall.

Eksempler

Observerer kun lønn dersom den er over reservasjonslønn,

$$w = \alpha_0 + \beta ed + \epsilon \quad (11.189)$$

$$w^r = \alpha_1 + \delta ed \rho m + u \quad (11.190)$$

$$s^* = w^r - w = \alpha_1 - \alpha_0 + (\delta - \beta)ed - \psi m + (u - \epsilon) \quad (11.191)$$

$$s^* = \gamma_0 + \gamma_1 ed + \gamma_2 m + v \quad (11.192)$$

hmm.

Kapittel 12

Kalkulus

Jeg trenger litt greier for optimering. Ta recap på basic theorem, derivasjon og integral, optimering i \mathbb{R} og \mathbb{R}^2 (både betinget og ubetinget), generalisering til flere dimensjoner, konkavitet, gradient, jacobi, hesse, lagrange, koblinger til lineær algebra (kvadratisk form, (lokalt) lineære transformasjoner,...)

12.1 Litt bakgrunn

12.1.1 Algebra

Det mulig å uttrykke påstander der sannhet avhenger av verdi til variabel x . En vanlig problemstilling med praktisk anvendelse er å finne sannhetsmengde til påstand. Delmengder av tallinjen er intervall eller kombinasjon av intervall (snitt/union). Avstand mellom to tall er absoluttverdi. Eksempler på påstander og deres sannhetsmengder:

$$|x| = D \iff x \in \{-D, D\} \iff x = -D \vee x = D \quad (12.1)$$

$$|x - a| = D \iff x = a - D \vee x = a + D \quad (12.2)$$

$$(2 - x)^{-1} < 3 \iff x \in (-\infty, \frac{5}{3}) \cap (-2, \infty) \quad (12.3)$$

Det tredje eksempel var mest sammensatt. Kan generelt løse ved å omskrive på form $\frac{a \cdot b}{c \cdot d} < 0$, finne ut for hvilke x -verdier hver av faktorene er negative og dermed finne sannhetsmengde til hele uttrykket. En annen vanlig problemstilling er å finne røttene til en funksjon, altså verdiene av x der $f(x) = 0$. En grunn til at dette dukker opp så ofte er at vi alltid kan omskrive $g(x) = h(x) \iff g(x) - h(x) = 0 \iff f(x) = 0$, der $f := g - h$.¹ En mulig fremgangsmåte her er å faktorisere $f(x) = a \cdot b \cdot c \cdot \dots$ og finne for hvilke x -verdier hver av faktorene er lik null.

¹Jeg tror addisjon og subtraksjon av funksjoner er veldefinert, men er ikke helt sikker.

12.2 Litt analyse

Det sies at kalkulus er studie av endring. Vi mapper tall, så det er jo litt interessant å se hvor mye verdien av output endres når vi gjør en liten endring i input. Historisk har folk brukt konseptet om *infinitesimal* endringer, men i siste halvdel av 1800-tallet kom teorien på litt tryggere grunn med følger og grenser (som man kan lære mer om i reell analyse).

12.2.1 Følger

Vi kan betrakte en følge $(x_1, x_2, \dots, x_n, \dots) := (x(n) : n \in \mathbb{N}) = (x_n)_{n \in \mathbb{N}}$ som en uendelig tuple der hvert positive heltall blir mappet til et objekt $x(n) := x_n$. Dette er en ganske generell datastruktur, men i kalkulus betrakter vi tilfelle der $x_n \in \mathbb{R}^N$. Vi er interessert i om følger konvergerer til et tall $\mathbf{r} \in \mathbb{R}^N$. Vi sier at $\lim x_n = \mathbf{r}$ eller $x_n \rightarrow \mathbf{r}$ dersom det er slik at uansett hvor lite vi gjør et nabolag om \mathbf{r} , så vil vi kunne finne en index N slik at alle elementene i følgen med høyere indeks befinner seg i nabolaget. Vi beskriver nabolaget med en såkalt ϵ -ball om \mathbf{r} ,

$$B_\epsilon(\mathbf{r}) := \{\mathbf{x} : d(\mathbf{x}, \mathbf{r}) < \epsilon\} \quad (12.4)$$

der $d(\cdot)$ er en norm, for eksempel den euklediske: $d(\mathbf{x}, \mathbf{r}) = \sqrt{(\mathbf{x} - \mathbf{r})'(\mathbf{x} - \mathbf{r})} := \|\mathbf{x} - \mathbf{r}\|$.

12.2.2 Rekker

En rekke er summen av leddene i en følge. Kan være entent endelig eller uendelig. Vi er interessert i om en uendelig rekke konvergerer. Tror vi kan betrakte summen av de n første leddene som element i en følge, $x(n) = \sum_{i=1}^n y_i$, og undersøke om $(x_n)_{n \in \mathbb{N}}$ konvergerer.

12.2.3 Grenser og kontinuitet

Bruker epsilon-delta til å definere grense.

$$\lim_{x \rightarrow a} f(x) = f(a) \iff \exists \delta |x - a| < \delta \implies |f(x) - f(a)| < \epsilon, \quad \forall \epsilon > 0 \quad (12.5)$$

Hvis grensen til f eksister i a så sier vi at funksjonen er kontinuerlig i a . Hvis funksjonen er kontinuerlig for alle $a \in I$ så er den kontinuerlig på intervallet I .² Det er også venstre- og høyrekontinuitet som jeg er litt usikker på hvordan man definerer formelt.

$$\lim_{x \rightarrow a} f(x) = f(a) \iff \lim_{x \rightarrow a^-} f(x) = f(a) \wedge \lim_{x \rightarrow a^+} f(x) = f(a) \quad (12.6)$$

²Sikkert noe tekniske greier om endepunktene.

12.2.4 Topologi

Vil inføre noen definisjoner til å beskrive mengden en funksjon er definert på.

Åpne mengder

En mengde er *åpen* dersom den ikke inkluderer grenseverdiene. Mer formelt er en mengde S åpen hvis det for hver $x \in S$ eksisterer en $\epsilon > 0$ slik at $B_\epsilon(x) \subset S$. Funksjoner på åpne mengder har ikke nødvendigvis noen ekstremverdier så de kan være vanskelig å optimere. Dette motiverer også bruk av inf og sup

Lukket mengde

En mengde er lukket dersom den inneholder grenseverdi til alle konvergente følger av elementer i mengden, $x_n \rightarrow x \implies x \in S$ dersom $x_n \in S \forall n \in \mathbb{N}$.

Begrenset mengde

En mengde er begrenset dersom det eksisterer et tall B slik at $\|x\| < B$ for alle $x \in S$.

Kompakt mengde

En mengde er kompakt dersom den er både lukket og begrenset.

12.3 Litt om funksjoner...

12.3.1 Real valued functions

En funksjon f er en real valuedfunksjon hvis den mapper til tallinjen; $f : A \rightarrow \mathbb{R}$. Vi er ofte interessert i å finne for hvilke element i A at funksjonen tar sin høyeste verdi.

- Maximizer av f på A er $\{a^* \in A | f(a^*) \geq f(a), \forall a \in A\}$
- Maksimumsverdi er da $f(a^*)$

Det er en utfordring at maxmizers ikke alltid eksisterer. Dette gjelder for eksempel for monotont voksende funksjoner på åpne intervall. For å håndtere dette kan vi definerer en supremum $s := \sup A$ der (1) $a \leq s, \forall a \in A$ og (2) det eksisterer en følge $(x_n), x_n \in A$, slik at $x_n \rightarrow s$

12.3.2 Inverse funksjoner

Hvis $f'(x) > 0$ for alle $x \in I$ er funksjonen strengt monotont voksende på intervallet slik at $b > a \implies f(b) > f(a)$. Da eksisterer det en funksjon f^{-1} med definisjonsmengde

$\{f(x) : x \in I\}$, verdimengde I og der $f^{-1}(f(x)) = x$. Det er en ny funksjon som tar output og mapper til input under f . Funksjonen f må være strengt monoton fordi ellers vil flere input mappe til samme output og dermed vil ikke den inverse tilfredstille definisjonen av en funksjon og er dermed ikke definert.

12.4 Lineær tilnærming av funksjoner

Funksjoner kan gjøres vilkårlig kompliserte og de kan potensielt være vanskelige å beskrive og manipulere. Mye av kalkulus handler om å finne en enklere (eg. lineær) representasjon av funksjonen som gir tilnærming av funksjonen i et omegn om $x^* \in S$.

