

STATISTICAL MEMOIR

Institute of Mathematics

University of Oslo

No. 2

November 1980

SYNSPUNKTER PÅ TIDSREKKEANALYSE

ved

Erling Sverdrup

Dette stensil inneholder notater som jeg gjorde i forbindelse med et foredrag som jeg holdt i Norsk Matematikkråd våren 1980. Selvsagt ble foredraget mindre omfattende. Fremstillingen bærer preg av å være beregnet på et bredt publikum av matematikere. Det kan være at jeg ikke har funnet det rette balansepunkt.

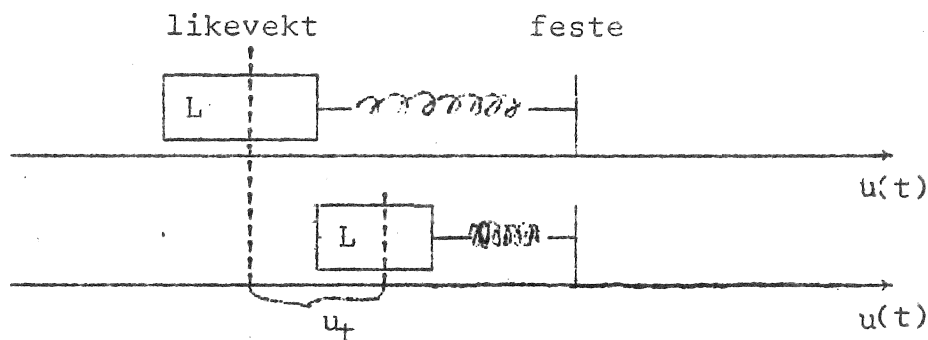
	<u>side</u>
I. Svingninger og treghet i tidsrekker	1
II. Spektrum og autokovarians	10
III. Stasjonære prosesser	13
IV. Statistiske metoder	
a: Prediksjon	23
b: Statistiske analyser	24
c: Periodiske utvalgsundersøkelser	29

I SVINGNINGER OG TREGHET I TIDSREKKER

Arbeider man med tidsrekker vil oppmerksomheten i mange tilfeller bli rettet mot svingefenomener, dvs. mer eller mindre regelmessige periodisiteter. Det gjelder enten man studerer empiriske materialer eller man lager modeller for tidsrekkenes forløp.

La oss se litt nærmere på denslags fenomener. De opptrer i fysikken ved f.eks. pendelsvingninger og elektromagnetiske svingninger; i det økonomiske livs konjunkturer; i geofysikken og i industrielle produksjonsprosesser. Fysikkens svingefenomener er vel de mest klassiske og mange har vel oppfattet f.eks. samfunnsviternes arbeider med fenomenene som et forlorent forsøk på å oppnå like elegante og enkle resultater. Det er de ikke. Jeg tror det kan være viktig å være oppmerksom på både likheter og vesentlige forskjeller mellom de modeller og problemstillinger som fysikere på den ene side og f.eks. geofysikere og økonomer arbeider med. Det vil da være klart at det er vesentlige strukturelle relasjoner i f.eks. det økonomiske liv som gir seg utslag i svingefenomener og at dette er viktig å studere.

La meg først se på situasjonen med en svingende springfjær



Loddet L har masse m . Avviket = avstanden fra likevektstilstanden på tidspunkt t er $u(t)$. Friksjonskoeffisienten er under bevegelsen er f . Vi har tre strukturelle relasjoner

$$\begin{aligned} \text{Tregghetskraften} &= m u''(t) = m \frac{d^2}{dt^2} u(t) \\ \text{Friksjonskraften} &= f u'(t) \end{aligned}$$

Den gjenopprettende kraft = $k u(t)$ (Hookes lov)

Dette er systemets indre krefter.

I tillegg er loddet utsatt for positive eller negative impulskrefter. Impulset på tidspunktet t er $K(t)$. Vi finner

$$K(t) = m u''(t) + f u'(t) + k u(t)$$

Impulset må overvinne tregghet, friksjon og press fra fjæren.

[Anta spesielt at $K(t) = 0$ overalt og at bevegelsen starter på t_0 med avvik $u(t_0)$ og hastighet $u'(t_0)$. Da finner vi såfremt $f^2 < 4 k m$

$$u(t) = \rho e^{-\frac{f}{2m}(t-t_0)} \sin[\lambda(t-t_0)+\phi]$$

hvor

$$\lambda = \sqrt{\frac{k}{m} - \left(\frac{f}{2m}\right)^2}$$

og ρ og ϕ er bestemt slik av initialbetingelsene

$$u(t_0) = \rho \sin \phi$$

$$u'(t_0) = \rho \left[-\frac{f}{2m} \sin \phi + \lambda \cos \phi\right]$$

Anta nå at $K(t)$ har følgende søyleform

$$K(t) = \begin{cases} 0 & \text{for } t < t_1 \\ K_1 & \text{for } t_1 \leq t \leq t_1 + \Delta \\ 0 & \text{for } t > t_1 + \Delta \end{cases}$$

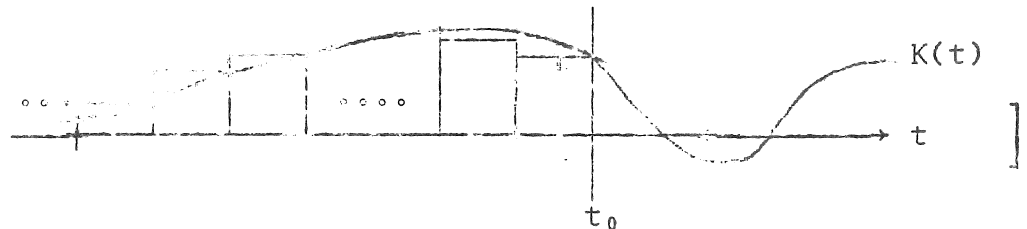
og at systemet er i ro i likevekt før t_1 . Da vil K_1 gi en akselerasjon $a_1 = K_1/m$ i intervallet $(t_1, t_1 + \Delta)$. Men $a_1 = (u'(t_1 + \Delta) - 0)/\Delta$. Herav

$$u'(t_1 + \Delta) = \Delta \cdot K_1/m, \quad u(t_1 + \Delta) = \frac{1}{2} \Delta^2 \cdot K_1/m \approx 0$$

asymptotisk for små Δ . Ved bruk av disse initialbetingelser med $t_0 = t_1 + \Delta$ finner vi $\phi = 0$ og $\rho = \Delta \cdot K_1/m \lambda$ og altså asymptotisk

$$u(t) = \frac{K_1 \Lambda}{m \lambda} e^{-\frac{f}{2m}(t-t_1)} \sin \lambda(t-t_1) ; t > t_1$$

En vilkårlig impulsfunksjon $K(t)$ kan vi tenke oss som en sum av slike søylefunksjoner som vi har betraktet



Vi får som løsning

$$u(t) = \int_{-\infty}^t \frac{K(\tau)}{m \lambda} e^{-\frac{f}{2m}(t-\tau)} \sin \lambda(t-\tau) d\tau$$

$$= \int_0^{\infty} \frac{K(t-\tau)}{m \lambda} e^{-\frac{f\tau}{2m}} \sin \lambda \tau d\tau$$

(Greens formel), såfremt $\lambda = \sqrt{\frac{k}{m} - \left(\frac{f}{2m}\right)^2}$ er reel.

Vi ser at impulset $K(t)$ er gjennom fjæren blitt transformert til en sammensetning av dempete svingninger, alle med samme periode $p = \frac{2\pi}{\lambda}$. Transformasjonen har form av et filter

$$u(t) = \int_{-\infty}^t K(\tau) \psi(t-\tau) d\tau$$

hvor formen av transformasjonen er bestemt av ψ dvs. bare av systemets strukturelle parametre m, f, k .

Hvis $K(t)$ oppfattes som en deterministisk funksjon er det lett å angi forutsetninger under hvilke nevnte løsning eksisterer og er entydig. Vi vil først og fremst tenke på situasjoner hvor $K(t)$ er tilfeldig impulser, dvs. en stokastisk prosess, det medfører at den simultane fordeling for $K(t_1), \dots, K(t_n)$ er gitt for alle n, t_1, \dots, t_n . Ofte er det naturlig å tenke seg at $K(t)$ er stasjonær; dvs. at fordelingen for $K(t+t_1), \dots, K(t+t_n)$ er uavhengig av t for alle n, t_1, \dots, t_n .

