

Lad  $\omega$  være en amerikaner. Definer

$F_i(\omega)$  = værdien af faktoren  $i$ 'te faktor for  $\omega$

$D(\omega)$  = dødsårsag for  $\omega$

$L(\omega)$  = dødsalder for  $\omega$

I det kommende undertrykkes afhængigheden af  $\omega$ .

## Faktorsandsynligheder

Vi har blandt andet brug for størrelserne

$P(F_i \in f)$ , for forskellige mængder  $f \in \mathcal{F}_i$

$P(F_{i_1} \in f_1, F_{i_2} \in f_2, \dots, F_{i_n} \in f_n)$  for  $f_i, i = 1, \dots, n$

hvor  $f_i$  er mængder der giver mening for deres tilhørende faktor. I filerne i mappen Factor\_frequencies ligger filer af typen.

$i_1$ 'te faktor	$\dots$	$i_n$ 'te faktor	freq
$f_1^0$	$\dots$	$f_n^0$	$P(F_{i_1} \in f_1^0, F_{i_2} \in f_2^0, \dots, F_{i_n} \in f_n^0)$
$f_1^1$	$\dots$	$f_n^1$	$P(F_{i_1} \in f_1^1, F_{i_2} \in f_2^1, \dots, F_{i_n} \in f_n^1)$
$\vdots$	$\dots$	$\vdots$	$\vdots$
$f_1^{n_k}$	$\dots$	$f_n^{n_k}$	$P(F_{i_1} \in f_1^{n_k}, F_{i_2} \in f_2^{n_k}, \dots, F_{i_n} \in f_n^{n_k})$
$\vdots$	$\dots$	$\vdots$	$\vdots$
$f_1^{n_1}$	$\dots$	$f_n^{n_1}$	$P(F_{i_1} \in f_1^{n_1}, F_{i_2} \in f_2^{n_1}, \dots, F_{i_n} \in f_n^{n_1})$

Alternativt kan man skrive  $F = (F_{i_1}, \dots, F_{i_n})$

## Incidents

Dernæst har vi sandsynlighederne for at dø af en dødsårsag i løbet af et år.

$$p_d(l) = P(D = d, L \leq l \mid L > l - 1), d \in \mathcal{D}, l \in \mathbb{R}_+$$

hvor  $\mathcal{D}$  er en mængde af alle dødsårsager i programmet. For beregning, kender vi  $p_d(l)$  som stykvis konstant funktion på mængderne

$$l_1, l_2, \dots, l_{22} = [0, 1), [1, 5), [5, 10), \dots, [95, 100), [100, \infty)$$

Vi er interesserede i vektoren

$$p_d(l_1) \quad p_d(l_2) \quad \dots \quad p_d(l_{22})$$

hvor vi, med den lidt misbrugte notation  $p_d(l_i)$ , mener  $p_d(l)$  for et  $l \in l_i$ . Disse estimeres med

$$p_d(l_i) \leftarrow \frac{\text{antal amerikanere døde af } d \text{ i aldersgruppen } l_i \text{ i år } Y_1}{\text{antal amerikanere i aldersgruppen } l_i \text{ i år } Y_2} =: \frac{a_{di}}{a_i}$$

(Lige nu har vi  $Y_1 = 2014, Y_2 = 2013$ ). Filerne af formen *ICDcode.txt* indikerer et  $d$  med deres titel og indholdet er

$$a_{d1} \quad a_{d2} \quad \cdots \quad a_{d22}$$

Og filen *population.txt* indeholder

$$a_1 \quad a_2 \quad \cdots \quad a_{22}$$

## Risk ratios

Lad  $F_1, \dots, F_k$  være nogle faktorer. Risk ratios i dette program fortolkes som

$$\text{RR}_d(f) := \frac{P(D = d, L \leq l \mid L > l - 1, (F_1, \dots, F_k) \in f)}{P(D = d, L \leq l \mid L > l - 1, (F_1, \dots, F_k) \in f_0)}$$

Egentlig skulle der et  $f_0$  på i notationen for  $\text{RR}_d(f)$ , for at indikere at det er riskratio med hensyn til baselinen  $f_0$ . Det er en antagelse, at riskratioen ikke afhænger af  $l$ . Risk ratioerne er kendt for en mængde af faktorinddelinger  $f \in \mathcal{F}$ . Vi kræver mere eller mindre at  $\mathcal{F}$  kan skrives på formen

$$\mathcal{F} = \mathcal{F}_1 \times \mathcal{F}_2 \times \cdots \mathcal{F}_k$$

hvor

$$\mathcal{F}_i = \{f_1^0, f_1^1, \dots, f_1^{n_i}\}$$

og så er de kodet ved hjælp af

$i_1$ 'te faktor	$\cdots$	$i_k$ 'te faktor	RR
$f_1^0$	$\cdots$	$f_k^0$	$\text{RR}_d((f_1^0, \dots, f_k^0))$
$f_1^1$	$\cdots$	$f_k^1$	$\text{RR}_d((f_1^1, \dots, f_k^1))$
$\vdots$	$\cdots$	$\vdots$	$\vdots$
$f_1^{n_1}$	$\cdots$	$f_k^{n_k}$	$\text{RR}_d((f_1^{n_1}, \dots, f_k^{n_k}))$
$\vdots$	$\cdots$	$\vdots$	$\vdots$
$f_1^{n_1}$	$\cdots$	$f_k^{n_k}$	$\text{RR}_d((f_1^{n_1}, \dots, f_k^{n_k}))$

## 0.1 Udregning

I første omgang vil vi gerne udregne

$$P(D = d, L \leq l \mid L > l - 1, F = f), \text{ for } f \in \mathcal{F} \quad (1)$$

hvor  $F$  er en vektor af faktorer og  $\mathcal{F}$  er den endelige mængde af faktorsammensætninger. (1) kan senere (i javascript-delen) bruges til at udregne

$$P(D = d, L \leq l \mid L > l - 1, F = f_{\text{personal}}) \quad (2)$$

hvor  $f_{\text{personal}}$  ikke (nødvendigvis) ligger i  $\mathcal{F}$ . Det gøres ved

$$(1) = \frac{\text{RR}_d(f)}{\sum_{f \in \mathcal{F}} \text{RR}_d(f) \cdot P(F = f)} p_d(l) \quad (3)$$

Der er dog nogle forhindringer før vi bare kan stoppe tallene fra filerne ind i (3).

- Vi ikke har nok information til at kende  $P(F = f), f \in \mathcal{F}$  100%.
- Vi har flere riskratio filer for samme  $d$ .
- Hvordan man skal tage højde for alders specifikke riskratios og alders specifikke  $P(F = f)$ 'er

Lad os tage dem i voksende sværhedsgrad

### Flere risk ratio filer

Hvis vi der er to riskratiofiler baseret på to vektorer af faktorer  $F^1, F^2$  og to tilhørende krydsmængder,  $\mathcal{F}^1$  og  $\mathcal{F}^2$ , så ville den mest rigtige måde at kombinere dem på være

$$\begin{aligned} P(D = d, L \leq l \mid L > l - 1, (F^1, F^2) = (f^1, f^2)) \\ = \frac{g(\text{RR}_d(f^1), \text{RR}_d(f^2))}{\sum_{f^1, f^2 \in \mathcal{F}^1 \times \mathcal{F}^2} g(\text{RR}_d(f^1), \text{RR}_d(f^2)) \cdot P((F^1, F^2) = (f^1, f^2))} p_d(l) \end{aligned} \quad (4)$$

hvor  $g$  er en passende interaktionsfunktion. Hvis  $F^1$  og  $F^2$  er uafhængige og  $g(x, y) = x \cdot y$ , kan man dog skrive det som

$$\begin{aligned} P(D = d, L \leq l \mid L > l - 1, (F^1, F^2) = (f^1, f^2)) \\ = \frac{\text{RR}_d(f^1)}{\sum_{f^1 \in \mathcal{F}^1} \text{RR}_d(f^1) \cdot P(F^1 = f^1)} \frac{\text{RR}_d(f^2)}{\sum_{f^2 \in \mathcal{F}^2} \text{RR}_d(f^2) \cdot P(F^2 = f^2)} p_d(l) \end{aligned} \quad (5)$$

Fordelen ved (5) er, at man kan lægge et tallene

$$\frac{\text{RR}_d(f^1)}{\sum_{f^1 \in \mathcal{F}^1} \text{RR}_d(f^1) \cdot P(F^1 = f^1)}, f^1 \in \mathcal{F}^1$$

i en fil og tallene

$$\frac{\text{RR}_d(f^2)}{\sum_{f^2 \in \mathcal{F}^2} \text{RR}_d(f^2) \cdot P(F^2 = f^2)}, f^2 \in \mathcal{F}^2$$

i en anden fil og tallene

$$p_d(l_i), i = 1, \dots, 22$$

i en tredje fil. Selvom betingelserne for at bruge (5) ikke er helt opfyldt, kan det måske alligevel være en god ide at bruge den.