12.4.1 Derivasjon

Hvis funksjonen er kontinuerlig så kan vi tegne punktene $\{(x, f(x)) : x \in I\}$ uten å løfte blyanten. Dermed kan vi tenke oss å lage en sekant som er en rett linje mellom to punkter på grafen $(x, f(x)), (x+h, f(x+h))$ og se hva som skjer med helningen når h går mot 0. Dette er den deriverte til funksjonen i x .

$$f'(x) := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (12.7)$$

for at den deriverte skal være definert må grenseverdien eksistere. Da er kontinuitet en nødvendig, men ikke tilstrekkelig betingelse. $f'(x)$ er også en funksjon av x og denne må være kontinuerlig i a for at $f'(a)$ skal være definert. Uansett, fra definisjonen over kan man utlede derivasjonsregler for alle (?) funksjonstyper og kombinasjoner av disse, så det er forholdsvis enkelt å derivere. Lite utvalg:

- Kjernerregel: $\frac{d}{dx} f(g(x)) = \frac{d}{dx} f(u) = \frac{df}{du} \frac{du}{dx}$
- Generalisering av kjernerregel:
- Noe om inverse funksjoner?

Når jeg jeg deriverer en funksjon $f : \mathbb{R} \rightarrow \mathbb{R}$ får jeg ut en ny funksjon $f' : \mathbb{R} \rightarrow \mathbb{R}$. Hvis jeg evaluerer den i et element av inputmengden får jeg $f'(x^*) = a \in \mathbb{R}$; et tall som angir grenseverdien av helningen til f i x^* . Jeg kan bruke dette til å beskrive funksjonen som gir verdier av tangentlinjen

$$g : s \mapsto f(x^*) + f'(x^*)(s - x^*) \approx f(x^* + s) \quad (12.8)$$

der jeg har bruke s istedet for Δx . Dette er en affine funksjon, men vi kan se på lineær representasjon ved å se direkte på endring av funksjonsverdi i stedet for nivå.

$$h : s \mapsto f'(x^*)(s - x^*) \approx f(x^* + s) - f(x^*) \quad (12.9)$$

Generelt kan vi skrive $h(s) = DF(x^*)s$. Hvis funksjonen $f : \mathbb{R}^N \rightarrow \mathbb{R}$ får vi lineær tilnærming

$$h : (s_1, \dots, s_N) \mapsto \sum \frac{\partial}{\partial x_n} F(\mathbf{x}^*) s_n = DF(\mathbf{x}^*) \mathbf{s} \quad (12.10)$$

der $DF(\mathbf{x}^*) := \left[\frac{\partial}{\partial x_1} F(x^*), \dots, \frac{\partial}{\partial x_N} F(x^*) \right]$ som også kalles Jacobi-matrisen til F i \mathbf{x}^* . Dette skiller seg fra gradienten som er en kolonnevektor. Vi kan enkelt utvide til en funksjon $F : \mathbb{R}^N \rightarrow \mathbb{R}^M$ ved å observere at vi kan behandle hver av komponentene separat.³

$$F(\mathbf{x}) = [F_1(\mathbf{x}), \dots, F_M(\mathbf{x})]' \in \mathbb{R}^M \quad (12.11)$$

For å finne lineær tilnærming er det bare å derivere hver av de M komponent funksjonene med hensyn på de N inputvariablene

$$DF(\mathbf{x}^*) = [DF_1(\mathbf{x}^*), \dots, DF_M(\mathbf{x}^*)]' \quad (12.12)$$

12.4.2 Taylor-tilnærming

Funksjoner kan være komplisert å arbeide med analytisk. Vi kan forsøke å finne en funksjon p der $p(x) \approx f(x)$ for x i et nabolag til a . Det kan være rimelig å kreve at $p(a) = f(a)$ og at $p'(x)|_a = f'(x)|_a$. Dette gir

$$f(x) \approx p(x) := f(a) + f'(a)(x - a) \quad (12.13)$$

Vi kan bruke dette til å beregne endring i y for små endringer i x

$$dy := p(x + dx) - p(x) \approx \Delta y := f(x + dx) - p(x) \quad (12.14)$$

Dette er en lineær tilnærming av funksjonen, også kalt første ordens taylor tilnærming. Hvis den n 'te deriverte til f er definert i a kan vi generelt definere n 'te ordens taylor tilnærming av f i a som

$$p(x) = f(a) + \sum_n^N \frac{1}{n!} f^{(n)}(a)(x - a)^n \approx g(x) \quad (12.15)$$

³Husk at det finnes ulike måter å betrakte samme transformasjon. Vi kan for eksempel tenke på $\cos(x^2) := f(x)$ eller som $g(h(x))$ der $h : x \mapsto x^2$ og $g : y \mapsto \cos(y)$.

det er mulig å vise at differansen $g(x) - p(x)$ er gitt ved lagranges feilledd

$$R_{n+1}(x) = \frac{1}{(n+1)!} f^{(n+1)}(c) x^{n+1} \quad (12.16)$$

Størrelsen på leddet avhenger av c som er vanskelig å finne. Vi kan velge c som gir størst absoluttverdi av $f^{(n+1)}(c)$ for å finne en øvre begrensning på feilen i intervallet vi tilnærmer funksjon.

12.5 Noen vanlige funksjoner

12.5.1 Polynomial

12.5.2 Eksponential og logaritmer

Det er mange størrelser som vokser med % av egen verdi. Populasjoner, penger i banken, mm. Dette kan beskrives med en eksponentialfunksjon $f(x) = ca^x$ der

- $f(x+1)/f(x) = ca^{x+1}/ca^x = a \iff f(x+1) = f(x)a$
- $f(3) = c \cdot a \cdot a \cdot a$
- $f(0) = c$.
- $f(\frac{m}{n}) = ca^{m/n}$
- $f(-x) = ca^{-x} = \frac{c}{a^x}$

derivert? valg av grunntall

Det er en monoton funksjon som er strengt voksende for $a > 1$ og strengt avtagende for $a < 1$. Hvis $a = 1$ er $f(x) = c$ for alle x . Med unntak av dette tilfelle har eksponentialfunksjoner en invers som kalles logaritmen

egenskaper, litt usikker på hvordan utleder

- $\log(xy) = \ln(x) + \ln(y)$
- $\log(\frac{x}{y}) = \ln(x) - \ln(y)$
- $\log(\frac{1}{x}) = 1 - \ln(x)$
- $\log(x^r) = r \cdot \ln(x)$

Naturlig logaritme

One base to rule them all.

$$a = e^{\ln(a)} \iff a^x = e^{\ln(a)x} = e^{\lambda x}, \quad \text{der } \lambda = \ln(a) \quad (12.17)$$

Log-lineær

Tror dette er når sammenhengen mellom variabler i lineær i logartimen, eks:

$$y = Ax^a \implies \log y = \log[Ax^a] = \log A + \log x^a = \log A + a \log x \quad (12.18)$$

12.5.3 Trigonometriske funksjoner

12.5.4 Sammensatte funksjoner

12.6 Kort om integral

Det fundamentale teoremet i kalkulus er

$$\int_0^x f(s)ds := F(x) \quad (12.19)$$

der F er integralet eller den antideriverte av f . Vi kan tenke oss at vi beveger oss på intervallet $[0, x]$ av tallinjen der variabelen s befinner seg og tar en slags vektet sum der hvert punkt blir vektet med $f(s)$. En alternativ fremgangsmåte er å eksplisitt betrakte en todimensjonal inputmengde som begrenset av $[0, f(s)]$ i t -retningen og $[0, x]$ i s -retningen. Vi vokter nå alle punkter likt slik at det ikke er noen vektingsfunksjon inne i integralet

$$\int_0^x \int_0^{f(s)} dt ds = \int_0^x t|_0^{f(s)} ds = F(s)|_0^x = F(x) \quad (12.20)$$

Vi kan selvfølgelig spesifisere en variabel som bestemmer nedre grense for intervallet. De analoge uttrykkene blir da

$$\int_a^b f(s)ds := F(a) - F(b) \quad (12.21)$$

og

$$\int_a^b \int_0^{f(s)} dt ds = \int_a^b t|_0^{f(s)} ds = F(s)|_a^b = F(a) - F(b). \quad (12.22)$$

Vi kan også innføre en annen nedre grense enn 0 i t -retningen. Med vektingstilnærming blir det differanse mellom vektene,

$$\int_0^x [f(s) - g(s)]ds = \int_0^x f(s)ds - \int_0^x g(s)ds = F(x) - G(x) \quad (12.23)$$

dersom $f(s) > g(s)$ for alle $s \in [0, x]$. Eventuelt må vi partisjonere intervallet av tallinjen og ta differansen separat. Kanskje mulig å bruke absoluttverdi... hm. Kan finne tilsvarende

med dobbeltintegral uten vektning av punktene,

$$\int_0^x \int_{g(s)}^{f(s)} dt ds = \int_0^x t|_{g(s)}^{f(s)} ds = \int_0^x [f(s) - g(s)] ds \quad (12.24)$$

12.7 Flervariable funksjoner

Betrakt en funksjon

$$f : \mathbb{R}^d \rightarrow \mathbb{R} \quad (12.25)$$

$$: \mathbf{x} = (x_1, \dots, x_d) \mapsto f(\mathbf{x}) \quad (12.26)$$

Generaliseringen av den deriverte til funksjonen er gradienten som angir den partiellderivate med hensyn på hver av variablene i inputvektoren.⁴

$$\nabla f : \mathbb{R}^d \rightarrow \mathbb{R}^d \quad (12.27)$$

$$: \mathbf{x} \mapsto \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_d} \right)' \quad (12.28)$$

Gradienten er en vektor som lever i inputspace. Eller i utgangspunktet er det en vektor av funksjoner. Det blir vektor av tall når vi angir spesifikk verdi av \mathbf{x} . Uansett hvor vi evaluere vil vektoren peke i retningen der funksjonen vokser raskest. Lengden på vektoren sier noe om hvor raskt den vokser.