Man vil gjerne også at impulsene til forskjellige tidspunkter skal være uavhengige, dvs. $K(t)$ er ren støy. Det er ikke helt enkelt å få til i det tidskontinuerlige tilfelle såfremt $K(t)$ skal være kontinuerlig. Det blir selvmotsigende å forlange kontinuitet og autouavhengighet. Det er imidlertid flere måter å velge $K(t)$ på slik at autouavhengigheten er tilnærmet oppfylt.

Men enten man velger $K(t)$ slik eller slik kan, under generelle forutsetninger, integraluttrykket for $u(t)$ gis en presis mening og være løsningen av differentiaalligningen. Det er bare å oppfatte deriverte og integralet som definert som grense i kvadratisk middel. For integralets vedkommende er det da som grenser for en Riemann-sum. At en stokastisk variabel X_m konvergerer til X i kvadratisk middel betyr at $E(X_m - X)^2 \rightarrow 0$.

La oss vende tilbake til det fysiske hovedpoeng. Svingende springfjær og svingende pendel er noe dagligdags for oss, det måtte vel en betydelig forsker som Huygens til for å oppfatte det som bemerkelsesverdig at dette var resultatet av Gallileis enkle prinsipper om momentum, tregnet osv. som i seg selv ikke sa noe om svingninger.

I det økonomiske liv har man ikke uten videre villet akseptere konjunktorene som noe selvfølgelig. Tvertimot har de blitt sett på som noe gåtefullt, i alle fall opptil 1930-årene. Man søkte mange slags forklaringer, bl.a. i solfleck aktiviteten. I dag tror jeg det er akseptert at det er selve strukturen i vårt økonomiske liv med relasjoner om vare- og pengeetterspørsel, sparing, investering osv. osv. som danner bakgrunn for forståelsen av dem, jfr. de mange intervjuer med de økonomiske vismenn, som ofte uttrykker seg i meget abstrakte ordelag, riktignok med henvisning til konkrete indikatorer ("Situasjonen er så anstrengt at man må vente et omslag").

Det typiske med situasjoner med pendel springfjær, elektromagnetiske svingninger er at man bruker modeller som gir relasjoner mellom posisjon, hastighet, aksellerasjon dvs. mellom posisjon på et tidspunkt og på en eller to tidspunkter like forut. Det er en slags "sluggishness" "treghet" ("seighet"?) i systemet, og det er dette som er hovedpoenget når det gjelder svingefenomenene.

I virkeligheten har man det samme forhold innen geofysikk, industrielle produksjonsprosesser, økonomi, medisinsk-demografiske undersøkelser, osv. Det er en åpenbar treghet i mekanismen. Den kan ikke neglisjeres, hverken ved estimeringer, ved statistiske analyser eller ved styringsproblemer. At denne iakttagelse er triviell er ingen innvending mot å ta hensyn til den. Å ta hensyn til den er sannelig ikke trivielt, som vi skal komme tilbake til.

Ved mange skipperskjønnsmessige statistiske analyser kommer tregheten inn. Valgprognoser er formelt utarbeidet uten å ta hensyn til tregheten, bortsett fra en stratifisering etter resultatet ved foregående valg. Men i siste instans trekkes treghetsbetraktninger inn når politikere og journalister kommenterer dem, uten å ta hensyn til de usikkerhetene som prognoseinstituttene angir.

Ofte finner man den type situasjoner. Den statistiske analyse er todelt. Først kommer den statistiske ekspert som har utledet sin analyseteknikk ut fra, i og for seg, sunne prinsipper fra en matematisk omhyggelig presisert modell, men som unnlater å ta hensyn til visse vesentlige trekk ved situasjonen. Så kommer den politiske økonom med sine relativt løse, men sikkert ofte gode, analyser som tar hensyn til faktisk a priori innsikt. Det burde naturligvis være statistikerens ansvar å trekke inn alle vesentlige omstendigheter i modellen fra første stund og så bygge sine analyser på dem.

For å konkretisere litt med et antydende eksempel, la meg ta sysselsetningen. (Med fare for å demonstrere naivitet og ukyndighet overfor økonometrikere.) Det er en størrelse som voktes av politikere, økonomer og bedriftsledere. Vi vet at det er en målsetting å holde den høy. Tenk bare på hva som skjer når arbeidstakere må sies opp. Straks settes en rekke tiltak i gang. Det begynner vel innen bedriften, rasjonalisering, salgskampanjer, intern refinansiering, så kommer bankene inn i bildet, og de offentlige myndigheter.

Hvis man antar at antall sysselsatte $S(t)$ er stasjonære med $ES(t) = \mu$, kan man kanskje beskrive $S(t)$'s variasjon med

$$mS''(t) + fS'(t) + k(S(t) - \mu) = I(t)$$

hvor altså $I(t)$ er stasjonær og kanskje tilnærmet ukorrelert. I så fall har vi som stasjonær løsning såfremt $0 < f < 2\sqrt{km}$,

$$S(t) = \mu + \frac{1}{m\lambda} \int_{-\infty}^t I(\tau) e^{-\frac{f}{2m}(t-\tau)} \sin \lambda(t-\tau) d\tau,$$

hvor $\lambda = \sqrt{\frac{k}{m} - \left(\frac{f}{2m}\right)^2}$.

En mer detaljert fortolkning av differentiallikningen er ikke så viktig fordi man ikke som ved den svingende springfjær, kan gjøre seg håp om å måle m , f , k ved separate forsøk; dvs. veie loddet (m), måle fjærens stramhet (k), bestemme friksjonen f .

La meg allikevel gjøre et forsøk. Anta at investeringene ("kreftene") på tidspunkt t , kan deles i den del av investeringen som beror på at sysselsetningen faller $\mu - S(t)$ under det normale og den som skyldes andre årsaker.

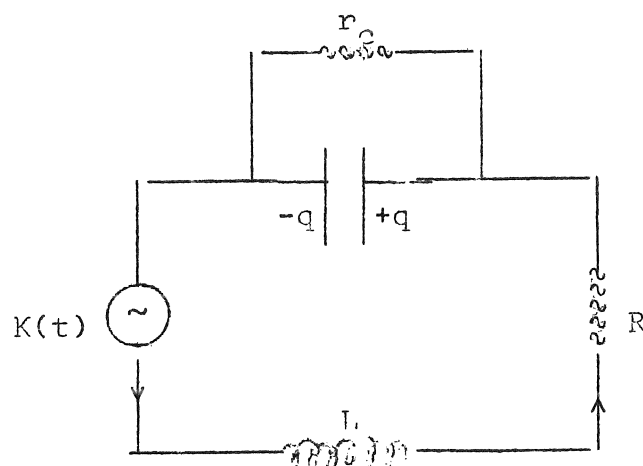
$$\text{Investering} = I(t) + k(\mu - S(t)).$$

Det siste ledd beror altså dels på det offentliges forpliktelse til å holde sysselsetningen høy, dels på private investeringer motivert ved at det er en fordel å investere når sysselsetningen er lav. Disse investeringer lar seg bare gjennomføre ved

endringen i sysselsettingen på tidspunktet t . La oss tenke oss at den momentane endring kan måles ved $S'(t)$ og $S''(t)$ eller summarisk ved $mS''(t) + fS'(t)$, slik at

$$I(t) + k(\mu - S(t)) = fS'(t) + mS''(t).$$

Et forsøk på å gi en dypere mening til de to ledd som friksjon eller inertia blir formodentlig kunstig. Det vesentlige er den a priori tro på sluggishness = treghet, så må konfrontasjonen med virkeligheten, dvs. aktuelle tidsrekkeverdier for $S(t)$ avgjøre realismen, samt muliggjøre estimeringen av k , m , f . (Jeg går da ut fra at $I(t)$ er en ikke-observerbar stokastisk variabel.) Legg forøvrig merke til at hvis variansen σ^2 for $I(t)$ er ukjent, og det vil den naturligvis være, så vil bare estimeringen av m/σ , f/σ , k/σ være mulig, dvs. parametrene m , f , k , σ kan ikke identifiseres utfra fordelingen for observasjonene. For å illustrere at denne mangel på identifiserbarhet ofte går dypere, la meg se på en elektrisk svingekrets bestående av induksjonsspole, kondensator og motstand.