## Vi kender ikke $P(F = f), f \in \mathcal{F}$

Her tænkes  $\mathcal{F}$  som den mængde hvor vi kender  $\text{RR}_d(f)$  hvis og kun hvis  $f \in \mathcal{F}$ .

Der er flere slags udfordringer her.

1. Det hænder at vi kun kender

$$P(F = f), f \in \mathcal{F}'$$

hvor  $\mathcal{F}' \not\subseteq \mathcal{F}$ .

2. Vi har  $\mathcal{F} = \mathcal{F}_1 \times \mathcal{F}_2$  og  $F = (F_1, F_2)$ , hvor  $F_1$  og  $F_2$  er hver deres faktor. Her hænder det at vi kun kender

$$P(F_i = f_i), f_i \in \mathcal{F}_i$$

for  $i = 1, 2$  og altså ingenting om den simultane fordeling af  $(F_1, F_2)$

3. Vi har  $\mathcal{F} = \mathcal{F}_1 \times \mathcal{F}_2 \times \mathcal{F}_3$  og  $F = (F_1, F_2, F_3)$ . Det hænder, at vi kun kender

$$\begin{aligned} P((F_1, F_2) = (f_1, f_2)), f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2 \\ P((F_1, F_3) = (f_1, f_3)), f_1 \in \mathcal{F}_1, f_3 \in \mathcal{F}_3 \end{aligned}$$

4. Vi har  $\mathcal{F} = \mathcal{F}_1$  og  $F = F_1$ . Det hænder, at vi kender

$$P((F_1, F_2) = (f_1, f_2)), f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2$$

men ikke (umiddelbart)

$$P(F_1 = f_1), f_1 \in \mathcal{F}_1$$

så vi så om sige har for meget

5. Vi har  $\mathcal{F} = \mathcal{F}_1$  og  $F = F_1$ . Det hænder, at vi slet ikke kender noget til

$$P(F_1 \in f)$$

Problemerne findes også i flere flere dimensioner og med kombinationer.

### En løsning af problem 4

Vi laver en funktion som marginaliserer, dvs

$$P(F_1 = f_1) = \sum_{f_2 \in \mathcal{F}_2} P((F_1, F_2) = (f_1, f_2))$$

### En løsning af problem 1

Hvis vi laver en funktion, som laver transformationen

$$\begin{aligned} w : \{P(F = f_1), P(F = f_2), \dots, P(F = f_n)\} \\ \mapsto \{P(F = f'_1), P(F = f'_2), \dots, P(F = f'_k)\} \end{aligned}$$

kan vi løse det første problem. Hvis vi antager

$$\bigcup_{i=1}^n f_i \subseteq \bigcup_{j=1}^k f'_j \tag{6}$$

kan man lave løsningen

$$\begin{aligned} P(F = f'_j) &= w(P(F = f_1), P(F = f_2), \dots, P(F = f_n)) \\ &= \sum_{i=1}^n \frac{P(F = f_i) \cdot |f_i \cap f'_j|}{|f_i|} \end{aligned}$$

hvor  $|\cdot|$  repræsenterer et mål. Det vil dog nok altid være muligt at bruge Lebesguemålet eller tælleområdet og nogle gange kan man måske være nødt til at bruge en mikstur af de to. Det ses foreksempel ved rygning, hvor der er en kategori, der hedder 0 cigaretter. Betingelsen (6) er nødvendig for at  $P(F = f'_j), j = 1, \dots, k$  summer til 1 (det er nemlig antaget at  $P(F = f_i), i = 1, \dots, n$  summer til 1).

### En løsning af problem 2,3 og 5

For at løse problem 2,3 og 5 kan man bruge tilpasning af marginaler. Man har den ønskede fordeling

$$P(F = f), f \in \mathcal{F}$$

hvor  $F = (F_1, \dots, F_k)$  og  $\mathcal{F} = \mathcal{F}_1 \times \dots \times \mathcal{F}_k$ , og man kender

$$\begin{aligned} P(F^1 = f^1), f^1 \in \mathcal{F}^1 \\ \vdots \\ P(F^m = f^m), f^m \in \mathcal{F}^m \end{aligned}$$

hvor  $\mathcal{F}^l = \mathcal{F}_{i_1^l} \times \mathcal{F}_{i_2^l} \times \dots \times \mathcal{F}_{i_{r_l}^l}$ . Man starter med en standard uniformfordeling

$$p(f) = \frac{1}{\#\mathcal{F}}, f \in \mathcal{F}$$

som estimer for  $P(F = f)$ . Dernæst *tilpasser man med marginalen*  $\mathcal{F}^1$

$$p^{ny}(f) = p(f) \frac{P(F^1 = f^1(f))}{\sum_{f': f^1(f') = f^1(f)} p(f')}$$

og derefter med marginalen  $\mathcal{F}^2, \mathcal{F}^3$  og så videre (indtil man når til hvad?).

### Forskellige credibilities

Vi har mængder,  $f \in \mathcal{F}$ , for hvilke vi vil finde  $P(F = f)$ .  $F$  kan her skrives  $(F_i)_{i \in I}$ , hvor  $I$  er en mængde i  $\{1, \dots, n\}$ . Vi kender da 'binnede' fordelinger af  $(F_i)_{i \in I_j}$  for  $j = 1, \dots, k$ , hvor  $I_j$  også er mængder i  $\{1, \dots, n\}$ . Credibilityscoren er en funktion  $c : \{I_j\}_{j=1, \dots, k} \rightarrow \mathbb{R}_+$ . De binnede fordelinger, der indgår i konstruktionen  $P(F = f)$  er

$$\left\{ (F_i)_{i \in I_j} \mid I_j \cap I \neq \emptyset \wedge (\nexists k : I \cap I_j \subseteq I \cap I_k \wedge c(I_k) > c(I_j)) \right\}$$

•

### Aldersspecifikke $P(F = f)$ 'er eller riskratioer

I princippet burde alle ovenstående udregninger laves separat for alle aldersgrupper. Det gøres også, og der er nogle genveje. Definer  $A$  til at være

den stokastiske variabel der angiver alderen på personen. De nye faktorsandsynlighedsfiler har specificerende kolonner ved

$i_1$ 'te faktor	$\dots$	$i_n$ 'te faktor	aldersgr.
$f_1^0$	$\dots$	$f_n^0$	$A_1$
$f_1^0$	$\dots$	$f_n^1$	$A_1$
$\vdots$	$\dots$	$\vdots$	
$f_1^0$	$\dots$	$f_n^{n_k}$	$A_1$
$\vdots$	$\dots$	$\vdots$	
$f_1^{n_1}$	$\dots$	$f_n^{n_k}$	$A_h$

og freq kolonnen indeholder

$$\begin{aligned}
 &\text{freq} \\
 &P(F_{i_1} \in f_1^0, F_{i_2} \in f_2^0, \dots, F_{i_n} \in f_n^0 \mid A \in A_1) \\
 &P(F_{i_1} \in f_1^0, F_{i_2} \in f_2^0, \dots, F_{i_n} \in f_n^1 \mid A \in A_1) \\
 &\vdots \\
 &P(F_{i_1} \in f_1^0, F_{i_2} \in f_2^0, \dots, F_{i_n} \in f_n^{n_k} \mid A \in A_1) \\
 &\vdots \\
 &P(F_{i_1} \in f_1^{n_1}, F_{i_2} \in f_2^{n_2}, \dots, F_{i_n} \in f_n^{n_k} \mid A \in A_h)
 \end{aligned}$$