Generaliseringen av den andrederiverte til funksjonen er hesse-matrisen

$$\mathbf{H}f : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d} \quad (12.29)$$

$$: \mathbf{x} \mapsto \mathbf{H}f(\mathbf{x}) \quad (12.30)$$

der

$$\mathbf{H}f(\mathbf{x})_{i,j} = \frac{\partial}{\partial x_j} \left(\frac{\partial f(\mathbf{x})}{\partial x_i} \right) \quad (12.31)$$

Funksjonen er strengt konkav hvis⁵:

$$\mathbf{x}' [\mathbf{H}f(\mathbf{x})] \mathbf{x} > 0, \quad \forall \mathbf{x} \neq \mathbf{0} \quad (12.32)$$

⁴Tror vel egentlig generaliseringen er jacobimatrisen som er et lineær map. Må avklare forhold mellom jacobi og gradient..

⁵Merk at den kvadratiske formen er generaliseringen for om en matrise er *positiv* eller *negativ*. Sier forhåpentligvis mer om dette i delen om lineær algebra.

12.8 Ubetinget optimering

12.8.1 Gradient descent

I maskinl ring har vi en kostnadsfunksjon $C : \Theta \rightarrow \mathbb{R}$ som angir sum av tap til hver av observasjonene i treningsdata for kandidatparameter θ . Jeg vil velge kandidaten som minimerer kostnaden. En mulighet er   finne gradientent $\nabla_{\theta}C$ og l se $\nabla_{\theta}C(\theta) = 0$. For   finne kandidater til optimum. Utfordringen er at $C(\cdot)$ kan v re vilk rlig komplisert og avhenge av mange variabler slik at vi ikke klarer   l se det ikke-line re ligningssystemet analytisk. Alternativet er da   g  frem numerisk. Hvis vi har et eksplisitt uttrykk for gradienten kan vi evaluere den i vilk rlig θ_{start} . Gradienten er vektor i parameterrommet som peker i retning der kostnaden vokser raskest og lengden avhenger av hvor bratt funksjonen er. Vi velger derfor   g  i helt motsatt retning,

$$\theta_{ny} = \theta_{start} - \eta \nabla_{\theta}C(\theta_{start}) \quad (12.33)$$

der η er den s kalte l ringsraten som skalerer gradienten og p virker hvor raskt vi beveger oss i parameterrommet. Det er viktig at den ikke er for h y slik at vi overskyter bunnpunktet og begynner   divergere, men b r heller ikke v re for lav slik at konvergens tar for lang tid. Bestemmer oss et treshold som avslutter algoritmen n r $\|\nabla_{\theta}\| < k$ siden vi i praksis ikke for den eksakt lik null.

Merk at dette er en lokal metode som bare kjenner helning til funksjon akkurat der den blir evaluert. Sikrer kun at det lokale minimum som den konvergerer mot ogs  tilsvarer det globale minimum dersom funksjonen er konveks. Skal n  se p  alternativ fremgangsm te som kan h ndtere funksjoner med plat  og flere lokale minimum.

Stokastisk gradient descent

Denne fremgangsm ten bruker kun subset av treningsdata n r den evaluerer hvordan kostnaden blir p virker av valg av parameter. N r vi evaluerer i ulike subsets f r vi ulike kostnadsfunksjoner slik at gradient hopper litt rundt om kring, men i gjennomsnitt s  vil den konvergere mot global minimum. Kan v re effektiv fordi den bruker mindre data og dermed kj rer raskere, selv om det tar flere steg og en litt mindre ryddig vei mot m let. Ogs  fordel at den kan hoppe ut at av blindveier.

12.9 Betinget optimering

Kapittel 13

Lineær algebra

13.1 Vektorer

En vektor er en tuple med reelle tall, $\mathbf{x} = (x_1, \dots, x_N)$, der $x_n \in \mathbb{R}$ for $n = 1, \dots, N$. De er to grunnleggende operasjoner som er definert på vektorer: skalering og summering. I tillegg er det definert noen andre konsept som ikke lager nye vektorer

- Indre produkt: $\langle \mathbf{x}, \mathbf{y} \rangle = \sum x_n y_n$
- Eukledisk norm : $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$

Det er noen grunnleggende ulikheter som gjelder for disse

- Triangelulikheten: $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$
- Cauchy-Schwarz-ulikheten: $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$

Det er mulig å uttrykke nye vektorer som lineær kombinasjon av eksisterende. Hvis vi har en mengde av vektorer $X = \{\mathbf{x}_1, \dots, \mathbf{x}_K\} \subset \mathbb{R}^N$ så finnes det en mengde

$$\{y \in \mathbb{R}^N | y = \sum \alpha_k x_k \text{ der } \alpha_k \in \mathbb{R} \text{ og } \mathbf{x}_k \in X \text{ for } k = 1, \dots, K\} \quad (13.1)$$

Denne mengden av vektorer som kan uttrykkes som lineær kombinasjon av vektorene i X betegnes som $span(X)$. Som vi skal se er dette vesentlig for å vurdere eksistens av løsning på lineære ligningssystem, altså om det eksisterer en \mathbf{x} slik at $y = \sum \alpha_k x_k$ for en gitt y . Løsningen eksisterer hvis og bare hvis $y \in span(X)$.

En mengde av vektorer er lineært uavhengige hvis

$$\sum \alpha_k x_k = \mathbf{0} \implies \mathbf{a} = \mathbf{0} \quad (13.2)$$

Dette impliserer også at ingen vektorer i mengden kan uttrykkes som en lineær kombinasjon av de resterende vektorene. Det er vesentlig for unikhhet av løsning siden hvis X er

en lineær uavhengig mengde vil

$$y = \sum \alpha_k x_k = y = \sum \alpha'_k x_k \implies y = \sum (\alpha_k - \alpha'_k) x_k = \mathbf{0} \implies \alpha_k = \alpha'_k \quad (13.3)$$

En mengde av lineært uavhengige vektorer Z utgjør en basis for $\text{span}(X)$ hvis $\text{span}(Z) = \text{span}(X)$. For øvring utgjør $\text{span}(X)$ et underrom av \mathbb{R}^N siden det er lukket under skalering og addisjon. Generelt er en mengde S et underrom av \mathbb{R}^N hvis det for alle $\alpha \in \mathbb{R}$

$$\bullet \mathbf{x} \in S \implies \alpha \mathbf{x} \in S$$

$$\bullet \mathbf{x}, \mathbf{y} \in S \implies \mathbf{x} + \mathbf{y} \in S$$

og dimensjonen til et underrom er antall vektorer i basisen.

13.2 Matriser

Kan bruke matriser til å *stacke* lineære ligningssystem,

$$a_{11}x_1 + a_{12}x_2 = y_1 \quad (13.4)$$

$$a_{21}x_1 + a_{22}x_2 = y_2 \quad (13.5)$$

tilsvarer

$$\begin{bmatrix} a_{11}x_1 + a_{12}x_2 \\ a_{21}x_1 + a_{22}x_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \quad (13.6)$$

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \quad (13.7)$$

$$\mathbf{A}\mathbf{x} = \mathbf{y} \quad (13.8)$$

Vi bruker \mathbf{A}_{ij} til å referere til komponent fra i 'te rekke og j 'te kolonne. Ofte har vi lyst til å gi en representasjon av \mathbf{A} uten å beskrive alle de individuelle komponentene. Vi kan da gruppere de inn i rekke- og kolonnevektorer. Jeg er litt usikker på hvilken notasjon jeg bruker.¹ Jeg kan da gi alternativ representasjon av matrisemultiplikasjon