Vi tenker oss at det er en viss strøm-lekasje i isolasjonen mellom platene i kondensatoren C , symbolisert ved en parallelkoblet ledning med motstand r_c .

I kretsen genereres en varierende elektromotorisk kraft $K(t)$ som gir en varierende ladning q på kondensatoren. Anta først ingen lekkasje. $K(t)$ skal overvinne spenningen p.g.a. den ohmske motstand RI , den elektrostatiske spenning over kondensatoren q/c og den selvinduserte motspenning $L\frac{dI}{dt}$ over spolen, hvor $I = \frac{dq}{dt}$,

dvs.

$$K(t) = RI + \frac{q}{c} + L\frac{dI}{dt}$$

dvs.

$$L\frac{d^2q}{dt^2} + R\frac{dq}{dt} + \frac{q}{c} = K$$

eller

$$L\frac{d^2I}{dt^2} + R\frac{dI}{dt} + \frac{1}{c}I = K'$$

som er helt analog til ligninger for svingende fjær. Med en lekkasje med motstand r_c får vi på den annen side

$$L\frac{d^2q}{dt^2} + (R + \frac{L}{cr_c})\frac{dq}{dt} + (\frac{1}{c} + \frac{R}{cr_c})q = K$$

eller

$$L\frac{d^2I}{dt^2} + (R + \frac{L}{cr_c})\frac{dI}{dt} + (\frac{1}{c} + \frac{R}{cr_c})I = K'$$

[Tre fundamentale lover ligger til grunn: 1) Over en platekondensator med to plater med ladning henholdsvis $+q$ og $-q$ coulomb er det en spenning på $U = q/c$ volt, når kapasiteten er c farad. c avhenger bare av kondensatorens fysiske utforming og er uavhengig av q . 2) Over en ledning med strømstyrke I er det en spenning RI , hvor R er den ohmske motstand. 3) Over en spole selvinduseres en motspenning $V = -L\frac{dI}{dt}$ som følge av endringen i strømstyrken I . L er en konstant som bare avhenger av spolens fysiske utformning. Dette er da de tre grunnleggende strukturrelasjoner. - Det følger av 2) at $q/c = ir_c$ hvor i er lekkasjestrømmen, mens strømmen I over spole-motstander må være lik summen av den strøm som går over kondensatoren og den som går som lekkasje-strøm, $I = i + \frac{dq}{dt} = \frac{1}{r_c}\frac{q}{c} + \frac{dq}{dt}$. Innsettes dette i den første ligningen

for $K(t)$ som vi skrev ned fåes nest siste ligning.]

Elektroingeniøren har små vanskeligheter med parametrene L, c, R, r_c . Så vidt jeg vet er motstandere spoler og kondensatorer alminnelig handelsvare som er påført parameterverdiene, dvs. de er målt med stor nøyaktighet ved separate forsøk. Han kan skifte dem ut og lett beregne virkningen av slike utskiftninger. Både styringsproblemet og estimeringsproblemet er greit.

La oss for å nærme oss en situasjon som vi ofte finner utenfor den klassiske fysikk, tenke oss at vi utelukkende var henvist til å observere I -s variasjon over tiden, gitt ved den siste ligning som vi skrev ned. K' tenker vi oss er en ikke observerbar stokastisk variabel som er nær opp til støy. Vi vil da finne estimatorer for

$$L, R + \frac{L}{cr_c}, \frac{1}{c} + \frac{R}{cr_c}$$

relativt til standardavviket σ for $K'(t)$, men herav ikke estimatorer for R, c, r_c . I tillegg til den store usikkerhet ved estimeringen får man et identifikasjonsproblem. L/σ kan studeres, men kapasiteten og de forskjellige Ohmske motstandere er "konfundert", selv relativt til σ . Men det er altså i økonomi, industriell produksjonsprosesser, geofysikk etc. at man vel ofte vil kunne stå overfor denslags problemer.

Situasjonen ved tidsrekkeanalyse er altså følgende,

1. Svingefenomener i tidsrekkene finner man på mange områder utenfor de klassisk-fysiske. Strukturen i tidsrekken kan i noen situasjoner beskrives ved en lineær differential ligning. Det er derfor verdt å beskjefte seg med det. Det er viktig for å oppnå effektive estimeringer og følsomme analysemetoder.

Men det vil ofte være slik at

2. Impulset, dvs. høyreleddet i ligningen vil være ikke observerbar, vi vil oppfatte den som stokastisk og nær opptil autouavhengige eller auto ukorrelert.

3. Koeffisientene i differentiaalligningen kan ikke måles med stor nøyaktighet ved separate eksperimenter, de må estimeres utfra selve tidsrekken.

4. Disse koeffisienter kan være funksjoner av visse bakenforliggende strukturelle parametre som kan brukes til å styre prosessen eller drøfte muligheter for prosessens fremtidige forløp. Ofte vil disse parametre ikke være identifiserbare, dvs. de kan ikke bestemmes entydig utfra ligningens koeffisienter.

Disse omstendighetene gjør ofte den statistiske behandling, dvs. estimeringen og analysen, til et meget vanskelig problem. Vi skal komme tilbake til det.

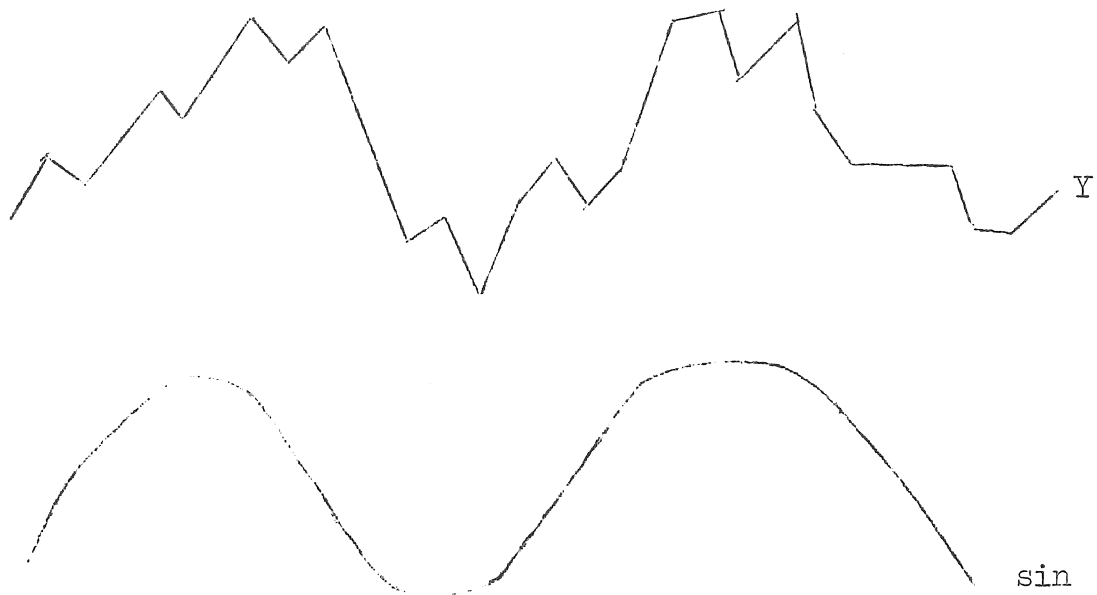
Men først skal vi si noe om det matematiske verktøy som brukes ved studiet av svingefenomener.

II. SPEKTRUM OG AUTOKOVARIANS.

La $Y(1), Y(2), \dots, Y(T)$ være en tidsrekke av observasjoner. For å undersøke om p grovt regnet kan ansees som en periodisitet for tidsrekken la oss sammenligne den med en funksjon

$$\sin(\lambda t + \varphi) = \sin\left(\frac{2\pi}{p}t + \varphi\right); \quad t = 1, 2, \dots, T$$

som har p som periode.