Det vil sige, at freq-kolonnen skal summe til  $h$ . Når man konstruerer  $P(F = f)$ , bør man konstruere

## 0.2 Javascript delen

Til hver cause hører noget data på formen

$$d, (p_d(l_i))_{i=1, \dots, 22} \text{ dataframe, risk ratio data}_d$$

hvor  $d$  bare er en string,  $(p_d(l_i))_{i=1, \dots, 22}$  er en data frame udskrevet som liste og risk ratio data'et har formen

$$(\text{Risk ratio datagruppe})_1^d, (\text{Risk ratio datagruppe})_2^d, \dots, (\text{Risk ratio datagruppe})_{n_d}^d$$

En risk ratio datagruppe - indekseret ved  $(j, d)$  - består af

$$(\text{norm}^{j,d}(l_i))_{i=1, \dots, 22}, (\text{risk ratio dataframe}_i^{j,d})_{i=1, \dots, k_{j,d}}, g_{j,d}$$

hvor  $(\text{norm}^j(l_i))_{i=1,\dots,22}$  er en liste af 22 tal som normaliserer for hver af de 22 aldersgrupper, list of risk ratio dataframes er en liste af risk ratio dataframes udskrevet som lister og  $g_{j,d}$  er en string, som siger hvilken interaktion der mellem  $(j, d)$  dataframesne. En risk ratio dataframe for en faktor  $F^{i,j,d}$  har formen

$$(r^{i,j,d}(f))_{f \in \mathcal{F}^{i,j,d}}$$

hvor  $\mathcal{F}^{i,j,d}$  er en endelig mængde af faktor levels. Dette  $r$  er blot en diskretisering af den underliggende riskfaktor funktion

$$R^{i,j,d}(\psi), \psi \in \Psi^{i,j,d}$$

hvor  $\Psi^{i,j,d}$  er en mængde af alle tænkelige værdier af faktoren  $F^{i,j,d}$ , og derfor kan den være uendelig. Vi vil lave en *polating* funktion,  $\text{pol}$ , til at evaluere funktionen

$$R^{i,j,d}(\psi) = \text{pol}((r^{i,j,d}(f))_{f \in \mathcal{F}^{i,j,d}}, \psi) \quad (7)$$

Lad nu  $\psi$  være alle en persons faktorværdier, og lad  $\phi^{i,j,d}$  være vektoren af faktorer der er relevante for den  $(i, j, d)$ 'te risk ratio fil, dvs. en delvektor af  $\psi$ . Lad  $a$  være en vilkårlig alder og lad  $l(a)$  være den af de 22 alderskategorier, som  $a$  falder i. Så defineres

$$P_d(a, \psi) = P_d(a) \cdot \prod_{j=1}^{n_d} \frac{g_{j,d}(R^{1,j,d}(\psi^{1,j,d}), \dots, R^{k_{j,d},j,d}(\psi^{k_{j,d},j,d}))}{\text{norm}^{j,d}(l(a))} \quad (8)$$

hvor  $P_d(a)$  er 'afdiskretiseringen' af  $p_d(l)$ . Det kunne måske defineres som

$$P_d(a) = \max(1, \text{pol}((p_d(l_i))_{i=1,\dots,22}, a))$$

Størrelsen  $P_d(a, \psi)$  tænkes at være det samme som (2), dvs.

$$P_d(a, \psi) = P(D = d, L \leq a \mid L > a - 1, F = \psi) \quad (9)$$

## Kombinationer af $P_d(a)$ 'er

Vi definerer nu

$$p(a, \psi) = \sum_{d \in \mathcal{D}} p_d(a, \psi)$$

hvor  $\mathcal{D}$  er mængden af alle dødsårsager. Hvis (9) faktisk gælder, er

$$p(a, \psi) = P(L \leq a \mid L > a - 1, F = \psi)$$



### Liste over mest interessante kombinationer

Lad nu  $\tilde{\mathcal{D}} \subseteq \mathcal{D}$  være en vilkårlig delmængde af dødsårsagerne. Alle de følgende størrelser kan varieres ved at tilføje  $d \in \tilde{\mathcal{D}}$  og/eller  $L \in [a, b]$  til betingningen

- Forventet levealder

$$E[L \mid \psi] = \sum_{a=0}^{\infty} a \cdot p(a, \psi) \prod_{b=0}^{a-1} (1 - p(b, \psi))$$

- Forventet antal år mistet til sygdom  $d$

$$E[L \mid \psi, D \in \mathcal{D} \setminus \{d\}] - E[L \mid \psi]$$

- Sandsynligheden for at dø af  $d$

$$P(D = d \mid \psi) = \sum_{a=0}^{\infty} p_d(a, \psi) \prod_{b=0}^{a-1} (1 - p(b, \psi))$$

- Ens overlevelseskurve

$$(P(L \geq a \mid \psi))_{a \in \mathbb{N}_0} = \left( \prod_{b=0}^{a-1} (1 - p(b, \psi)) \right)_{a \in \mathbb{N}_0}$$

### Forklare døden ved hjælp af ens faktorer

Vi vil nu også lave en optimal  $\psi$ -værdi som man kan måle brugerens  $\psi$  op imod. Dette er dog en udfordring, fordi en faktor som forårsager en sygdom kan hæmme fremkomsten af en anden sygdom. Det gælder foreksempel rygning, lungekræft og Parkinson's. Vi definerer derfor  $\psi_0$  så  $\psi_0^{i,j,d}$  minimerer  $R^{i,j,d}(\cdot)$ . Det vil sige at i en lungekræftssammenhæng har  $\psi_0$  rygning på 0, mens i en Parkinson's sammenhæng har  $\psi_0$  rygning på 2 (cig/day). Antag af  $\psi$  er gjort op af  $n$  faktorer. Vi definerer nu for et  $J \subseteq \{1, \dots, n\}$  faktorværdivektoren  $\psi_{0,F_J}$ . For en indang  $i$  er

$$(\psi_{0,F_J})_i = (\psi_{0,\{F_j\}_{j \in J}})_i = \begin{cases} (\psi)_i & \text{hvis } i \in J \\ \left( \arg \min_{\tilde{\psi}: \tilde{\psi}_j = (\psi)_j \text{ for } j \in J} R(\tilde{\psi}) \right)_i & \text{ellers} \end{cases} \quad (10)$$

Så vi har  $\psi_\emptyset = \psi_0$  og  $\psi_{0,F_{\{1, \dots, n\}}} = \psi$ . Betragt nu følgende dekomposition

$$\begin{aligned}
P(D = d \mid F = \psi) &= \sum_{a=0}^{\infty} p_d(a, \psi) \prod_{b=0}^{a-1} (1 - p(b, \psi)) \\
&= \sum_{a=0}^{\infty} (p_d(a, \psi) - p_d(a, \psi_0)) \prod_{b=0}^{a-1} (1 - p(b, \psi)) \\
&\quad + \sum_{a=0}^{\infty} p_d(a, \psi_0) \prod_{b=0}^{a-1} (1 - p(b, \psi))
\end{aligned}$$

Definitionen af  $\psi_0$  gør begge de to led positive. Vi fortolker det første led som sandsynligheden for død på grund af ens faktorerwaardier, mens det andet led er *death by chance/age/destiny*. Bemærk at andet led ikke er identisk med  $P(D = d \mid F = \psi_{0,F})$  fordi parameteren  $\psi$  stadig indgår i leddet. Vi er nu interesserede i at dekomponere

$$p_d(a, \psi) - p_d(a, \psi_0)$$

for da kan vi dekomponere  $P(D = d \mid F = \psi)$  i flere led. Fra nu af holder vi  $a$  og  $d$  fast og ser på ovenstående størrelse og definerer  $U_* = p_d(a, \psi) - p_d(a, \psi_0)$ .