¹En mulighet er å bruke \mathbf{A}_i til å betegne rekke og \mathbf{A}_j for kolonne. Et mulig problem er at det er ambiguitet om jeg betrakter \mathbf{A}_i som en kolonnevektor; altså transponerte av rekken i matrisen. Også problem om jeg vil referere til spesifikk tall. Kan bruke $\mathbf{a}_{\bullet 1}$ og $\mathbf{a}_{1\bullet}$ til å referere til henholdsvis første kolonne og rekke. hmm

med utgangspunkt i disse blokkene,

$$\mathbf{Ax} = \begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \mathbf{a}_1 x_1 + \mathbf{a}_2 x_2 \quad (13.9)$$

$$= \begin{bmatrix} \mathbf{a}_{1\bullet} \\ \mathbf{a}_{2\bullet} \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{a}_{1\bullet} \mathbf{x} \\ \mathbf{a}_{2\bullet} \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{a}'_1 \mathbf{x} \\ \mathbf{a}'_2 \mathbf{x} \end{bmatrix}, \quad \mathbf{A} := \begin{bmatrix} \mathbf{a}'_1 \\ \mathbf{a}'_2 \end{bmatrix} \quad (13.10)$$

der jeg i siste likhet lar \mathbf{a}_1 være rekke i matrisen på kolonneform.² Det er også poeng at man kan blokkpartisjonere matriser på andre måter og gjøre operasjon på blokkene så lenge de er kompatible, men det må bli en annen gang.

$$A = \begin{bmatrix} a & \cdots \\ \vdots & \end{bmatrix} \quad (13.11)$$

13.2.1 Derivasjon

Jeg vil nå begynne å derivere med hensyn på vektor. Det er litt som å ta partiell derivert med hensyn på hver av komponentene i en matrise og stacke de oppå hverandre,

$$\frac{d}{d\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_K} f(\mathbf{x}) \end{bmatrix} \quad (13.12)$$

Litt usikker på dimensjonene. Går ut i fra at den deriverte har samme dimensjon som \mathbf{x} . Noen regneregler³

$f(\mathbf{x})$	$\frac{d}{d\mathbf{x}} f(\mathbf{x})$
\mathbf{Ax}	\mathbf{A}
$\mathbf{a}'\mathbf{x}$	\mathbf{a}
$\mathbf{x}'\mathbf{a}$	\mathbf{a}
$\mathbf{x}'\mathbf{x}$	$2\mathbf{x}$
$\mathbf{x}'\mathbf{Ax}$	$2\mathbf{Ax}$

Eksempel: minste kvadrat

Har tapsfunksjon

$$L = \frac{1}{N} \sum (y_n - \mathbf{x}'_n \mathbf{b})^2 \quad (13.13)$$

²Notasjon kan bli ganske forvirrende

³Merk at de resuserer til vanlige derivasjonsregler dersom matriser og vektor bare har én komponent

Vil ha en vektor der n 'te komponent er $\mathbf{x}'_n \mathbf{b}$. Tilsvarende \mathbf{Xb} der $\mathbf{X}_{n\bullet} = \mathbf{x}'_n$. Kan da skrive det på matriseform,

$$L = \frac{1}{N}(\mathbf{y} - \mathbf{Xb})'(\mathbf{y} - \mathbf{Xb}) \quad (13.14)$$

$$= \frac{1}{N}(\mathbf{y}' - \mathbf{b}'\mathbf{X}')(\mathbf{y} - \mathbf{Xb}) \quad (13.15)$$

$$= \frac{1}{N}(\mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbf{Xb} - \mathbf{b}'\mathbf{X}'\mathbf{y} + \mathbf{b}'\mathbf{X}'\mathbf{Xb}) \quad (13.16)$$

$$= \frac{1}{N}(\mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{Xb} + \mathbf{b}'\mathbf{X}'\mathbf{Xb}) \quad (13.17)$$

Ser bort i fra skaleringen $\frac{1}{N}$ og deriverer med hensyn på \mathbf{b} ,

$$\frac{dL}{d\mathbf{b}} = 2\mathbf{X}'\mathbf{y} - 2\mathbf{X}'\mathbf{Xb} = \mathbf{0} \quad (13.18)$$

$$\implies \mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (13.19)$$

gitt at $\mathbf{X}'\mathbf{X}$ er invertibel.

13.3 Lineære transformasjoner

En funksjon $T : \mathbb{R}^K \rightarrow \mathbb{R}^N$ er lineær transformasjon hvis

$$T(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha T\mathbf{x} + \beta T\mathbf{y} \quad (13.20)$$

for alle $\mathbf{x}, \mathbf{y} \in \mathbb{R}^K$ og $\alpha, \beta \in \mathbb{R}$. Vi skriver funksjonen med stor bokstav og uten parentes fordi den oppfører seg som en matrise. Skal vise senere at det er én-til-én korrespondanse mellom matriser og lineære transformasjoner mellom vektorrom. Dette er en veldig grei egenskap til lineære transformasjoner som gjør de er mye brukt i anvendt matematikk. Definisjonen over kan generaliseres til

$$T\left(\sum \alpha_k \mathbf{x}_k\right) = \sum \alpha_k T\mathbf{x}_k \quad (13.21)$$

dette impliserer at

$$T\mathbf{x} = \sum \alpha_k T\mathbf{e}_k \quad (13.22)$$

slik at $\text{rng}(T) = \text{span}(V)$, der $V = \{T\mathbf{e}_1, \dots, T\mathbf{e}_K\}$. For $T : \mathbb{R}^N \rightarrow \mathbb{R}^N$ så er det ekvivalens mellom ikke-singularitet og masse greier. Det er en ideell situasjon siden det alltid eksisterer en unik løsning. I praksis må vi ofte finne tilnærmet løsning fordi hvis $T : \mathbb{R}^K \rightarrow \mathbb{R}^N$, der $K < N$, så kan det være slik at $\mathbf{y} \notin \text{rng}(T)$. Altså finnes det ingen \mathbf{x} slik at $T\mathbf{x} = \mathbf{y}$. Vår beste løsning da er å finne \mathbf{x}' som minimerer $\|\mathbf{y} - T\mathbf{x}'\|$. For å minimere avstand mel-

lom to en vektor og et underrom får vi bruk for ortogonale projeksjoner, fordi løsningen er å gå den strakeste vegen.

13.4 Ortogonale projeksjoner

To vektorer er ortogonale hvis $\langle \mathbf{x}, \mathbf{y} \rangle = 0 \iff \mathbf{x} \perp \mathbf{y}$. Dette konseptet generaliserer også til andre objekt som vi kan definere indre produkt på, som f.eks. tilfeldige variabler. En vektor kan også være ortogonal på en mengde S

$$\mathbf{x} \perp S \iff \mathbf{x} \perp \mathbf{z} \text{ for alle } \mathbf{z} \in S \quad (13.23)$$

Mengden av alle vektorer som er ortogonal på mengden S utgjør dets ortogonale komplement

$$S^\perp = \{\mathbf{x} | \mathbf{x} \perp \mathbf{z}, \text{ for alle } \mathbf{z} \in S\} \quad (13.24)$$

En mengde av vektorer X er ortogonale hvis vektorene er parvise ortogonale $\mathbf{x}_j \perp \mathbf{x}_k$, $j \neq k$. Den er i tillegg ortonormal hvis $\|\mathbf{x}\| = 1$ for alle $\mathbf{x} \in X$. Dette er en veldig grei egenskap siden det gjør det enkelt å finne vektene i lineære kombinasjoner

$$\mathbf{y} = \sum \alpha_k \mathbf{x}_k = \sum \langle \mathbf{y}, \mathbf{x}_k \rangle \mathbf{x}_k \quad (13.25)$$

Dette gjør de velegnet som basiser for vektorrom. Det eksisterer alltid ortonormal basiser og vi kan bruke algoritmer (eg. Gram-Schmidt) for å konstruere.

Uansett, det store resultatet er ortogonale projeksjons theoremet! La S være et underrom av \mathbb{R}^N og $\mathbf{y} \in \mathbb{R}^N$. Vi vil finne

$$\hat{\mathbf{y}} = \arg \min_{\tilde{\mathbf{y}} \in S} \|\tilde{\mathbf{y}} - \mathbf{y}\| \quad (13.26)$$

Theorem sier da at dette har en unik løsning der $\mathbf{y} - \hat{\mathbf{y}} \perp S$. Merk at for alle andre $\mathbf{z} \in S$ så er

$$\|\mathbf{y} - \mathbf{z}\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{y}} - \mathbf{z}\|^2. \quad (13.27)$$

Som et eksempel kan vi projekte \mathbf{y} på $\mathbf{1}$. Vil finne $\hat{\mathbf{y}} \in \text{span}(\mathbf{1})$ som minimerer avstand til \mathbf{y} og vet at $\langle \mathbf{y} - \alpha \mathbf{1}, \mathbf{1} \rangle = 0$. Kan da finne $\alpha = \bar{y}_N$. Mer generelt kan vi betrakte projeksjon

på et vektorrom S med flere dimensjoner. La $S = \text{span}(X)$, der $X := \{\mathbf{x}_1, \dots, \mathbf{x}_K\}$.

$$\mathbf{y} - \hat{\mathbf{y}} \perp \mathbf{x}_k, \quad k = 1, \dots, K \quad (13.28)$$

$$\implies \mathbf{x}'_k(\mathbf{y} - \mathbf{X}\beta) = 0, \quad k = 1, \dots, K \quad (13.29)$$

$$\implies \mathbf{X}'(\mathbf{y} - \mathbf{X}\beta) = \mathbf{0} \quad (13.30)$$

$$\implies \beta = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \quad (13.31)$$

Det eksisterer en lineær transformasjon som utfører projeksjonen, $\mathbf{P} : \mathbf{y} \mapsto \hat{\mathbf{y}} = \text{proj}_S \mathbf{y}$. Fra utledningen over følger det at $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}$. Denne transformasjonen har en del egenskaper som følger ganske intuitivt fra at det er en projeksjon. Ofte er vi også interessert i residualen, som vi kan finne med $\mathbf{M} := \mathbf{I} - \mathbf{P}$.