Vi forskyver sinusfunksjonen, dvs. varierer φ slik at den samvarierer best mulig med $Y(t)$, målt ved korrelasjonen eller kovariansen

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T (Y(t) - \bar{Y})(\sin(\lambda t + \varphi) - \overline{\sin(\lambda t + \varphi)}) &= \\ &= \frac{1}{T} \sum_{t=1}^T (Y(t) - \bar{Y})\sin(\lambda t + \varphi) = \frac{1}{T} K(\varphi) \end{aligned}$$

dvs. at vi maksimerer $K(\varphi)$ m.h.p. φ . Med $\hat{\varphi}$ som maksimerende verdi, finner vi da

$$K(\hat{\varphi}) = [(\sum(Y(t) - \bar{Y})\cos \lambda t)^2 + (\sum(Y(t) - \bar{Y})\sin \lambda t)^2]^{\frac{1}{2}}$$

Ved å studere

$$K(\hat{\varphi}) = Q(\lambda) = Q\left(\frac{2\pi}{p}\right)$$

som funksjon av p kan man få et inntrykk av hvilke periodiske komponenter som $Y(t)$ grovt regnet er sammensatt av. Med en periode av lengde p vil det være $\nu = \frac{1}{p} = \frac{\lambda}{2\pi}$ perioder pr.

tidsenhet

$$R^2(\lambda) = \left(\frac{2}{T} Q(\lambda)\right)^2 ; \quad \lambda = \frac{2\pi}{p} = 2\pi\nu$$

som funksjon av p kalles periodegrammet, som funksjon av ν kalles det spektret. Ofte er det forøvrig

$$I_T(\lambda) = \frac{T}{8\pi} R^2(\lambda) = \frac{1}{2\pi T} Q^2(\lambda)$$

som kalles det empiriske spektrum (fordi det har visse egenskaper som estimator). Markerte perioder p i tidsrekken slår ut med en høy verdi av $R^2\left(\frac{2\pi}{p}\right)$ eller $I_T\left(\frac{2\pi}{p}\right)$.

En annen måte å studere forekomsten av perioden i tidsrekken på er å studere autokovariansen

$$C_p = \frac{1}{T-p} \sum_{t=1}^{T-p} (Y(t) - \bar{Y})(Y(t+p) - \bar{Y})$$

som funksjon av p . Den gir uttrykk for korrelasjoner mellom Y -verdier i en tidsavstand p fra hverandre. C_p er åpenbart stor når p er en periode i tidsrekken.

En enkel manipulasjon gir følgende algebraiske sammenheng mellom periodogram og autokovarians

$$R^2(\lambda) = \frac{4}{T} \sum_{h=-(T-1)}^{T-1} \left(1 - \frac{|h|}{T}\right) C_h \cos \lambda h = \frac{4}{T^2} \left| \sum_{t=1}^T (Y(t) - \bar{Y}) e^{it\lambda} \right|^2$$

hvor vi setter $C_p = C_{-p}$ når $p < 0$.

Hvis vi multipliserer med $\cos \lambda k$ og integrerer fra $-\pi$ til π , finner vi

$$\frac{4}{T} \left(1 - \frac{|k|}{T}\right) C_k = \frac{1}{\pi} \int_{-\pi}^{\pi} \cos \lambda k R^2(\lambda) d\lambda ; \quad |k| < T.$$

Det er derfor en en-entydig sammenheng mellom autokovarians og

periodogram, forsåvidt kunne man glemme det ene og konsentrere seg om det andre. Det har imidlertid vist seg bekvemt å operere med begge karakteristikker. Men de er åpenbart rent deskriptive.

III. STASJONÆRE PROSESSER.

For å kunne arbeide rasjonelt med metodeproblemene ved estimeringer og statistiske analyser, trenger man en modell. I første omgang oppfatter man da observasjonene $Y(1), \dots, Y(T)$ som et utsnitt av en tidsfunksjon $Y(t)$ med et eller annet doméne for t . Et rom av funksjons former $Y = \{Y(t)\}_t$ er definert. Man tar utgangspunkt i at for alle $n, t_1, \dots, t_n, y_1, \dots, y_n$ har utsagnet $\prod_{j=1}^n (Y(t_j) \leq y_j)$ en sannsynlighet. Et sannsynlighetsmål P er dermed definert over den minste sigma algebra over disse utsagn. Funksjonene $Y(t) = Y(t; Y)$ er dermed målbare og vi kan definere forventninger og autokovarians

$$\eta(t) = EY(t) = \int Y(t; Y) dP, \quad \sigma(t, s) = \int (Y(t; Y) - \eta(t))(Y(s; Y) - \eta(s)) dP$$

såfremt integralene er endelige. Som tidligere nevnt sier vi at prosessen Y er stasjonær hvis utsagnene

$$\prod_{j=1}^n (Y(t_j) \leq y_j) \quad \text{og} \quad \prod_{j=1}^n (Y(t_j + t) \leq y_j)$$

har samme sannsynlighet for alle n, t, t_1, \dots, t_n , dvs.

$$Y(t_1), \dots, Y(t_n) \quad \text{og} \quad Y(t_1 + t), \dots, Y(t_n + t) \quad \text{har}$$

samme sannsynlighetsfordeling. Stasjonæritet innebærer altså at det bare er den innbyrdes konstellasjon mellom t_1, \dots, t_n som spiller rolle, ikke deres plassering på tidsaksen. Det er ingen

historiske begivenheter. Åpenbart innebærer stasjonæritet at $\eta(t) = \eta$ er uavhengig av t og at $\sigma(t,s) = \sigma(t-s)$ kun avhenger av $|t-s|$, og naturligvis $\text{var} Y_t = \sigma(0)$. Når disse egenskaper er oppfylt, snakker man om "annenordens stasjonæritet".

De har som umiddelbar konsekvens at $\sigma(k)$ er en positiv semidefinit funksjon av k . Hvis nå doméne for t , og dermed for k , er de hele tall $\dots -1, 0, 1, \dots$ så medførte dette (Herglotz) at det eksisterer en ikke avtagende $F(\lambda) \geq 0$, kontinuerlig fra høyre $F(\pi) = \sigma(0)$, $F(-\pi-0) = 0$, slik at

$$\sigma(k) = \int_{-\pi}^{\pi} \cos \lambda k dF(\lambda)$$

Ved en Fourier utvikling av $F(\lambda)$ (samt $\frac{\sigma}{2\pi}(\lambda+\pi)$) finner vi

$$F(\lambda) = \frac{\sigma}{2\pi}(\pi+\lambda) + \frac{1}{\pi} \sum_{k=1}^{\infty} \frac{\sigma(k)}{k} \sin \lambda k.$$

Det kan nå vises at det kumulert av det tidligere omtalte empiriske spektrum $I_T(\lambda)$ er en asymptotisk forventningsrett estimator for $F(\lambda)$ (i kontinuitetspunktene)

$$\lim_{T \rightarrow \infty} E \int_{-\pi}^{\lambda} I_T(v) dv = F(\lambda).$$

Dessuten

$$\lim E C_k = \sigma(k).$$

Dette rettferdiggjør karakteriseringen av prosessen ved hjelp av prosessens spektrum $F(\lambda)$ og autokovarians.

Når doméne for tiden t er hele den reelle tall-linje og $\sigma(k)$ er kontinuerlig, er spektret $F(\lambda)$ definert ved

$$\sigma(k) = \int_{-\infty}^{\infty} e^{ik\lambda} dF(\lambda)$$

i samsvar med Bochner Khintchins som sikrer at en slik $F(\lambda)$ eksisterer siden $\sigma(k)$ er positiv semidefinit.

I konkrete situasjoner vil modellen bak den observerte tidsrekke vanligvis være spesifisert ved noe mer enn stasjonærhets-egenskapen, f.eks. ved de differentiaalligninger som jeg innledningsvis omtalte.

Man kan nå avlede autokovariansfunksjonene og spektrene av modellforutsetningene, i mange tilfeller eksplisitt med relativt enkle analytiske uttrykk. Har man en samling av dem, kan man konfrontere dem med de tilsvarende empiriske funksjoner, altså autokovarianser og spektrere med sikte på å plukke ut en brukbar modell. Slike modelldiskusjoner vil selvfølgelig i høy grad hvile på a priori overveielser. Som på så mange andre områder av statistikken er det lite utviklet teori for denslags modelldiskusjoner.