Hvis  $F$  bare er en enkelt faktor dvs.,  $F = (F_1)$ , er den fuldt dekomponeret. Men hvis  $F$  er en vektor  $F = (F_1, F_2, \dots, F_n)$  ville vi ideelt have positive tal

$$S(F_J) = S(\{F_j\}_{j \in J}), J \subseteq \{1, \dots, n\}$$

sådan at

$$U_* = \sum_{J \subseteq \{1, \dots, n\}, J \neq \emptyset} S(F_J) \quad (11)$$

og hver af  $S(F_J)$ 'erne havde en pæn fortolkning; gerne sådan at

$$p_d(a, \psi_{0,F_I}) - p_d(a, \psi_0) = \sum_{J \subseteq I, J \neq \emptyset} S(F_J) \quad (12)$$

for alle  $I \subseteq \{1, \dots, n\}$ . Definer for simplicitet  $U(F_I) = p_d(a, \psi_{0,F_I}) - p_d(a, \psi_0)$ .

## Overvejelse

Betragt de simple riskratiointeraktioner

For at undgå komplicerede subskripts definerer vi  $S(F_{i_1}, \dots, F_{i_k}) = S(F_{\{i_1, \dots, i_k\}})$ . Ved at bruge (12) på Tabel 0.1 får man

	$f_1^1$	$f_1^2$
$f_2^1$	1.0	1.1
$f_2^2$	1.1	2.1

Tabel 0.1: Her er der en positiv interaktion mellem  $F_1$  og  $F_2$ 

	$f_1^1$	$f_1^2$
$f_2^1$	1.0	2.0
$f_2^2$	2.0	2.1

Tabel 0.2: Her er der en negativ interaktion mellem  $F_1$  og  $F_2$ 

$$\begin{aligned}\tilde{S}(F_1) &= 0.1 \\ \tilde{S}(F_2) &= 0.1 \\ \stackrel{(12)}{\Rightarrow} \tilde{S}(F_1, F_2) &= 0.9\end{aligned}$$

hvor  $\tilde{S} = S/p_d(a, \psi_0)$ . (12) virker ikke lige så god på Tabel 0.2 for da er

$$\begin{aligned}S(F_1) &= 1.0 \\ S(F_2) &= 1.0 \\ \stackrel{(12)}{\Rightarrow} S(F_1, F_2) &= -0.9\end{aligned}$$

Og den negative værdi besværliggør fortolkningen, så det er måske problematisk at bruge (12).

### Løsningsforslag

Man kan godt få lidt mening ud af de negative værdi der kan forekomme under systemet i (12). Systemet medfører at

$$S(F_I) = U(F_I) - \sum_{J \subseteq I, J \neq \emptyset, J \neq I} S(F_J) \quad (13)$$

hvilket ikke ikke leder til en cirkeldefinition, fordi alle led i summen har færre elementer end antallet af elementer i  $I$ . For at få mening ud af  $S$ 'erne foreslås følgende fortolkninger/omregninger, som inddeles i forskellige kompleksitet

- (i) 1. ordensfortolkning,  $S^{(1)}$ . Dette er den simpleste fortolkning for her undgår vi at have interaktioner mellem faktorårsagerne. Det kan opnås ved at normalisere  $S(F_i)$ 'erne

$$S(F_i) = U(F_i)$$

$$S^{(1)}(F_i) = \frac{S(F_i)}{\sum_{j=1}^n S(F_j)} \cdot U_*$$

for  $i = 1, \dots, n$ . Man kunne sige  $S^{(1)}(F_J) = 0$  for  $|J| > 1$ . Faktoren  $\frac{U_*}{\sum_{j=1}^n S(F_j)}$  ganges på for at  $S^{(1)}(F_i)$ 'erne fungerer som en fuld dekomposition.

- (ii) Næst kigger vi på interaktioner mellem par af faktorårsagerne, der udregnes som

$$S(F_i, F_j) = U(F_i, F_j) - S(F_i) - S(F_j).$$

Problemet er her, at det er muligt at  $S(F_i, F_j) < 0$ . Vi ved dog at

$$U(F_i, F_j) > \max[S(F_i), S(F_j)] \quad (14)$$

og dermed at hvis  $S(F_i, F_j) < 0$  så er

$$|S(F_i, F_j)| < S(F_i) \quad (15)$$

Et negativt  $S(F_i, F_j)$  betyder at niveauerne af  $F_i$  og  $F_j$  ikke forstærker hinanden. Tværtimod, de dækker over de samme sygdomstilfælde. Antag fx

$$D = \{\text{drukne i havnen}\}$$

$$F_1 = \{\text{genstande drukket om ugen}\}$$

$$F_2 = \{\text{max genstande drukket på en enkelt dag om ugen}\}$$

Folk der drikker meget har større sandsynlighed for at falde i havnen og drukne fordi de ikke kan komme op. Derfor forventer vi at højere værdier af  $F_1$  og  $F_2$  øger sandsynligheden for  $D$ . Vi forventer at man har en endnu højere risiko for at drukne i havnen hvis man gør begge dele, men vi forventer ikke at den er meget højere end hvis man kun gjorde en af tingene. Vi kan derfor godt være i situationen som i Tabel 0.2. Man kan forstå det som at nogle af de dødsfald som ville være sket på grund af  $F_1$ , stadig ville være sket på grund af  $F_2$  - og omvendt.

Derfor foreslås en ny slags faktorårsag; et ELLER led, noteret med  $T(F_1, F_2)$ . For  $n = 2$  ville vi definere det således:

$$T(F_1, F_2) = \begin{cases} |S(F_1, F_2)| & \text{hvis } S(F_1, F_2) < 0 \\ 0 & \text{ellers} \end{cases} \quad (16)$$

$$\tilde{S}(F_1, F_2) = \begin{cases} 0 & \text{hvis } S(F_1, F_2) < 0 \\ S(F_1, F_2) & \text{ellers} \end{cases} \quad (17)$$

$$\tilde{S}(F_i) = S(F_i) - T(F_i, F_j), \quad \text{for } i, j = 1, 2, j \neq i \quad (18)$$

$$C = \frac{U_*}{\sum_{i \neq j}^n \tilde{S}(F_i, F_j) + T(F_i, F_j) + \sum_i \tilde{S}(F_i)} \quad (19)$$

$$S^{(2)}(F_i) = C\tilde{S}(F_i) \quad (20)$$

$$S^{(2)}(F_i, F_j) = C\tilde{S}(F_i, F_j) \quad (21)$$

$$T^{(2)}(F_i, F_j) = CT(F_i, F_j). \quad (22)$$

Men der er et problem; hvis  $n > 2$  er der mange forskellige par af faktorer og så ved man ikke hvilket  $j$  man skal vælge i (18). For at undgå de problemer kan vi udvide ELLER-leddene til at inkludere så mange led at det ikke bliver et problem. Dertil definerer vi for hvert  $i = 1, \dots, n$  mængden af inhiberende faktorer af anden orden

$$\mathcal{H}_i^2 = \left\{ j \in \{1, \dots, n\} : \exists k, i_1, \dots, i_k \in \{1, \dots, n\} \text{ hvor} \right. \\ \left. S(F_j, F_{i_1}) < 0, S(F_{i_1}, F_{i_2}) < 0, \dots, \right. \\ \left. S(F_{i_{k-1}}, F_{i_k}), < 0, S(F_{i_k}, i) < 0 \right\}.$$

Alle  $S$ 'er mellem faktorer inden for  $\mathcal{H}_i^2$  og faktorer uden for  $\mathcal{H}_i^2$ , er positive, hvilket vil sige at der er ingen ELLER-led at tage højde for. Definer nu  $V(F_J) = U(F_J) - \sum_{j \in J} S(F_j)$ . Man derfor vælge faktorårsagsstørrelserne som

$$T(F_{\mathcal{H}_i^2}) = \begin{cases} |V(F_{\mathcal{H}_i^2})| & \text{hvis } V(F_{\mathcal{H}_i^2}) < 0, |\mathcal{H}_i^2| > 1 \\ 0 & \text{ellers} \end{cases} \quad (23)$$

$$\tilde{S}(F_i, F_j) = \begin{cases} 0 & \text{hvis } i \in \mathcal{H}_j^2 \\ S(F_i, F_j) & \text{ellers} \end{cases} \quad (24)$$

$$\tilde{S}(F_{\mathcal{H}_i^2}) = \begin{cases} V(F_{\mathcal{H}_i^2}) & \text{hvis } V(F_{\mathcal{H}_i^2}) > 0, |\mathcal{H}_i^2| > 2 \\ 0 & \text{ellers} \end{cases} \quad (25)$$