La nå S ha en ortonormal basis U . Merk at $\mathbf{u}'_j \mathbf{u}_k = 0, j \neq k$ og lik 1 hvis $j = k$. Det følger da at $\mathbf{U}'\mathbf{U} = \mathbf{I}$ slik at $\mathbf{P} = \mathbf{U}\mathbf{U}'$ og $\mathbf{P}\mathbf{y} = \sum \mathbf{u}_k \langle \mathbf{u}_k, \mathbf{y} \rangle$. Merk generelt at $\mathbf{y} \in \text{span}(X)$ alltid kan skrives som $\mathbf{X}\mathbf{b}$ for noen \mathbf{b} , men det kan være vanskelig å finne vektene i den lineære kombinasjonen. Med ortonormal basis blir dette enklere fordi vektene ikke bidrar i samme retninger og enhetslengde gjør det enkelt å finne riktig skalering... Har ikke helt intuisjonen på dette, men gir sann omtrent mening.

13.5 Kvadratisk form

Den kvadratiske formen til en matrise \mathbf{A} er $\mathbf{x}'\mathbf{A}\mathbf{x}$. Matrisen er positiv (semi-)definit hvis $\mathbf{x}'\mathbf{A}\mathbf{x}(\geq) > 0$ for alle \mathbf{x} . Det er litt analog til om matrisen er positiv eller negativ...

Merk at matriser på formen $\mathbf{B}'\mathbf{B}$ alltid er positiv (semi-)definit fordi

$$\mathbf{x}'\mathbf{B}'\mathbf{B}\mathbf{x} = (\mathbf{B}\mathbf{x})'(\mathbf{B}\mathbf{x}) = \|\mathbf{B}\mathbf{x}\|^2 \geq 0 \quad (13.32)$$

For hver positiv (semi-)definit \mathbf{A} er det mulig å dekomponere $\mathbf{A} = \mathbf{B}'\mathbf{B}$ der \mathbf{B} ikke er unik. Kan for eksempel bruke Cholesky-dekomponering. Tror dette er litt analog til kvadratroten av positiv skalar.

Kapittel 14

Matematisk tenkemåte

Det matematiske språket har en logisk oppbygning. Logikk er læren om hva som gjør argumentasjon gyldig. For veldefinerte spørsmål i matematiske modeller vil det være entydige svar fordi konklusjonen følger logisk fra informasjonen som er tilgjengelig. Mer generelt vil konklusjoner avhenge av skjulte antagelser... slik at man trekker ulike konklusjoner. Kan knytte dette til modell (liten verden versus stor verden). Viktig for konstruksjon av argument og da spesielt matematiske bevis.

Det matematiske språket medfører og fasiliterer abstraksjon. Vi fokuserer på de vesentlige egenskapene og finner en enklere representasjon. Vi kan benytte oss av matematiske objekt med gitt definisjon som avgrenser hvilke størrelse som tilhører en gitt type objekt. Kan deretter finne egenskaper til objekter av gitt type og definere operasjoner på disse. Det gir oss et generelt rammeverk som vi kan anvende på mange ulike problemstillinger. Grunnleggende objekter er mengder, tupler og tall. Videre skal vi se på relasjoner, hvorav funksjoner er et viktig spesialtilfelle. Til slutt skal vi også se på grafer som kan beskrive en mer hierarkisk struktur (ikke bare hvorvidt det er kobling men hvor mange ledd unna samt .. mer info.)

14.1 Logikk

Logikk handler mye om hva som følger. Hvis noe er sant, så kan vi vise at andre ting også må være sanne. Jeg er fra Bergen og Bergen er en by i Norge. Dermed følger det logisk at jeg er fra Norge. Jeg vil bruke logikk til å konstruere gyldige argument der jeg viser at en konklusjon må være sanne gitt at premissene er sanne. Det kan også være nyttig å manipulere uttrykk som gir det en annen representasjon, men bevarer samme meningsinnhold. Vi begynner med klassisk utsagnslogikk som kun består av sammensetninger av enkle utsagn med logiske bindeord. Deretter vil jeg se på førsteordens logikk som blant annet bruker begreper fra mengdelære og gir oss et mer fleksibelt språk til å representere utsagn og sammenhenger mellom disse. På denne måten kan vi (nesten) konstruere hele det matematiske språket fra bunnen av.

14.1.1 Utsagnslogikk

Vi vil finne en representasjon av meningsinnholdet i språklige yttringer. De består grunnleggende av *utsagn* som enten kan være sanne eller usanne. Disse kan representeres med en *utsagnsvariabel*. Dette er plassholder for *valuasjon* av enten 1 hvis sann eller 0 hvis usann. Det kan også være andre slags yttringer, for eksempel "hurra!", men dette påvirker ikke gyldigheten til argumentasjonen og vi kan derfor se bort i fra dette. Videre vil yttringene bestå av *logiske bindeord* som og, eller, ikke og hvis, så. Disse representeres med henholdsvis $\wedge, \vee, \neg, \rightarrow$. Med utsagnsvariabler og logiske bindeord kan vi konstruere *formler*. Mengden av formler \mathcal{F} er definert som den induktive tillukningen av utsagnsvariabler under de logiske bindeordene.

Vi kan betrakte de logiske bindeordene som funksjoner $f : \{0, 1\}^2 \rightarrow \{0, 1\}$.¹ Vi kan konstruere sannhetstabeller som gir valuasjon av formelen for alle kombinasjoner av valuasjoner av de atomære formulene som kun består av utsagnsvariabel. For en gitt valuasjon kan det være praktisk å bruke et tre til å representere hierarkisk struktur av det sammensatte formelen og propagere sannhetsverdi oppover.

Nødvendighet og tilstrekkelighet

Jeg vil ha noen begreper for å beskrive hvordan valuering av en formel påvirker hva vi vet om en annen formel. I beskrivelsen under bruker jeg utsagn, men merk at disse er atomære formler og argumentene gjelder generelt for formler.

- Et utsagn P er *tilstrekkelig* for Q dersom Q må være sann dersom P er sann, altså $P \rightarrow Q$. Dette kan også uttrykkes som *hvis P , så Q* . Merk at tilstrekkelig ikke impliserer *nødvendig*. Det kan derfor være slik at Q er sann uten at P er det. Kun en valuasjon av Q gir oss derfor ingen informasjon om sannhetsinnholdet i P .
- En påstand P er *nødvendig* for Q dersom det ikke er mulig at Q er sann uten at også P er sann, altså $P \leftarrow Q$. Dette kan også uttrykkes som *Q bare hvis P* . Merk at nødvendig ikke impliserer tilstrekkelig. Kun en valuasjon av P gir oss derfor ingen informasjon om Q .
- Hvis påstanden P er både nødvendig og tilstrekkelig for Q er påstandene ekvivalente i betydning av at begge med nødvendigheten enten er sanne eller usanne samtidig, altså $P \leftrightarrow Q$. Dette kan også uttrykkes som *Q hvis og bare hvis P* .

Logisk ekvivalens

To formler er logisk ekvivalente dersom de har samme sannhetsverdi for alle mulige tilordninger av sannhetsverdier til de atomiske påstandene de består av. En mulighet for

¹Med unntak av \neg som er $f : \{0, 1\} \rightarrow \{0, 1\}$.