La meg nå et øyeblikk vende tilbake til de differentiaalligninger som jeg omtalte innledningsvis. Det kan vises at slike lineære differentiaalligninger med konstante koeffisienter og stasjonære høyre-sider $K(t)$, eller mer generelt

$$Y^{(p)}(t) + \beta_1 Y^{(p-1)}(t) + \dots + \beta_p Y(t) = K(t)$$

har en entydig løsning såfremt alle røttene i den tilordnede algebraiske ligning har realdeler $\neq 0$. I så fall kan løsningen gis som en filteroperasjon ψ på K ,

$$Y(t) = \int_{-\infty}^{+\infty} \psi(t-\tau)K(\tau)d\tau$$

hvor integralet eksisterer i kvadratisk middel og ψ er bestemt

av β_1, \dots, β_p alene. Nå er det åpenbart meningsløst å la $Y(t)$ på tidspunktet t være avhengig av de fremtidige impulser $K(\tau)$; $\tau > t$; $Y(t)$ kan ikke være bestemt av hvilke ytre påvirkninger som vil opptre i fremtiden. Derfor må vi ha

$$Y(t) = \int_{-\infty}^t \psi(t-\tau)K(\tau)d\tau$$

dvs. $\psi(s) = 0$ for $s < 0$. Nødvendig og tilstrekkelig betingelse for det viser seg å være at realdelen av alle røttene i den tilordnede algebraiske ligning er < 0 .

Det var den finske aktuar Karhunen, en elev av Cramer, som først klarla disse forhold i et arbeide fra 1946. La meg i den forbindelse komme med noen historiske bemerkninger. Klarleggingen av begrepet stasjonær prosess skyldes så vidt jeg vet Khintchine fra 1934. Han utviklet også de grunnleggende sammenhenger mellom autokorrelasjon og spektrum. Det var klassisk matematisk analyse, velkjent er Bochner-Khintchine's setning. I 1942 kom Cramer's resultat som gikk ut på at en annen ordens stasjonær process $Y(t)$ (med $EY(t) = 0$) kunne fremstilles på formen

$$Y(t) = \int_{-\infty}^{+\infty} e^{it\lambda} d\zeta(\lambda)$$

hvor $\zeta(\lambda)$ opptil en additiv konstant er entydig bestemt ved at tilvekstene $\zeta(\lambda) - \zeta(\lambda')$ og $\zeta(\kappa) - \zeta(\kappa')$ over ikke overlappende intervaller (λ', λ) , (κ', κ) er ukorrelerte. Integralet skal oppfattes som grense i kvadratisk middel for en Riemann-Stieltje-sum. Med $\zeta(-\infty) = 0$, har vi

$$E\zeta(\lambda) = 0, \text{ var}(\zeta(\lambda) - \zeta(\lambda')) = F(\lambda) - F(\lambda')$$

hvor $F(\lambda)$ er det tidligere omtalte kumulerte spektrum. Dette resultat er epokegjørende fordi det beskjeftiger seg med selve sampelstien $Y(t)$, ikke matematiske karakteriseringer av fordelingene for $Y(t)$. Det kaster prinsipielt lys over begrepet tilfeldighet. Det bekrefter det som vi vel alle intuitivt føler at tilfeldighet innebærer, at det er noe ved prosessen som opptrer uavhengig. Det er også viktig for operatorformål, omtrent som en genererende funksjon, Fourier-transform eller Laplace-transform. Således bygger Karhunen's utledning på Cramér's resultat.

Det statistiske innhold av Cramér's resultat er naturligvis hovedsaken. Det er imidlertid interessant å merke seg at det matematiske innhold er et spesialtilfelle av et resultat i funksjonalanalysen som var funnet av Stone allerede i 1932, men som Cramér ikke kjente i 1942.

Siden det er så mange renmatematikere i denne forsamling, antar jeg at det kan være av interesse å se sammenhengen. Av hensyn til dem som, i likhet med meg, ikke til daglig arbeider med funksjonalanalytiske problemer, la meg først rekapitulere noen elementære ting. For visse typer av matriser A er det slik at det eksisterer en matrise C slik at

$$CAC^{-1} = \Lambda \quad \text{og herav} \quad CA^t C^{-1} = \Lambda^t$$

hvor Λ er en diagonal matrise med egenverdiene for A som diagonalelementer. La $\lambda_1, \dots, \lambda_r$ være de r forskjellige egenverdier $r \leq p$, hvis A er $p \times p$. Vi kan da skrive

$$\Lambda = \sum \lambda_j P_j, \quad \text{hvor} \quad \sum P_j = I$$

og hvor alle P_j er diagonale med 1 og 0 i diagonalene på en

ikke overlappende måte. Herav får en vilkårlig kolonne x

$$CAx = \sum \lambda_j P_j Cx.$$

Altså: etter skifte av koordinatsystem er Ax en linearkombinasjon projeksjonene av x på et fullstendig system av ortogonale plan. Oppfatter vi nå A som en transformasjon som har matrisene A eller CA som forskjellige fremstillingsformer, og x som en vektor med kolonnematrisene x og Cx som forskjellige fremstillingsformer, kan vi skrive

$$Ax = \sum \lambda_j P_j x, A^t x = \sum \lambda_j^t P_j x.$$

Lar vi spesielt A være en ortogonal, eller om man for bekvemhets skyld opererer i det komplekse plan, en unitær transformasjon $A = U$, vet vi at alle egenverdier ligger på enhetssirkelen og vi kan skrive $e^{i\lambda_j t}$ istedenfor λ_j , dvs.

$$U^t = \sum_{j=1}^r e^{i\lambda_j t} P_j,$$

hvor

$$\sum P_j = I, \quad P_i P_j = 0, \quad i \neq j$$

Mer generelt i et Hilbertrom har vi

$$U^t = \int_{-\pi}^{+\pi} e^{i\lambda t} dE_\lambda$$

hvor U er unitær, dvs. dens adjungerte U^* er lik dens inverse $U^* = U^{-1}$ og alle E_λ er projeksjoner, dvs. har egenverdier 0 eller 1, og slik at

$$E_{-\pi} = 0, \quad E = I, \quad (E_\lambda - E_\mu)E_\nu = 0; \quad \lambda > \mu > \nu.$$

Det adjungerte A^* av en transformasjon A er definert ved at $(Ax, y) = (x, A^*y)$, med (x, y) som indre produkt. Dette svarer til transponering i tilfelle med reelle matriser og da svarer U til en ortogonal matrise. (Erstatt (x, y) med $y'x$ og (Ax, y) med $y'Ax$, osv.)

Nåvel, dette regnes i dag som et klassisk og sentralt resultat i funksjonalanalysen. Stone's resultat går nå ut på følgende. La U_t , t vilkårlig reell, være en parametrisk gruppe av unitære transformasjoner slik at

$$U_{t+s} = U_t U_s, \quad U_0 = I$$

slik at altså $U_{-t} = U_t^{-1} = U_t^*$. Da har vi at det eksisterer en spektralfamilie $(E_\lambda)_{-\infty < \lambda < \infty}$ slik at

$$U_t = \int_{-\infty}^{\infty} e^{it\lambda} dE_\lambda$$

hvor

$$1) \quad E_{-\infty} = 0, \quad E_\infty = I, \quad E_{\lambda+0} = E_\lambda$$

$$2) \quad (E_\lambda - E_\mu)E_\nu = 0, \quad \lambda > \mu > \nu.$$

En forutsetning man gjør er f.eks. at $(U_t x, y)$ er en kontinuerlig funksjon av t . Meningen med integralet fremgår ved å se på det indre produkt,

$$(U_t x, y) = \int_{-\infty}^{+\infty} e^{it\lambda} d(E_\lambda x, y)$$

hvor $(E_\lambda x, y)$ nå er en kompleks funksjon og integralet er Lebesgue-Stieltjesk.