Størrelserne  $T(F_{\mathcal{H}_i^2})$  fanger nu den evt. negative interaktion der er mellem grupper af faktorer. Ud fra definitionen af  $\mathcal{H}_i^2$  er vi dog ikke sikre på at  $V(F_{\mathcal{H}_i^2}) < 0$ , og i det tilfælde tager vi den med som en positiv interaktion i  $\tilde{S}$ . Da vi tog højde for den negative interaktion i (18) havde vi uligheden i (15) til at sikre os at  $\tilde{S}(F_i)$  blev positiv. Det har vi ikke længere, hvorfor det foreslås at vi gør det følgende:

a. Hvis

$$T(F_{\mathcal{H}_i^2}) \cdot \frac{2}{|\mathcal{H}_i^2|} \leq \min_{j \in \mathcal{H}_i^2} S(F_j) \quad (26)$$

sættes

$$\tilde{S}(F_i) = S(F_i) - T(F_{\mathcal{H}_i^2}) \cdot \frac{2}{|\mathcal{H}_i^2|} \quad (27)$$

Det vil sige at den negative interaktion mellem variablene  $F_{\mathcal{H}_i^2}$  ikke er større end at den kan repræsenteres med en bar der kan tages uniformt fra alle faktorer i  $F_{\mathcal{H}_i^2}$ .

b. Hvis

$$T(F_{\mathcal{H}_i^2}) \cdot \frac{2}{|\mathcal{H}_i^2|} > \min_{j \in \mathcal{H}_i^2} S(F_j) \text{ og} \quad (28)$$

$$2 \cdot T(F_{\mathcal{H}_i^2}) \leq \sum_{j \in \mathcal{H}_i^2} S(F_j) \quad (29)$$

definerer vi flere, nye  $T$ -led der kan “tages” uniformt fra deres marginale  $S(F_i)$ 'er. Lad  $S^{(k)}$  være den  $k$ 'te mindste værdi fra mængden  $\{S(F_i)\}_{i \in \mathcal{H}_i^2}$  og lad  $h^{(k)}$  være det tilsvarende faktorindeks. Dvs.  $S(F_{h^{(k)}}) = S^{(k)}$ . Så defineres

$$\tilde{T}(F_{\mathcal{H}_i^2}) = S^{(1)} \cdot \frac{|\mathcal{H}_i^2|}{2} \quad (30)$$

Derefter defineres

$$\tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}) = \left( S^{(k+1)} - S^{(k)} \right) \frac{|\mathcal{H}_i^2| - k}{2} \quad (31)$$

for  $k = 1, \dots, k_0$ , hvor  $k_0 \equiv k_0(i)$  er det største tal der ville opfylde denne ulighed:

$$\sum_{k=1}^{k_0} \tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}) < T(F_{\mathcal{H}_i^2}) \quad (32)$$

Til sidst defineres

$$\tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k_0+1)}\}}) = T(F_{\mathcal{H}_i^2}) - \sum_{k=1}^{k_0} \tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}) \quad (33)$$

Ideen er, at vi nu har dekomponeret  $T(F_{\mathcal{H}_i^2})$  i flere mindre ELLER-led. På grund (35) ved vi at  $l_0 < k$ . Det kan ske at  $|\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(l)}\}| = 1$  i hvilket tilfælde det sidste led ikke er et rigtigt ELLER-led, men blot et normalt led.

c. Hvis

$$T(F_{\mathcal{H}_i^2}) \cdot \frac{2}{|\mathcal{H}_i^2|} > \min_{j \in \mathcal{H}_i^2} S(F_j) \text{ og} \quad (34)$$

$$2 \cdot T(F_{\mathcal{H}_i^2}) > \sum_{j \in \mathcal{H}_i^2} S(F_j) \quad (35)$$

kan vi ikke længere repræsentere den fulde negative interaktion med ELLER-led. En oplagt måde at imødekemme det, er at ændre  $k_0$  til at blive udregnet fra (32) til at være det største tal der opfylder

$$\sum_{k=1}^{k_0} \tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}) < \frac{\sum_{j \in \mathcal{H}_i^2} S(F_j)}{2} \quad (36)$$

og dermed sætte

$$\tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k_0+1)}\}}) = \frac{\sum_{j \in \mathcal{H}_i^2} S(F_j)}{2} - \sum_{k=1}^{k_0} \tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}) \quad (37)$$

Tilfældet b. kan ses som en generalisering af a. og c. er en generalisering af b. Derfor beskriver jeg kun normaliseringen for c.-tilfældet (hvor  $\frac{\sum_{j \in \mathcal{H}_i^2} S(F_j)}{2}$  erstattes af  $\min\left(\frac{\sum_{j \in \mathcal{H}_i^2} S(F_j)}{2}, T(F_{\mathcal{H}_i^2})\right)$  i (36) og (37)).

Normaliseringen bør være

$$\tilde{S}(F_i) = S(F_i) - \sum_{\substack{k \in \{0, \dots, k_0+1\}: \\ i \in \mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}} \tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}) \frac{2}{|\mathcal{H}_i^2| - k} \quad (38)$$

$$C = \frac{U_*}{\frac{\sum_{i \neq j}^n \tilde{S}(F_i, F_j) + \sum_i \tilde{S}(F_i) + \left( \frac{\sum_{\{H_i^2\}_{i=1, \dots, n}} (\tilde{T}(F_{\mathcal{H}_i^2}) + \tilde{S}(F_{\mathcal{H}_i^2}))}{\sum_{k=1}^{k_0(i)} \tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}})} \right)}{}} \quad (39)$$

$$S^{(2)}(F_i) = C \tilde{S}(F_i) \quad (40)$$

$$S^{(2)}(F_i, F_j) = C \tilde{S}(F_i, F_j) \quad (41)$$

$$T^{(2)}(F_{\mathcal{H}_i^2}) = C \tilde{T}(F_{\mathcal{H}_i^2}) \quad (42)$$

$$S^{(2)}(F_{\mathcal{H}_i^2}) = C \tilde{T}(F_{\mathcal{H}_i^2}), \quad |\mathcal{H}_i^2| > 2 \quad (43)$$

$$(44)$$

og

$$T^{(2)}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}) = C \tilde{T}(F_{\mathcal{H}_i^2 \setminus \{h^{(1)}, \dots, h^{(k)}\}}) \quad (45)$$

for  $k = 1, \dots, k_0(i)$ . Se evt. nedenstående eksempel

### Eksempel

Betragt følgende Risk-ratio tabel for en sygdom  $d$  og en alder  $A$ :

	$F_1 = a$		$F_1 = b$	
	$F_2 = a$	$F_2 = b$	$F_2 = a$	$F_2 = b$
$F_3 = a$	1.0	2.0	1.8	1.8
$F_3 = b$	1.5	1.9	1.9	2.1

Antag  $\psi = (b, b, b)$ . Vi har tydeligvis  $\psi_0 = (a, a, a)$ . Vi vil nu dekomponere

$$p_d(A, \psi) - p_d(A, \psi_0) \quad (46)$$

som er det samme som at dekomponere

$$\frac{p_d(A, \psi) - p_d(A, \psi_0)}{p_d(A, \psi_0)} = R(\psi) - R(\psi_0) \quad (47)$$



hvor  $R$  er ovenstående risk-ratio-tabel. Fremover antager vi derfor uden tab af generalitet at sandsynlighederne  $S, T, U, V$ , osv. er riskratios. Fra (13) får vi

$$S(F_1) = R(\psi_{0,F_1}) - R(\psi_0) = R(b, a, a) - R(a, a, a) = 1.8 - 1.0 = 0.8$$

$$S(F_2) = R(\psi_{0,F_2}) - R(\psi_0) = R(a, b, b) - R(a, a, a) = 1.9 - 1.0 = 0.9$$

$$S(F_3) = R(\psi_{0,F_3}) - R(\psi_0) = R(a, a, b) - R(a, a, a) = 1.5 - 1.0 = 0.5$$