å bevise ekvivalens er derfor å konstruere sannhetstabellene. En alternativ fremgangsmåte er å bruke et resonement. For å vise at formlene F og M er ekvivalente må vi

1. Anta at F er sann og vise at da må M være sann.
2. Anta at M er sann og vise at da må F være sann.

I denne argumentasjonsrekken kan vi få bruk for å manipulere uttrykk slik at de får ny representasjon, men samme meningsinnhold. Noen viktige ekvivalenslover

- Distributativ: $A \vee (B \wedge C) \leftrightarrow (A \vee B) \wedge (A \vee C)$
- DeMorgan: $\neg(A \vee B) \leftrightarrow (\neg A \wedge \neg B)$
- Assosiativ: $A \vee (B \vee C) \leftrightarrow (A \vee B) \vee C$
- Kommutativ: $A \wedge B \leftrightarrow B \wedge A$

Fra de grunnleggende konnektivene kan vi definere nye konnektiver som gir oss enklere måte å representere formler med samme sannhetsinnhold.

- $P \wedge Q \leftrightarrow \neg(\neg P \vee \neg Q)$
- $P \rightarrow Q \leftrightarrow \neg(P \wedge \neg Q)$
- $(P \leftrightarrow Q) \leftrightarrow (P \rightarrow Q) \wedge (Q \rightarrow P)$

Logisk konsekvens

Hvis vi godtar et premiss eller gjør antagelser om at noen påstander er sanne, så vil det være andre påstander som følger med nødvendighet av dette; de må, av logisk konsekvens, også være sanne. Mer presist er F en logisk konsekvens av M dersom F alltid er sann når M er sann. Et argument er gyldig dersom konklusjonen er en logisk konsekvens av mengden med antagelser. Et eksempel på gyldig argument er

$$\frac{P \vee Q \quad \neg P}{\therefore Q}$$

Noen flere begreper

Vi kan innføre litt begreper om hvordan sannhetsinnhold til en formel avhenger av valuering til de atomære formlene (utsagnsvariabler) den består av.

- Oppfyllbar dersom det eksisterer en valuering der formelen er sann.
- Falsifiserbar dersom det eksisterer en valuering der den ikke er sann.

- Tautologi dersom den alltid er sann (ikke falsifiserbar).
- Motsigelse dersom den alltid er usann (ikke oppfylld).

14.1.2 Første ordens logikk

Vi skal gå fra utsagnslogikk til første ordens logikk, også kalt for predikatlogikk. Dette gir et mer fleksibelt rammeverk til å uttrykke påstander om egenskaper til elementer i en mengde, samt si ting om sammenheng mellom ulike elementene. Vi kan bruke relasjon til å uttrykke påstand om element, $P(x)$. Vi kan la x være en *variabel* som representerer eller er plassholder for element i universet Ω . Uttrykket $P(x)$ er en predikat siden det inneholder en *fri variabel* og vi kan ikke si om det er sant eller usant. Sannhetsmengden til predikatet er $S = \{x \in \Omega : P(x) = 1\} \subset \Omega$. Vi kan i prinsippet beskrive hele sannhetsmengden, men ofte vil bare uttrykke egenskaper ved den, for eksempel $S \neq \emptyset$ eller $S \neq \Omega$. For å gjøre dette får vi bruke for kvantorer.

Kvantorer

Vi vil ofte si noe om innholdet i sannhetsmengden til påstander. Vi har to såkalte *kvantorer*:

1. Universalkvanten: $\forall x P(x)$, betyr at påstand er sann for alle x i universet vi betrakter
2. Eksistenskvanten: $\exists x P(x)$, betyr at det eksisterer minst én x der påstand er sann

I tillegg brukes $\exists! x P(x)$ som betyr at det eksisterer én og bare én x som gjør påstand sann. Dette er mer en notasjonell konvensjon. Vi kan også omformulere uttrykk med quantifiers for å finne ekvivalente representasjoner. Merk at

- $\forall x P(x) \Leftrightarrow \neg \exists x \neg P(x)$. Hvis det er sant for alle x kan det ikke eksistere en x der det ikke er sant. Medfører $\neg \forall x P(x) \Leftrightarrow \exists x \neg P(x)$.
- $\exists x P(x) \Leftrightarrow \neg \forall x \neg P(x)$. Hvis det finnes minst én x der det er sant, så kan det ikke være usant for all x . Medfører $\neg \exists x P(x) \Leftrightarrow \forall x \neg P(x)$.

Merk også at rekkefølgen til *quantifiers* har betydning dersom de er ulike, men dersom de er like kan vi lese $\exists x \exists y$ som at det eksisterer x og y slik at (...), og rekkefølgen har ikke betydning. Merk også at dersom vi bruker quantifiers om to størrelser fra samme univers må vi spesifisere eksplisitt dersom $x \neq y$. Vi kan også bruke denne notasjonen for å si at et predikat $P(\cdot)$ er sant for alle elementer i en mende S :

$$\forall x \in S P(x) \Leftrightarrow \forall x (x \in S \Rightarrow P(x))$$

Formelle definisjoner

Vi vil fortsatt betrakte sannhetsverdi til formler, men disse består ikke lenger bare av induktiv tillukning av atomære formler (utsagn) under logiske bindeord. For det første har vi introdusert to nye logiske symboler med entydig betydning. Deretter vil vi også åpne for at språket vi bruker har en *signatur* som består av *konstanter*, funksjoner og relasjoner som vi selv definerer meningsinnholdet til. Konstantene er navngitte representasjoner av elementene i universet. Funksjonene definerer operasjon som lar oss representere elementer i universet uten at de har en egen konstant. Relasjonene tar *termer* (konstanter eller variabel) som argument, der antall argument er gitt ved *ariteten* til relasjonen. En relasjon R med aritet 2 evaluert i konstante (a, b) kan skrives som $(a, b) \in R \iff aRb \iff R(a, b)$. Det tar verdi i $\{0, 1\}$ og er derfor et utsagn.

Mengden av formler \mathcal{F} under et språk er induktivt definert. Basismengde av såkalte *atomære formler* som er relasjon på termer. Hvis mengden er lukket under funksjoner f, g, \dots på elementene i basismengde så er $f(t_1, \dots, t_n)$ også en term. I tillegg har vi at

1. $\phi, \psi \in \mathcal{F} \rightarrow \phi \wedge \psi, \phi \vee \psi, \dots \in \mathcal{F}$
2. $\phi \in \mathcal{F} \implies \forall x \phi, \exists x \phi \in \mathcal{F}$

Første ordens språk

Kan si litt om språk generelt og ta formell definisjon av første orders språk her.

14.2 Mengdelære

En mengde er en uordnet kolleksjon av distinkte objekter. Dette medfører at rekkefølge og antall forekomster ikke har betydning, slik at $\{a, a, b\} = \{b, a\}$. I eksempelet ble mengdene representert ved å liste opp objektene. Alternativt kan innholdet beskrives ved

$$A = \{x : P(x)\}, \quad (14.1)$$

der $P : S \rightarrow \{0, 1\}$ er et medlemskriterium og S er universet av objekter vi betrakter. Et objekt er da element i mengden A dersom påstanden $P(\cdot)$ er sann (tar verdi 1) når det blir evaluert for det objektet. Ettersom mengder er så fleksible vil vi også ha en mer fleksibel måte å konstruere de. En alternativ måte er å bruke elementer fra en annen mengde I som indeks når vi konstruerer

$$P = \{p_i : i \in I\} \quad (14.2)$$

Rangering og operasjoner

En mengde B er en delmengde av A hvis alle elementene i B også er element i A

$$B \subset A \iff x \in B \implies x \in A \quad (14.3)$$

To mengder er like dersom de inneholder akkurat de samme elementene

$$A = B \iff A \subset B \wedge B \subset A \quad (14.4)$$

Vi kan også konstruere nye mengder ved å utføre *operasjoner* på eksisterende

- Union: $A \cup B = \{x \in S : x \in A \text{ eller } x \in B\}$
- Interseksjon eller snitt: $A \cap B = \{x \in S : x \in A \text{ og } x \in B\}$
- Differanse: $A \setminus B = \{x \in S : x \in A \text{ og } x \notin B\}$
- Komplement: $A^C = \{x \in S : x \notin A\}$
- Symmetrisk differanse: $A \triangle B = \{x \in S : (x \in A \wedge x \notin B) \cup (x \notin A \wedge x \in B)\}$

Univers av tall

Mengder kan i utgangspunktet inneholder alle slags objekter, men i praksis liker vi å jobbe med mengder av tall.

- \mathbb{R} , mengden av reelle tall, det vil si alle tall på tallinjen.
- \mathbb{Z} , menden av heltal $\{\dots, -1, 0, 1, \dots\}$
- \mathbb{Q} , mengden av rasjonelle tall, det vil si tall som kan skrives som brøk av heltall.
- \mathbb{N} , mengden av naturlige tall, $\{0, 1, \dots\}$. Merk at noen ikke inkluderer 0 som naturlig tall.

Det er også vanlig å avgrense disse mengdene, for eksempel ved å kun betrakte positive reelle tall. Det kan betegnes som \mathbb{R}^+ .