Anvendbarheten av dette resultat er omfattende etter hva jeg

har forstått. Jeg skal se på hvorledes det kan brukes i forbindelse med stasjonære prosesser. La $Y = \{Y(t)\}_t$ hvor nå $Y(t)$ kan være kompleks. Hver funksjonsverdi $Y(t)$ er trivielt en målbar funksjon av Y

$$Y(t) = Y(t; Y)$$

Vi betrakter mengden av alle slike funksjoner av Y , og utvider den til å omfatte alle lineærkombinasjoner $z(Y) = \sum_{j=1}^m a_j Y(t_j, Y)$ og kompletterer den ved å medta alle grenser for Cauchy-sekvenser av slike lineærkombinasjoner. Med et indreprodukt

$$(z, w) = E z(Y) \overline{w(Y)}$$

har vi et Hilbertrom. Spesielt blir altså

$$(Y(t), Y(s)) = E Y(t) \overline{Y(s)}$$

autokovariansen, som vi altså forutsetter kun avhenger av $t-s$.

Vi definerer nå for enhver reell h transformasjon U_h ved

$$U_h Y(t) = Y(t+h)$$

med utvidelse til en vilkårlig z . Det er åpenbart av stasjonæritetsegenskapen at U_h er isometrisk, dvs. den bevarer det indre produkt. Men den er til og med unitær, fordi

$$(U_h Y(t), Y(s)) = E U_h Y(t) \overline{Y(s)} = E Y(t+h) \overline{Y(s)} = E Y(t) \overline{Y(s-h)}$$

ved stasjonæritetsegenskapen. På den annen side er

$$(U_h Y(t), Y(s)) = (Y(t), U_h^* Y(s)) = E Y(t) \overline{U_h^* Y(s)}$$

ved definisjonsegenskapen ved adjungering.

Herav finner vi da

$$\overline{Y(s-h)} = \overline{U^*Y(s)}$$

dvs.

$$U^*Y(s) = Y(s-h) = U_h^{-1}Y(s)$$

altså $U_h^* = U_h^{-1}$. Altså kan vi bruke Stone's resultat spesielt på $Y(0)$

$$Y(t) = U_t Y(0) = \int_{-\infty}^{+\infty} e^{it\lambda} dE_\lambda Y_0 = \int_{-\infty}^{+\infty} e^{it\lambda} d\zeta(\lambda)$$

hvor vi har innført $\zeta(\lambda) = E_\lambda Y_0$. Men siden $(E_\lambda - E_\mu)E_\nu = 0$ ($\lambda > \mu > \nu$) og $E_\lambda^* = E_\lambda$ finner vi

$$0 = ((E_\lambda - E_\mu)E_\nu z, w) = (E_\nu z, (E_\lambda - E_\mu)w) = E\zeta(\nu)(\overline{\zeta(\lambda) - \zeta(\mu)})$$

altså er ζ en prosess med ukorrelerte tilvekster og vi er tilbake til Cramer's resultat.

La meg nå vende tilbake til et problem jeg nevnte tidligere om impulset som en uavhengig prosess i tilfellet med kontinuerlig tid. En mulighet er å ta utgangspunkt i Wienerprosessen $W(t)$ som er en slags kontinuerlig versjon av tilfeldig gang. Denne prosess har en lang og problemfylt historie. Prosessen må vel tilskrives Einstein. Hans hypotese om de Brownske molekylære bevegelser gikk ut på at variansen for $W(t)$ måtte være $t\sigma^2$. Dette gjorde en ting av Avogadro's tall mulig. Nøyaktigere beregning/ Det ble dessuten forlangt at den skulle være kontinuerlig med uavhengige tilvekster. Mange strevet med å gjennomskue konsekvensene av disse forutsetninger, Wiener, Gnedenko, Doob. Det lyktes tilslutt å bevise at disse forutsetninger måtte medføre at prosessen var normal, dvs. Gaussisk, noe som alle følte måtte være riktig. Resultatet skyldtes formodentlig amerikaneren Donsker.

Av at $W(t) - W(s)$ og $W(s)$; $t > s$; er uavhengige, følger at $W(t) = W(t) - W(s) + W(s)$ og $W(s)$ har kovarians $\text{cov}(W(t), W(s)) = \sigma^2$; $t \geq s$. Vi kunne na velge $K(t) = W(t) - W(t-d)$, men den brutte overgang fra avhengighet til uavhengighet er litt kunstig. Setter man istedet

$$K(t) = e^{-\beta t} W(e^{2\beta t}),$$

får man for $t \geq s$

$$\text{cov}(K(t), K(s)) = e^{-\beta|s-t|} \sigma^2.$$

Her er det et markert fall fra $s = t$ til $s \neq t$. Dessuten er Fourier-transformen = spektraltettheten lik

$$\frac{\sigma^2}{\pi} \frac{2\beta}{\beta^2 + \lambda^2}$$

som for stor β er nesten flat som ved en uavhengig prosess.

For de spesifikke modeller som jeg hittil har omtalt er det så vidt jeg vet, lite utviklet av statistisk inferens teori, dvs. man har ikke utviklet noen estimeringsteori eller statistisk analyseteori for å konfrontere modellene med faktiske observasjoner la oss si $Y(t_1), \dots, Y(t_N)$

Men de tilsvarende diskrete modeller har vært gjenstand for stor oppmerksomhet i de siste 10 år, det er de modeller som fremkommer hvis man tenker seg at prosessen $\{Y(t)\}$ er definert med t heltallig og slik at man i de nevnte differentiaalligninger erstatter $Y'(t)$ med $Y(t) - Y(t-1)$, $Y''(t)$ med $Y(t) - 2Y(t-1) + Y(t-2)$ osv. Da fremkommer en differensligning av typen

$$Y(t) = \varphi_1 Y(t-1) + \varphi_2 Y(t-2) + \dots + \varphi_p Y(t-p) + K(t)$$

hvor $K(t)$ er stasjonær. Det kan vises at denne ligning har en entydig, stasjonær løsning hvis ingen av røttene i den tilordnede algebraiske ligning $1 = \varphi_1 \beta + \dots + \varphi_p \beta^p$ ligger på enhetssirkelen. Denne betingelse er også nødvendig såfremt spektret for $K(t)$ er absolutt kontinuerlig. Løsningen skrives på formen

$$Y(t) = \sum_{j=-\infty}^{\infty} \varkappa(j) \psi(t-j)$$

hvor $\psi(t)$ er bestemt av strukturen, dvs. $\varphi_1, \dots, \varphi_p$. Som nevnt, er det meningsløst å la $Y(t)$ være avhengig av de fremtidige impulser, derfor må vi ha

$$Y(t) = \sum_{j=-\infty}^t K(j) \psi(t-j).$$

Dette vil være tilfelle hvis og bare hvis alle røttene ligger utenfor enhetssirkelen. Det er to spesielle forutsetninger om $K(t)$ som har vært gjenstand for behandling fra et inferens synspunkt. Det ene er tilfellet hvor $K(t)$ -ene er uavhengig. Da en AR-process, kalles $Y(t)$ / det annet er tilfellet hvor $K(t)$ er et glidende gjennomsnitt (moving average), dvs. $= V(t) + \theta_1 V(t-1) + \dots + \theta_q V(t-q)$ hvor $V(t)$ -ene er uavhengige stasjonære. Da kalles $y(t)$ en ARMA-process, (i begge tilfeller er spektret absolutt kontinuerlig).

IV. STATISTISKE METODER.

IV a. Prediksjoner.

De statistiske anvendelser bygger nå på at det i modellen er innebygget en ARMA-prosess

En populær bruk av ARMA prosessen er for rent prediksjonsformål. Man antar at enten Y_t selv eller $Z_t = \Delta^d Y_t$; for $d = 1, 2$

(dvs. $Y_t - Y_{t-1}$ eller $Y_t - 2Y_{t-1} + Y_{t-2}$) følger en stasjonær ARMA prosess. Man antar at prosessen løper fritt uten inngripen og uten strukturelle endringer. Man har observert et utsnitt av den

$$Y_1, \dots, Y_T$$

og ønsker å prediktere Y_{t+h} ; $h > 0$; ved hjelp av en prediktor $\hat{Y}_T(h)$ som er en funksjon av Y_1, \dots, Y_T . Prediktoren finner vi ved å minimere $E(Y_{T+h} - \hat{Y}_T(h))^2$. Med $d = 0$, f.eks. finner vi da rekursivt

$$\hat{Z}_T(h) = \varphi_1 \hat{Z}_T(h-1) + \dots + \varphi_{h-1} \hat{Z}_T(1) + \varphi_h Z_T + \dots + \varphi_p Z_{T+h-p}$$

Usikkerheten ved denne estimeringen for gitte $\varphi_1, \dots, \varphi_p$ kan lett undersøkes, men i tillegg kommer estimeringen av $\varphi_1, \dots, \varphi_p$ på grunnlag av Y_1, \dots, Y_T som resulterer i en usikkerhet som må studeres. Denslags predikteringer for korte perioder fremover har gitt en del gode resultater. Men det er klart at anvendelsesmulighetene for en slik prediksjonsteori er relativt begrenset.