Vi har også

$$U_* = R(\psi) - R(\psi_0) = R(b, b, b) - R(a, a, a) = 2.1 - 1.0 = 1.1$$

Derfor er

$$C = \frac{U_*}{\sum_{i=1}^2 S(F_i)} = \frac{1.1}{2.2} = 0.5$$

Dermed får vi

$$S^{(1)}(F_1) = 0.8 \cdot 0.5 = 0.4$$

$$S^{(2)}(F_2) = 0.9 \cdot 0.5 = 0.45$$

$$S^{(3)}(F_3) = 0.5 \cdot 0.5 = 0.25$$

Når vi husker at  $R(a, a, a) = 1.0$  får vi sandsynligheden at dø af

$$\text{Ukendt grund: } \frac{1.0}{1.0 + 0.4 + 0.45 + 0.25} \approx 48\%$$

$$F_1: \frac{0.4}{1.0 + 0.4 + 0.45 + 0.25} \approx 19\%$$

$$F_2: \frac{0.45}{1.0 + 0.4 + 0.45 + 0.25} \approx 21\%$$

$$F_3: \frac{0.25}{1.0 + 0.4 + 0.45 + 0.25} \approx 12\%$$

- For at udregne 2. ordens interaktioner starter vi igen med (13)

$$\begin{aligned}
S(F_1, F_2) &= R(\psi_{0, F_{\{1,2\}}}) - R(\psi_0) - 2.2 \\
&= R(b, b, a) - R(a, a, a) - 2.2 \\
&= 1.8 - 1.0 - 2.2 \\
&= -1.4 \\
S(F_2, F_3) &= R(a, b, b) - R(a, a, a) - 2.2 \\
&= -1.3 \\
S(F_1, F_3) &= R(b, a, b) - R(a, a, a) - 2.2 \\
&= -1.3
\end{aligned}$$

Altså er

$$\mathcal{H}_1^2 = \mathcal{H}_2^2 = \mathcal{H}_3^2 = \{1, 2, 3\}$$

Så udregner vi  $T(F_{\mathcal{H}_i^2})$

$$\begin{aligned}
R(b, b, b) - R(a, a, a) - \sum_{i=1}^3 S(F_i) &= 2.1 - 1.0 - 2.2 \\
&= -1.1
\end{aligned}$$

Da  $-1.1 < 0$  sætter vi  $T(F_{\mathcal{H}_i^2}) = |-1.1| = 1.1$  og da  $1.1 > \min_{j \in \mathcal{H}_1^2} S(F_j) = S(F_3) = 0.5$  bruger vi formel (30) til at udregne  $\tilde{T}(F_{\mathcal{H}_i^2})$

$$\begin{aligned}
\tilde{T}(F_{\mathcal{H}_i^2}) &= S^{(1)} \cdot \frac{|\mathcal{H}_i^2|}{2} \\
&= 0.5 \cdot \frac{3}{2} \\
&= 0.75
\end{aligned}$$

Nu tjekker vi om vi er færdige med at udregne  $\tilde{T}$ -led ved at undersøge kriterium (32), dvs. sammenligne summen af  $\tilde{T}$ 'erne og  $T(F_{\mathcal{H}_i^2})$ . Vi har  $0.75 < 1.1$ , så vi udregner endnu et  $\tilde{T}$ -led vha. (31).

$$\begin{aligned}
\tilde{T}(F_{\mathcal{H}_i^2 \setminus \{3\}}) &= \left( S^{(2)} - S^{(1)} \right) \frac{|\mathcal{H}_1^2| - 1}{2} \\
&= 0.8 - 0.5 \\
&= 0.3
\end{aligned}$$

Da  $0.3 + 0.75 < 1.1$ , fortsætter vi selvom det sidste  $\tilde{T}$ -led kun indeholder en faktor og derfor ikke er et almindeligt ELLER-led.

$$\begin{aligned}\tilde{T}(F_{\mathcal{H}_i^2 \setminus \{3,1\}}) &= \left( S^{(3)} - S^{(2)} \right) \frac{|\mathcal{H}_1^2| - 2}{2} \\ &= (0.9 - 0.8) \cdot \frac{1}{2} \\ &= 0.05\end{aligned}$$

Vi har  $0.05 + 0.75 + 0.35 \not\prec T(F_{\mathcal{H}_i^2})$ , så vi kan godt stoppe nu og behøver ikke lave nogen korrektion fra c. tilfældet fordi  $0.05 + 0.75 + 0.35 = T(F_{\mathcal{H}_i^2})$ . Korrektion fra c. ville have været at omdefinere  $\tilde{T}(F_{\mathcal{H}_i^2 \setminus \{3,1\}})$  til

$$\frac{\sum_{j \in \mathcal{H}_i^2} S(F_j)}{2} - \tilde{T}(F_{\mathcal{H}_i^2 \setminus \{3\}}) - \tilde{T}(F_{\mathcal{H}_i^2})$$

hvilket også havde givet 0.05. Til sidst, normaliseringerne

$$\begin{aligned}\tilde{S}(F_1) &= S(F_1) - \tilde{T}(F_{\mathcal{H}_1^2}) \frac{2}{3} - \tilde{T}(F_{\mathcal{H}_1^2 \setminus \{3\}}) \\ &= 0.8 - 0.5 - 0.3 \\ &= 0\end{aligned}$$

og

$$\begin{aligned}\tilde{S}(F_2) &= S(F_2) - \tilde{T}(F_{\mathcal{H}_2^2}) \frac{2}{3} - \tilde{T}(F_{\mathcal{H}_2^2 \setminus \{3\}}) - 2\tilde{T}(F_{\mathcal{H}_2^2 \setminus \{3,1\}}) \\ &= 0.9 - 0.5 - 0.3 - 0.1 \\ &= 0.0\end{aligned}$$

og

$$\begin{aligned}\tilde{S}(F_3) &= S(F_3) - \tilde{T}(F_{\mathcal{H}_2^2}) \frac{2}{3} \\ &= 0.5 - 0.5 \\ &= 0\end{aligned}$$

Det vil sige vi har følgende fortolkning af sandsynligheden for at dø af

$$\text{Ukendt grund: } \frac{1.0}{2.1} \approx 48\% \quad (48)$$

$$F_1, F_2 \text{ eller } F_3: \frac{0.75}{2.1} \approx 36\% \quad (49)$$

$$F_1 \text{ eller } F_2: \frac{0.3}{2.1} \approx 14\% \quad (50)$$

$$F_2: = \frac{0.05}{2.1} \approx 2\% \quad (51)$$

## Forslag til polerende funktion

For hver dødsårsag,  $d$ , er der en række riskratiofiler (indekseret med  $j$ ). I formel (7) diskuterede vi udregningen af en interpolerende funktion

$$R^{i,j,d}(\psi) = \text{pol}((r^{i,j,d}(f))_{f \in \mathcal{F}^{i,j,d}}, \psi).$$

Her er  $i$  en (vistnok overflødig) indeksering der angiver aldersgruppen.  $\mathcal{F}^{i,j,d}$  er de binnede faktorer i risk ratio-filen og  $\psi$  er et niveau af faktorer (som ikke er binnede). For at gøre det efterfølgende mere simpelt bruger vi notationen

$$R(\psi) = \text{pol}(r, \psi)$$

hvor vi altså tænker at  $r$  må afhænge af  $\mathcal{F}^{i,j,d}$  som vi omdøber til  $\mathcal{F}$ . For hvert niveau i  $\mathcal{F}$  udregner vi et midtpunkt for de kontinuerte variable,  $\text{mid} : \mathcal{F} \mapsto \mathcal{S}$ . Her skal også tages nogle valg. Hvis  $f \in \mathcal{F}$  er på formen

$$\{0\} \times [15, 32] \times [\text{“Yes”}] \times [4, \infty)$$

er det tyer det meget naturligt at midtpunkterne for de tre første faktorer er 0, 16 og “Yes” hhv. Det sidste interval er her lidt tricky da det ikke har noget midtpunkt. Jeg ser to valg

- (A) Vælg tallet 4
- (B) Kig på det foregående interval, som måske er på formen  $[1, 4)$ , noter at dets midtpunkt er 1.5 væk fra 4, og sig at midtpunktet for  $[4, \infty)$  bør være  $4 + 1.5 = 5.5$