Intervaller utgjør viktige delmengder. Eksempler på intervall er

$$A = \{x | a < x < b\} = (a, b) \subset \mathbb{R} \quad (14.5)$$

$$B = \{x | a \leq x \leq b\} = [a, b] \subset \mathbb{R} \quad (14.6)$$

der det første er åpnet og det andre er lukket. Det eksisterer en presis definisjon på om en mengde er åpen eller lukket som avhenger av om den inkluderer endepunktene. Vi kan beskrive nye mengder med utgangspunkt i eksisterende mengder.

Kardinalitet

Kan si at $|A|$ er antallet elementer i A . En mengde er uendelig dersom det ikke eksisterer et tall som representerer kardinaliteten til mengden. To mengder har samme kardinalitet, $|A| = |B|$, dersom det eksisterer en bijektiv funksjon mellom elementene i mengdene. Dette gjør det mulig å sammenligne kardinalitet til uendelige mengder. En mengde A er tellbar dersom det finnes en injektiv transformasjon $f : A \rightarrow \mathbb{N}$.

Identiteter

Vil knytte operasjon på mengder til logiske konnektiver... ekvivalens.

Merk at det finnes ulike måter å gi ekvivalente representasjoner av samme uttrykk. For eksempel er $A \setminus B = \{x \in S \mid x \in A \text{ og } x \notin B\} = x \in S \wedge x \in A \wedge x \notin B$. Tror jeg.. de har samme sannhetsmengde i hvertfall.. Vi kan analysere den logiske formen til uttrykk om mengder og betrakte sammenhengen mellom reglene for operasjoner på mengder og reglene for ekvivalens mellom uttrykk. Eksempel

$$x \in A \setminus (B \cap C) \quad (14.7)$$

$$P \wedge \neg(Q \wedge R) \quad (14.8)$$

$$P \wedge (\neg Q \vee \neg R) \quad (14.9)$$

$$(P \wedge \neg Q) \vee (P \wedge \neg R) \quad (14.10)$$

$$x \in A \setminus B \cup A \setminus C \quad (14.11)$$

der jeg brukte at $x \in A$ (osv.) er påstander og dermed kan representeres med bokstav. Sannhetsverdi avhenger av variabel x så jeg kunne også ha betegnet det med $P(x)$.

Mengder av mengder

Det kan i prinsippet være alle mulige typer objekter, inkludert andre mengder. Mengder av mengder betegnes ofte som en familie. Et eksempel på dette er *power set* til en mengde A som består av alle delmengdene til A .

$$\mathcal{P}(A) = \{x : x \subseteq A\} \quad (14.12)$$

Vi kan ha en kolleksjon av mengder som potensielt er uendelig. La $\mathcal{F} = \{A_i : A_i = [\frac{1}{i}, 1], i \in \mathbb{Z}^+\}$. Denne kolleksjonen er monotont voksende fordi $j > i \implies A_i \subset A_j$. Vi kan ta unionen av alle mengdene og se om de konvergerer til en gitt mengde,

$$\cup_{i=1}^{\infty} A_i = (0, 1] \quad (14.13)$$

Vi sier at en følge av hendelser A_1, A_2, \dots er stigende dersom $A_n \subset A_{n+1}, n \geq 1$ og avtagende hvis $A_{n+1} \subset A_n, n \geq 1$. Grenseverdien til slike følger er definert ved

- $\lim_{n \rightarrow \infty} A_n = \bigcup_{i=1}^{\infty} A_i$, hvis stigende
- $\lim_{n \rightarrow \infty} A_n = \bigcap_{i=1}^{\infty} A_i$, hvis avtagende

Partisjonering

En partisjonering av en mengde S er en mengde av ikke-tomme delmengder $S_k \subset S$ der

1. Unionen av delmengdene utgjør mengden, $\bigcup S_k = S$
2. Delmengdene er disjunkte slik at snittet av distinkte delmengder er tomt, $S_k \cap S_j = \emptyset, k \neq j$

Kan innføre begrep som rangere partisjonering... Finere.

Kartesisk produkt

Ofte vil vi betrakte samling av element der hvert element består av komponenter som kommer fra ulike mengder. For å håndtere dette definerer vi en tuple som en endelig samling av objekt der rekkefølge og antall forekomster har betydning. Et eksempel på en tuple er (a_1, \dots, a_N) . To tupler a og b er like dersom $a_n = b_n$ for $n \in 1, \dots, N$. Mengder av tupler blir ofte konstruert av kartesisk produkt av en mengde mengder. Dette betegnes også som kryssproduktet av mengdene og består av alle mulige tupler der komponent n kommer fra den n 'te mengden.

$$A_1 \times \dots \times A_N = \times_{n=1}^N A_n = \{(a_1, \dots, a_N) | a_n \in A_n \text{ for } n = 1, \dots, N\} \quad (14.14)$$

I praksis bruker vi ofte kryssprodukt av mengder av reelle tall der hvert element er en vektor. Vektorrommet $\mathbb{R}^N = \mathbb{R} \times \dots \times \mathbb{R}$ der $\mathbf{x} = (x_1, \dots, x_N) \in \mathbb{R}^N$. Husker at intervall er viktige delmengder av \mathbb{R} . Dette kan generaliseres til \mathbb{R}^N som rektangler som består av kartesisk produkt av intervall

$$I = \times_{n=1}^N [a_n, b_n) = \{(x_1, \dots, x_N) | a_n \leq x_n < b_n \text{ for } n = 1, \dots, N\} \quad (14.15)$$

14.3 Relasjoner

En relasjon R fra mengden S til mengden T er en delmengde av det kartesiske produktet av mengdene, $R \subset S \times T$. Vi kan bruke en alternativ notasjon for å beskrive medlemskap i relasjonen, $(a, b) \in R \iff aRb$. For relasjonen $<$ innebærer dette at $(a, b) \in < \iff a < b$. Noen navngitte relasjoner:

- Identitetsrelasjonen: $R = \{(x, x) : x \in S\}$
- Tom relasjon: $R = \emptyset$
- Universell relasjon: $R = \{(x, y) : x \in S, y \in T\}$

Noen begreper for å beskrive egenskaper en relasjon kan oppfylle

- Refleksiv: $(x, x) \in R$ for alle $x \in S$
- Symmetrisk: $(x, y) \in R \implies (y, x) \in R$
- Transitiv: $(x, y) \in R \wedge (y, z) \in R \implies (x, z) \in R$
- Antisymmetrisk: $xRy \wedge yRx \implies x = y$
- Irrefleksiv: $(x, x) \notin R, \forall x \in S$

Vi har litt terminologi om relasjoner avhengig av hvilke av egenskapene over de oppfylle. En relasjon er en *ekvivalensrelasjon* dersom den er refleksiv, symmetrisk og transitiv. Den utgjør en *ordning* dersom den er refleksiv, antisymmetrisk og transitiv. Ordningen er total hvis xRy eller yRx for alle $x \in S, y \in T$, slik at alle elementene inngår i relasjonen. Ellers er ordningen partiell.²

14.3.1 Tillukning

Kanskje flytte dette til etter definisjon av operasjon?

- En mengde M er lukket under en operasjon R hvis $xRy \in M$ for alle $x, y \in M$.
- Tillukningen av M under R er den minste mengden som inkluderer ... hmm
- Tillukningen av en relasjon R med hensyn på en egenskap er den minste mengden R' der $R \subset R'$ som oppfylle denne egenskapen.

Ekvivalensmengder

Hvis \sim er en ekvivalensrelasjon kan vi for hver $s \in S$ definere ekvivalensklassen til s som $[s] = \{t \in S : t \sim s\}$. Mengden av ekvivalensklassene kalles kvotientmengden til S under \sim ,

$$S/\sim := \{[s] : s \in S\} \quad (14.16)$$

Kan merke at det utgjør en mengde av delmengder av S . Hvis \sim er identitetsrelasjonen så er $S/\sim = \{\{s\} : s \in S\}$. Vi kan vise at mengden av ekvivalensklasser alltid utgjør

²Jeg vil knytte dette til ordning og ekvivalens på tall for intuisjon, og deretter se hvordan det kan utvides til andre mengder...

en partisjonering av S . Denne partisjoneringen er ofte enklere å jobbe med siden den grupperer elementer som er ekvivalente i henhold til en gitt relasjon slik at vi kan behandle de likt.³ Det er et eksempel på abstraksjon som innebærer en transformasjon til en annen representasjon som lar oss fokusere på det som er vesentlig.