IV b. Statistiske analyser.

Langt interessantere er de statistiske analyser av tidsrekkene som nå er i ferd med å komme, men hvor mange problemer ennå står uløst. Slike analyser kan belyse hvorledes forskjellige faktorer kan innvirke på prosessen; og dermed bringe på det rene viktige strukturelle sammenhenger av interesse for fremtidige disposisjoner.

Vi kan tenke på industrielle produksjonsprosesser hvor det tilføres råstoffer. Råstoffenes kvalitet og blandingsforhold er kanskje ikke under fullstendig kontroll eller kan ikke

observeres direkte, men vi vet at det er en treg forandring over tiden. Det fører til en treg endring også av det ferdige produkt. Samtidig er det faktorer som innvirker på prosessen som man kan observere eller har kontroll over. Man ønsker å studere disse faktorerers virkning. Vi kan tenke på smelteverk, valseverk eller cellulose-produksjon. - Men det behøver ikke nødvendigvis være industrielle produksjonsprosesser. Det kan være en demografisk-medisinsk studie av frekvensen av visse sykdommer og hvorledes disse påvirkes av miljømessige forhold. Infeksjonssykdommer er jo så åpenbart autoregressive. Det kan være luftforurensning og hvorledes denne påvirkes av fyringsoljen i de private hjem og i visse nøkkelbedrifter i området; samt av været.

Vi har altså en prosess $Y(t)$ som gjennomgår trege endringer på grunn av innvirkningen fra visse ikke-observerbare faktorer, derfor antar vi en ARMA prosess. Dessuten har vi visse observerbare faktorer $z_1(t), \dots, z_q(t)$ som påvirker $Y(t)$, men ikke påvirkes av $Y(t)$. Vi velger derfor modellen

$$Y(t) = \varphi_1 Y(t-1) + \dots + \varphi_p Y(t-p) + \gamma_1 z_1(t) + \dots + \gamma_s z_s(t) \\ + V(t) + \theta_1 V(t-1) + \dots + \theta_q V(t-q),$$

hvor $\{V(t)\}$ er en uavhengig prosess med $EV(t) = 0$, $\text{var} V(t) = \sigma^2$. Vi antar alle $z_j(t)$ ikke-stokastiske.

For å gjøre det klart hvorledes $z_j(t)$ gir uttrykk for virkningen, la oss tenke oss at de g første z_j gir uttrykk for virkningen av temperaturen. Vi har g kontrollerbare temperaturnivåer S_1, \dots, S_g , f.eks. i valsemassen i et valseverk. Vi lar $z_i(t) = 1$ hvis temperaturen er S_i , ellers 0. Da er altså virk-

ningen på Y , $\gamma_1, \dots, \gamma_g$ alt ettersom temperaturen er $S_1 < \dots < S_g$. Hvilken temperatur er best? Hvis Y gir uttrykk for kvalitet er vi altså interessert i å komme med utsagn om hvilke temperaturer som gir bedre resultat enn andre i den utstrekning slike utsagn er mulig med den usikkerhet og det begrensede observasjonsmateriale som er tilstede. Konklusjonen kan f.eks. ha formen av å si at S_5, S_6, S_7 gir bedre resultat enn de andre temperaturer, men man tør ikke innlate seg på en innbyrdes rangering av S_5, S_6, S_7 . Hovedpoenget er at man er interessert i om lineærformen $\gamma_i - \gamma_j$ er > 0 . Er man interessert i om det er en opptrappet virkning av temperaturøkning, vil man være interessert i om $\gamma_i - 2\gamma_{i+1} + \gamma_{i+2} > 0$, såfremt S_i -ene er ekvidistante. Ved siden av temperaturen kan man være interessert i en annen faktor, f.eks. justering av valseene i h posisjoner P_1, \dots, P_h . La oss da for bekvemhets skyld bruke doble fotskrifter, altså skrive $\gamma_{ij}z_{ij}(t)$. La $z_{ij}(t)$ være 1 hvis temperaturen er S_i og valseposisjonen P_j , ellers 0. γ_{ij} gir uttrykk for bidrag til kvalitet ved kombinasjonen (S_i, P_j) . Man kan nå ved siden av egenvirkningene av temperatur S og posisjon P være interessert i samspillet mellom temperatur og posisjon. Hvis posisjon P_2 er gunstigere enn P_1 , $\gamma_{i2} - \gamma_{i1} > 0$ for alle S_i , vil da denne tendens gjøre seg sterkere gjeldende ved lavere enn ved høyere temperatur? Er f.eks. $\gamma_{12} - \gamma_{11} > \gamma_{22} - \gamma_{21}$? Man er interessert i utsagn av typen $\sum_{ij} \gamma_{ij} > 0$. Vi er rett og slett interessert i å se på tallene for å finne ut noe interessant som vi tør påstå er reelle virkninger.

Vi vender tilbake til den generelle modell og er altså interessert i en konklusjon som går ut på å peke ut "kontraster"

$\sum_{i=1}^r f_i \gamma_i$ som > 0 . Dette er hva jeg mener med statistisk analyse.

Jeg vil nå gå ut fra at vi har en ren AR-process, dvs.

$$\theta_1 = \dots = \theta_q = 0.$$

Spørsmålet er nå: Med de gitte observasjoner $Y(1), \dots, Y(T)$, hvilke kontraster $\sum_{i=1}^r f_i \gamma_i$ skal vi erklære > 0 ? Statistikerens oppgave er å utvikle en metode som for alle $Y(1), \dots, Y(T)$ peker ut de kontraster som vi tør påstå er > 0 . Metoden må være slik at det er liten sannsynlighet for å komme med gale utsagn og stor sannsynlighet for å komme med riktige utsagn. Det første innebærer at

$$Q = \Pr \left[\bigcup_{f: \sum \gamma_i f_i < 0} (Y(1), \dots, Y(T) \mid \text{det påstås } \sum \gamma_i f_i > 0) \right]$$

skal være liten, vi forlanger

$$Q \leq \epsilon = \text{f.eks. } 0.05 \text{ eller } 0.01.$$

La oss se litt på konstruksjonen av en metode som er justert på denne måten.

Vi finner først minste kvadrater estimatorer

$$\hat{\varphi}_1, \dots, \hat{\varphi}_p, \hat{\gamma}_1, \dots, \hat{\gamma}_s \text{ ved å minimere}$$

$$\sum_{t=p+1}^T [Y(t) - \varphi_1 Y(t-1) - \dots - \varphi_p Y(t-p) - \gamma_1 z_1(t) \dots - \gamma_s z_s(t)]^2$$

mhp. φ og γ . Herav finnes også en estimator

$$\hat{\sigma}^2 = \frac{1}{T-p-s} \sum_t [Y(t) - \hat{\varphi}_1 Y(t-1) - \dots - \hat{\varphi}_p Y(t-p) - \hat{\gamma}_1 z_1(t) - \dots]^2$$

for $\sigma^2 = \text{var}V(t)$. For at ikke det hele skal bli for svevende, la oss si litt om den konkrete numeriske fremgangsmåte. Vi danner oss matrisen med elementer

$$\frac{1}{T} \sum_{t=1}^T z_i(t) z_j(t); \quad i, j = 1, 2 \dots s$$

og så inverterer vi den. Matrisen som består av de r første rekker og kolonner i den inverterte matrise, betegner vi med m . Denne inverterer vi igjen til m^{-1} . Det kan nå vises at