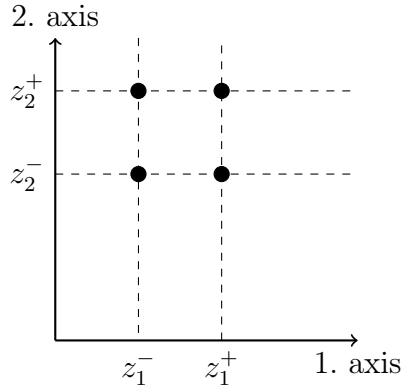
Generelt har vi kun brug for at interpolere og finde midtpunkter, når faktorniveauerne er intervaller. Det antages derfor i det efterfølgende. Når vi har valgt midtpunkterne er det nemmere at bruge standard interpolationsmetoder til at estimere  $R$ . Hvis vi brugte dem direkte ville vi have

$$R(\text{mid}(f)) = r(f) \quad \forall f \in \mathcal{F} \tag{52}$$

Men vi ville ikke vide med sikkerhed om følgende ligning var sand:

$$\int_f R(x) dx \stackrel{?}{=} r(f) \quad \forall f \in \mathcal{F} \tag{53}$$

Vi ville hellere have at det var omvendt; altså at (53) gjaldt altid, mens vi er lidt ligeglade med (52). For at opnå det foreslås det følgende: Lad  $g : (\mathbb{R}^n \times \mathbb{R}^{(n+1) \cdot k}) \mapsto \mathbb{R}$  være en interpolationsmetode der bruger  $k = \#\mathcal{F}$  punkter af formen  $(x_i, y_i)$ , hvor  $x_i \in \mathbb{R}^n, y_i \in \mathbb{R}$  til at kunne forudsige en



Figur 0.1: A hypercube in 2 dimensions encapsulated by  $(z_1^-, z_2^-)$ ,  $(z_1^+, z_2^-)$ ,  $(z_1^-, z_2^+)$ , and  $(z_1^+, z_2^+)$ .

$y$ -værdi når man kommer med et nyt  $x$ -punkt. Vi har allerede  $x_i$ -punkterne fra vores udregning af midtpunkterne, men vi mangler  $y_i$ -værdierne. Dem vil vi bestemme ud fra (53), dvs. løse ligningerne

$$r(f) = \int_f g(x, \{(x_1, y_1), \dots, (x_k, y_k)\}) dx \quad \forall f \in \mathcal{F} \quad (54)$$

med hensyn til  $y_1, \dots, y_k$ . En specielt attraktiv interpolationsfunktion kunne være den  $n$ -dimensionelle lineære interpolation, hvilket vil sige at  $g$  inddeles området i et antal begrænsede  $n$  dimensionale hyper-rektangler. Hver hyperrektangel er afgrænset af  $2^n$  punkter fra  $\{x_1, \dots, x_k\}$ . Vi kan skrive de  $2^n$  punkter

$$(z_1^\pm, z_2^\pm, \dots, z_n^\pm). \quad (55)$$

Et eksempel ses i Figur 0.1. Lad  $y_1, \dots, y_{2^n}$  være de korresponderende  $y$ -værdier for hjørnerne i hyperkuben. Så kan man skrive interpolationsoverfladen som

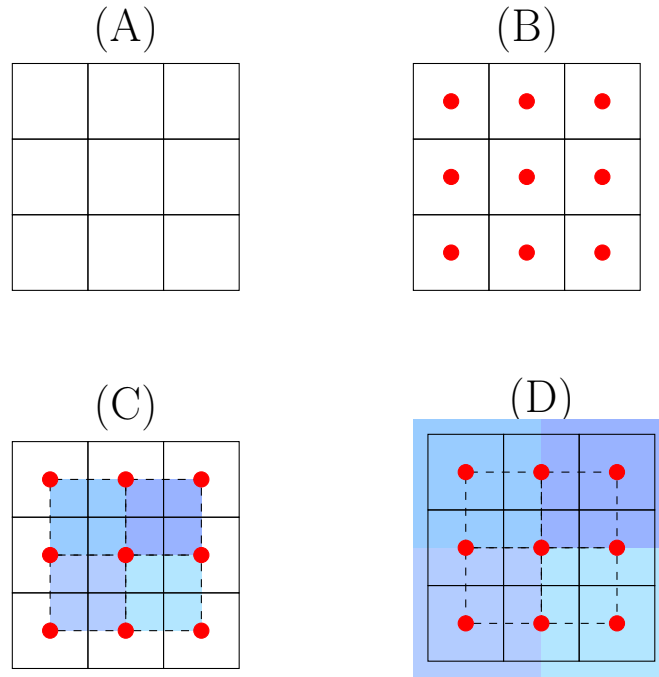
$$f(z) = f((z_1, \dots, z_n)) = \sum_{J \subseteq \{1, \dots, n\}} a_J \prod_{j \in J} z_j \quad (56)$$

hvor

$$\begin{pmatrix} a_\emptyset \\ a_{\{1\}} \\ \vdots \\ a_{\{1, \dots, n\}} \end{pmatrix} = \begin{pmatrix} \vec{v}(z_1^-, z_2^-, \dots, z_n^-) \\ \vec{v}(z_1^+, z_2^-, \dots, z_n^-) \\ \vdots \\ \vec{v}(z_1^+, z_2^+, \dots, z_n^+) \end{pmatrix}^{-1} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{2^n} \end{pmatrix} \quad (57)$$

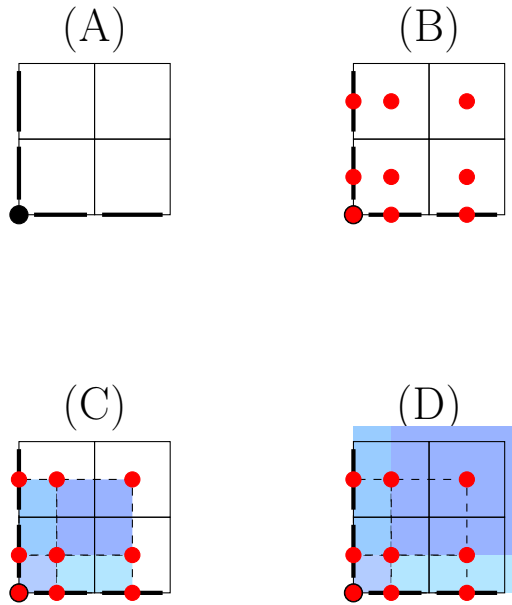
hvor

$$\vec{v}(w_1, \dots, w_n) = (1 \quad w_1 \quad w_2 \quad \dots \quad w_n \quad w_1 w_2 \quad w_1 w_3 \quad \dots \quad \prod_{i=1}^n w_i).$$



Figur 0.2: (A) viser faktorerne i  $\mathcal{F}$ , når der er tale om et grid af endelige faktorlevels med positivt lebesguemål (i  $\mathbb{R}^2$ ). (B) viser midtpunkterne. (C) viser de 4 interpolationsflader der kan laves af de 9 punkter. (D) viser ekstrapoleringen af de 4 interpolationsområder som potentialt går ud i hele  $\mathbb{R}^2$ .

Det gør at hvert  $a$  i (56) kan skrives som en linearkombination af  $y_i$ 'er, og dermed kan de  $k$  ligninger i (54) skrives som en linearkombinationer af  $y_1, \dots, y_k$ . Med andre ord kan (54) (formodentlig) løses entydigt. Hvordan hyperkuberne skal lægges er illustreret i 2 dimensioner i Figur 0.2.