14.3.2 Funksjoner

En binær relasjon fra S til T er en funksjon f dersom hvert $x \in S$ blir assosiert med nøyaktig én $y \in T$. Med andre ord, så vil det for alle $x \in S$ være nøyaktig én $y \in T$ der $(x, y) \in f$. En vanligere notasjon for å beskrive medlemskap i funksjonen er $f(x) = y$, der x er *argumentet* til funksjonen og y er *verdien*. For å henvise til selve funksjonen f som et matematisk objekt liker jeg å bruke

$$\begin{aligned} f : S &\rightarrow T \\ &: x \mapsto f(x) \end{aligned}$$

siden det både viser *definisjonsmengden* og *verdimengden*, samt regelen som tilordner argument til verdi. Vi sier at *bildemengden* av f er $\{f(x) : x \in S\}$ som er alle mulige verdier funksjonen kan ta. Vi sier også at *bildet* av en delmengde \mathcal{X} under f er $\{f(x) : x \in \mathcal{X} \subset S\}$. Vi har litt fleksibilitet i valg av verdimengde så lenge bildemengden er delmengde av denne. Jeg vil nå innføre litt terminologi for å beskrive egenskaper til funksjoner

- Injektiv: funksjonen er *en-til-en* slik at ulikt mapper til ulikt, $x \neq y \implies f(x) \neq f(y)$.
- Surjektiv: funksjonen er *på* verdimengden slik at det tilsvarer bildemengden, $\forall y \in T$ så eksisterer det $x \in S$ slik at $f(x) = y$.
- Bijektiv: både injektiv og surjektiv, altså både på og en-til-en.

Merk at hvis en funksjon er injektiv eksisterer det alltid en invers funksjon $f^{-1} : B \rightarrow A$ der $f^{-1}(b) = a \iff f(a) = b$. Hvis den også er surjektiv vil definisjonsmengden til f^{-1} tilsvare verdimengden til f .

Operasjoner

En funksjon fra $S^n \rightarrow S$ kan betegnes som en operasjon...

³For en dørvakt kan det for eksempel være relevant å dele folk inn etter hvor fulle de er og de er tilstrekkelig gamle...

Funksjoner som matematiske objekt

Kan putte de som element i mengder, definerer norm på de, rangere de, definere operasjon på de som gir ny funksjon.. bruke dette til å beskrive sammensatte funksjoner.. Se litt på analogi til mengde..

14.4 Grafer

Grafer er viktig i matematisk modellering fordi det gir en representasjon som fanger opp essensen av strukturer med relasjoner mellom objektene. En graf G består av en mengde V av *noder* og en mengde E av *kanter* $\{u, v\}$, der $u, v \in V$. Vi sier at to noder er naboer hvis de forbindes med en kant. Vi kan innføre litt mer terminologi:

- En kant som forbinder en node med seg selv kalles en løkke
- To kanter er parallelle dersom de forbinder de samme nodene
- Grafen er enkel dersom den ikke har løkker eller parallelle kanter
- I en retningsgraf er kantene tupler (u, v)
- Grafen er tom hvis $E = \emptyset$
- Grafen er komplett hvis alle nodene er forbundet med alle de andre nodene

Vi kan ha lyst til å bevege oss rundt i grafen og jeg vil innføre litt terminologi for å beskrive dette.

- En *vandring* med lengde n er en sekvens av noder og kanter $(v_0, e_1, v_1, \dots, e_n, v_n)$ der $e_i := \{v_{i-1}, v_i\}$. I enkle grafer er det tilstrekkelig å skrive nodene. Vi kan også forenkle notasjonen slik at vandringen er beskrevet av $v_0 v_1 \dots v_n$.
- Vandringen utgjør en *sti* dersom ingen av nodene blir besøkt mer enn én gang.
- Vandringen er lukket dersom $v_0 = v_n$. Den utgjør en *krets* dersom den er både lukket og en sti.
- Hva er teknisk definisjon på en *sykel*?
- Grafen er *sammenhengende* dersom det er mulig å gjennomføre en vandring mellom v_i og v_j for alle $v_i, v_j \in V$.

14.4.1 Trær

Definisjonen er på grafer er ganske fleksibel. Vi kan påføre mer struktur ved å spesifisere egenskapene til kantene. Vi sier at grafen er et *tre* dersom det er sammenhengende og asyklisk. Terminologi: blad og rot. Har litt egenskaper... blant annet kun én sti mellom to noder.

14.4.2 Vektete grafer

14.5 Induksjon

Induktiv definisjon av mengde: Basismengde og tillukning av basismengde under operasjon. Måte å avgrense uendelig mengder..

14.6 Kombinatorikk

Kombinatorikk kalles gjerne for kunsten å telle. I mange situasjoner kan det være aktuelt å beregne for mange ulike valgmuligheter vi har eller hvor mange ulike måter noe kan gjøres på. Det er blant annet aktuelt i sannsynlighetsregning og for å beregne hvor raskt kompleksiteten til en algorithme vil vokse.

14.6.1 Multiplikasjonsprinsippet

Hvis vi kan betrakte situasjonen som en sekvens av uavhengige valg, i betydningen av at antallet muligheter i hvert valg ikke avhenger av de andre valgene, vil at det totale antallet mulige valg være produktet av antallet muligheter i hvert valg. Det tilsvarer antallet elementer i det kartesiske produktet av valgmulighetene,

$$|A_1 \times \dots \times A_K| = |A_1| \cdot \dots \cdot |A_K| \quad (14.17)$$

Vi kan bruke dette prinsippet til å utlede antallet mulige delmengder av en mengde A , siden vi kan betrakte det som en sekvens av valg der vi i hvert valg bestemmer om vi skal inkludere et element eller ikke. Det samlede antallet valgmuligheter er da $2^{|A|}$.

14.6.2 Permutasjoner

En permutasjon kan betegnes som en tilordning av elementer i en rekkefølge. Ofte kan vi være interessert i antall mulige permutasjoner av elementer i en mengde A . La $|A| = n$. Da er antallet permutasjoner $n \cdot (n - 1) \cdot \dots \cdot 1 := n!$. Vi kan igjen betrakte det som et sekvens av valg og antallet tilsvarer igjen produktet av antall muligheter i hver valg, men merk at valgene ikke er uavhengige og at antallet muligheter blir redusert med én for

hver gang. Hvis vi vil finne alle permutasjoner av elementene i mengden med lengde k , blir det

$${}^nP_k = n \cdot (n-1) \cdot \dots \cdot (n-(k-1)) = \frac{n!}{(n-k)!} \quad (14.18)$$

14.6.3 Kombinasjoner

Andre ganger er vi ikke opptatt av rekkefølge og kun interessert i hvor mange ulike delmengder med gitt kardinalitet vi kan konstruere fra en mengde.

$$\binom{n}{k} := \frac{{}^nP_k}{n!} = \frac{n!}{(n-k)!k!} \quad (14.19)$$

14.7 Informasjonsteori

Hva er informasjon og hvordan kan vi kvantifisere det? Intuitivt så vil en melding inneholde informasjon hvis det reduserer usikkerheten til mottaker. Det kan for eksempel være tilbakemelding om karakter på en eksamen. Vi kan tenke at det inneholder informasjon fordi det avgrenser mulige alternativer.

Vi kan formalisere dette i en modell der informasjon i en melding avhenger av antall mulige meldinger som kunne ha blitt sendt. Mengden informasjon avhenger altså av kontekst og ikke av innhold i den spesifikke meldingen. Hver melding kan representeres med en bitstreng. Vi kan bruke minste antall bits som er nødvendig for å representere alle de mulige meldingene som måleenhet for informasjon. Husk at med k bits kan vi representere 2^k ulike meldinger. For å representere N ulike alternativer trenger vi $2^k = N \iff k = \log_2(N)$ bits.

14.7.1 Entropi

Vi kan betrakte modellen over som et spesialtilfelle der alle de ulike meldingene er like sannsynlige. Modellen kan utvides ved å vekte alternativene med sannsynligheten for at den meldingen blir sendt. Hvis sannsynlighetsmassen er konsentrert på få alternativer er vi i utgangspunktet mindre usikre og det er derfor mindre reduksjon i usikkerhet av melding og dermed inneholder den mindre informasjon. Merk igjen at mengden informasjon ikke avhenger av selve innholdet i melding, så (tror ikke) det er mer informasjon dersom det er melding om lite sannsynlig alternativ.

Vi definerer entropi som

$$H = - \sum p_n \log(p_n) \quad (14.20)$$

der høyere entropi medfører høyere usikkerhet og dermed mer informasjon i melding.

Dersom alle alternativer er like sannsynlig tilsvarer det antall bits i modellen over,

$$H = - \sum \frac{1}{N} \log\left(\frac{1}{N}\right) = -\log\left(\frac{1}{N}\right) = \log(N) \quad (14.21)$$

og hvis vi allerede vet hvilke alternativ det blir så er

$$H = 1 \log(1) + 0 = 0. \quad (14.22)$$

Entropi er altså en egenskap ved sannsynlighetsmassefunksjoner og den er høyere jo mer spredt sannsynlighetsmassen er.