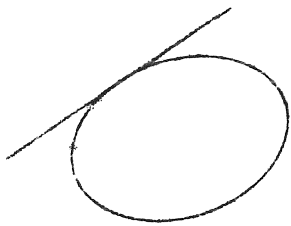
$$W = \sum_{i,j=1}^r (\hat{Y}_i - \gamma_i)(\hat{Y}_j - \gamma_j)(m^{-1})_{ij} / \hat{\sigma}^2$$

asymptotisk χ^2 -kvadrat fordelt med r frihetsgrader når $T \rightarrow \infty$. (Visse forutsetninger om hvorledes z_{it} oppfører seg når $t \rightarrow \infty$ må gjøres). Det betyr da at hvis w er $1-\epsilon$ fraktilen for denne fordeling, så er $\Pr(W > w) = \epsilon$. I det r -dimensionale rom for $\hat{Y}_1, \dots, \hat{Y}_r$, la oss nå se på tangentplanene til den r -dimensjonale ellipsoide $W \leq w$. Det viser seg da at disse plan er gitt ved

$$\sum_{i=1}^r f_i (\hat{Y}_i - \gamma_i) = [w \sum_{i,j} f_i f_j m_{ij}]^{\frac{1}{2}} = K_f$$

Nå er det lett å se at punktet $(\hat{Y}_1, \dots, \hat{Y}_r)$ ligger utenfor ellipsoiden hvis og bare hvis det eksisterer et f slik at

$$\sum_{i=1}^r f_i (\hat{Y}_i - \gamma_i) > K_f$$



La oss nå se på kontrastene $\sum f_i \gamma_i$. Det er naturlig å se på deres estimatorer $\sum f_i \hat{Y}_i$ og påstå at $\sum f_i \gamma_i > 0$ hvis $\sum f_i \hat{Y}_i$ er stor. La oss spesielt velge K_f -ene som kritiske punkter.

Av hva vi har sagt følger at

$$W > w \iff \bigcup_f \left(\sum_{i=1}^r f_i (\hat{Y}_i - \gamma_i) > K_f \right).$$

Herav

$$\Pr\left(\bigcup_f \left(\sum_{i=1}^r f_i (\hat{Y}_i - \gamma_i) > K_f\right)\right) = \epsilon,$$

Herav følger da spesielt at hvis $z_1(t), \dots, z_r(t)$ er uten noensomhelst virkning på $Y(t)$, dvs. $\gamma_1 = \dots = \gamma_r = 0$, da vil sjansen for å begå den bommert å påstå minst en "interessant" kontrast være ϵ for enhver $\gamma_{r+1}, \dots, \gamma_r, \varphi_1, \dots, \varphi_p, \sigma^2$. Men vi er ikke interessert i en såkalt null-hypotese $\gamma_1 = \dots = \gamma_r = 0$, vi er interessert i sjansen for å begå feil for en vilkårlig $\gamma_1, \dots, \gamma_r, \dots, \gamma_s, \sigma^2, \varphi_1, \dots, \varphi_p$. Nå kan vi imidlertid bevise at sjansen for å begå en feil er høyst ϵ , asymptotisk.

For en vilkårlig $(\gamma_1, \dots, \gamma_r)$ har vi nemlig

$$\begin{aligned} Q &= \Pr(\text{begå en feil}) = \Pr\left(\bigcup_{f: \sum_i f_i \gamma_i \leq 0} \left(\sum_{i=1}^r f_i \hat{Y}_i \geq K_f\right)\right) = \\ &= \Pr\left(\bigcup_{\sum_i f_i \gamma_i \leq 0} \left(\sum_{i=1}^r f_i (\hat{Y}_i - \gamma_i) + \sum_{i=1}^r f_i \gamma_i \geq K_f\right)\right) \leq \\ &\leq \Pr\left(\bigcup_{\sum_i f_i \gamma_i \leq 0} \left(\sum_{i=1}^r f_i (\hat{Y}_i - \gamma_i)\right)\right) \end{aligned}$$

siden leddet $\sum_i f_i \gamma_i$ i ulikheten er ≤ 0 . Men tar vi nå unionen over alle f , får vi naturligvis at det siste uttrykk er

$$\leq \Pr\left(\bigcup_f \left(\sum_{i=1}^r f_i (\hat{Y}_i - \gamma_i)\right)\right) = \epsilon$$

slik at vi har bevist at $Q = \Pr(\text{begå feil}) \leq \epsilon$.

IV c. Periodiske utvalgsundersøkelser.

Vi nevnte tidligere at en stadig tilbakevendende utvalgstelling er et eksempel på en situasjon hvor statistikerne ikke effektivt tar hensyn til den faktiske a priori innsikt ved konstruksjon av den statistiske metode. Slike utvalgstellinger ut-

føres av Statistisk Sentralbyrå for registrering av husdyrhold osv. i landbruket, for studium av antall som sysselsettes eller blir arbeidsledige osv. osv. Meningsmålingsinstituttene undersøkelser av partitilhørighet er også et eksempel. For illustrasjonens skyld, la oss se på sistnevnte situasjon; uten at det som sies bokstavelig må oppfattes som en anbefaling; til det trengs en grundigere utredning.

La ξ_t betegne antall høyrestemmer på tidspunktene $t = 1, 2, \dots, n$. Det er altså disse tall vi er interessert i. På tidspunktet t spørres et tilfeldig utvalg av stemmeberettigede. Det relative antall høyrefolk blant disse finnes og multipliseres med antall stemmeberettigede i hele befolkningen. Derved fremkommer et foreløpig estimat X_t for ξ_t .

Med gitte ξ_t kan sannsynlighetsfordelingen for (X_1, \dots, X_n) finnes. X -ene vil være stokastisk mer eller mindre avhengige eller uavhengige alt ettersom det foretas en hel eller delvis utskiftning fra utvalg til utvalg. Fordelingen for X_1, \dots, X_n vil vi nå oppfatte som en betinget fordeling gitt ξ_1, \dots, ξ_t , siden vi skal oppfatte ξ_t som stokastisk. Åpenbart er det en viss treghet i endringen i den sanne ξ_t ; mange høyrefolk er og forblir høyrefolk. Vi går derfor ut fra at $\xi_t - \xi_{t-1}$ eller $\xi_t - 2\xi_{t-1} + \xi_{t-2}$ tilfredsstiller en ARMA prosess som vi har beskrevet, med normalt fordelte impulser.

La nå V_t være differensen mellom det estimerte X_t og det sanne uobserverte ξ_t . Vi vet at ved moderat store utvalg kan V_t antas normalt fordelt for gitt ξ_t . Anta, litt forenklet, at V_t og ξ_t er uavhengige. Numeriske studier med simulerte tall tyder på at dette kan være en brukbar tilnærming. Dermed kan den

simultane fordeling for $(X_1, \dots, X_n, \xi_1, \dots, \xi_n)$ utledes. La oss nå se på estimeringen av antall høyrestemmer ξ_n på siste tidspunkt. Det er nå klart at hele tidsrekken X_1, \dots, X_n ikke bare X_n vil belyse hvor stor ξ_n er. Som estimator for ξ_n vil vi med kjente $\varphi_1, \dots, \varphi_p, \theta_1, \dots, \theta_q$ bruke $\hat{\xi}_n = E(\xi_n | X_1, \dots, X_n)$, som kan finnes ut fra simultanfordelingen for (X_1, \dots, ξ_n) . Nå vil denne $\hat{\xi}_n$ være avhengig av parametrene $\varphi_1, \dots, \varphi_p, \theta_1, \dots, \theta_q$ som i virkeligheten er ukjente. Disse kan imidlertid estimeres og innsettes i $\hat{\xi}_n$, som resulterer i en estimator $\hat{\hat{\xi}}_n$. Estimeringen av $\varphi_1, \dots, \theta_q$ foregår f.eks. ved å finne simultanfordelingen for X_1, \dots, X_n og så finne sannsynlighetsmaksimerings-estimatorer. Ut fra denne fordelingen kan så fordelingen for $\hat{\xi}_n$ finnes og brukes til å vurdere metodens usikkerhet. Vel, så enkelt kan metoden skisseres. Den praktiske gjennomføring vil være ganske formidabel, men det er ting det arbeides med i Statistisk Sentralbyrå.