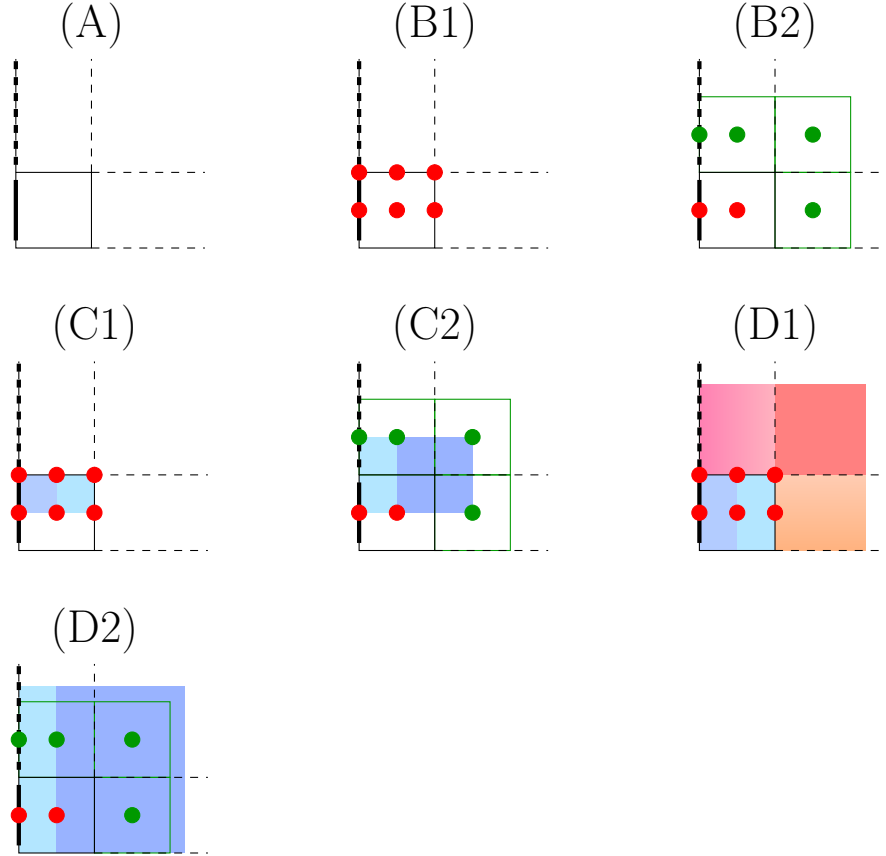


Figur 0.3: (A) viser faktorerne i  $\mathcal{F}$ , når der er tale om et endeligt grid af faktorlevels hvor der dog er punkter og linjer med positivt mål (markeret med sort prik og fed linje, hhv.). (B) viser midtpunkterne. (C) viser de 4 interpolationsflader der kan laves af de 9 punkter. (D) viser ekstrapoleringen af de 4 interpolationsområder som potentialt går ud i hele  $\mathbb{R}^2$ . Læg mærke til at y-værdien i nedre venstre hjørne er kendt, og at der derfor ikke er grund til at opstille integralligningen for den.

### Løsning af Integralsystem

Vi løser nu (54). Vi ser fra (56) at vi bør integrere

$$\begin{aligned}
 & \int_{[z_1^-, z_1^+] \times \dots \times [z_n^-, z_n^+]} z_1 \cdot z_2 \cdots z_m \, d(z_1, \dots, z_n) \\
 &= \int_{[z_2^-, z_2^+] \times \dots \times [z_n^-, z_n^+]} \frac{1}{2} \left( (z_1^+)^2 - (z_1^-)^2 \right) z_2 \cdots z_m \, d(z_2, \dots, z_n) \\
 &= \dots \\
 &= \frac{1}{2^m} \prod_{i=1}^m \left( (z_i^+)^2 - (z_i^-)^2 \right) \int_{[z_{m+1}^-, z_{m+1}^+] \times \dots \times [z_n^-, z_n^+]} 1 \, d(z_1, \dots, z_n) \\
 &= \frac{1}{2^m} \prod_{i=1}^m \left( (z_i^+)^2 - (z_i^-)^2 \right) \prod_{i=m+1}^n (z_i^+ - z_i^-) \tag{58}
 \end{aligned}$$



Figur 0.4: (A) viser faktorerne i  $\mathcal{F}$ , når der er tale om et grid af faktorlevels hvor der er linjer med positivt mål og faktorniveauer der breder sig ud mod uendelig (markeret med fed linje og stiplet linje, hhv.). (B1) viser en metode for midtpunkterne, mens (B2) viser en anden. I (B2) er der også sat markeringer af et pseudo-område for (B2). (D1) og (D2) viser, hvad vi ser som den mest naturlige måde at udvide interpolationsfladerne. I (D1) angiver det store røde felt at interpolationsværdien bliver sat til konstanten  $r(f)$ . De to gradientfelter (hhv. orange og lyserød) indikerer at interpolationsværdien bliver udregnet ved at projicere ned på kanten af den endelige hyperkube.



under antagelsen at  $z_i^+ > z_i^-$ . I det mere generelle tilfælde, hvor  $z_i^+ \geq z_i^-$  antager vi at når  $z_i^+ = z_i^-$  integrerer vi med hensyn til dirac-målet koncentreret i  $z_i^+$ . I så fald bliver (58) i stedet

$$H([(z_1^-, z_1^+), \dots, (z_n^-, z_n^+)], \{1, \dots, m\}) := \prod_{i=1, z_i^- \neq z_i^+}^m \left( \frac{1}{2}(z_i^+)^2 - \frac{1}{2}(z_i^-)^2 \right) \prod_{i=1, z_i^- = z_i^+} z_i^+ \prod_{i=m+1, z_i^- \neq z_i^+}^n (z_i^+ - z_i^-) \quad (59)$$

Nu udvider vi notationen lidt. Lad

$$\text{mid}(f) = (\text{mid}_1(f_1), \dots, \text{mid}_n(f_n)) \quad (60)$$

være midtpunktet for faktorniveau  $f = (f_1, \dots, f_n) \in \mathcal{F}$ . Vi antager at  $\mathcal{F}$  er et grid, dvs.

$$\mathcal{F} = \left\{ [a_1^-(j_1), a_1^+(j_1)] \times \dots \times [a_n^-(j_n), a_n^+(j_n)] : \right. \\ \left. j_1 \in \{1, \dots, k_1\}, \dots, j_n \in \{1, \dots, k_n\} \right\} \quad (61)$$

hvor  $a_i^\pm(j)$  er passende konstanter i  $\mathbb{R} \cup \{-\infty, \infty\}$ , der opfylder

$$a_i^+(j) = a_i^-(j+1), \quad j = 1, \dots, k_i - 1$$

Per definition af midpoint ved vi

$$a_i^-(j) \leq \text{mid}_i([a_i^-(j), a_i^+(j)]) \leq a_i^+(j) \quad (62)$$

Vi forkorter notationen yderligere med

$$z(i, j) = \text{mid}_i([a_i^-(j), a_i^+(j)]).$$

og dermed kan et midtpunkt skrives som

$$(z(1, j_1), z(2, j_2), \dots, z(n, j_n)) \quad (63)$$

for passende  $j_1, \dots, j_n$ . En interpolationsflade (noteret  $f$  i (56)), kan nu skrives mere præcist som

$$h^{j_1, \dots, j_n}(z) = \sum_{J \subseteq \{1, \dots, n\}} a_J(j_1, \dots, j_n) \prod_{j \in J} z_j \quad (64)$$

for  $j_i \in \{1, \dots, k_i - 1\}$ , hvilket giver  $(k_1 - 1) \dots (k_n - 1)$  forskellige interpolationsflader. Interpolationsfladen i (64) gælder umiddelbart kun for  $z$  i mængden

$$[z(1, j_1), z(1, j_1 + 1)] \times \dots \times [z(n, j_n), z(n, j_n + 1)] \quad (65)$$

Dog udvider vi ofte denne mængde, se Figur 0.2, 0.3 og 0.4. For at løse ligningssystemet (54) med (B2)-midpoints sætter vi

$$\tilde{z}(i, j) = \begin{cases} a_i^-(1) & \text{hvis } j = 1, a_i^-(j) > -\infty \\ a_i^+(1) - (a_i^+(2) - a_i^-(2)) & \text{hvis } j = 1, a_i^-(j) = -\infty \\ a_i^+(k_j) & \text{hvis } j = k_j, a_i^+(k_j) < \infty \\ a_i^+(k_j) + (a_i^+(k_j - 1) - a_i^-(k_j - 1)) & \text{hvis } j = k_j, a_i^+(j) = \infty \\ z(i, j) & \text{ellers.} \end{cases} \quad (66)$$

og lader (64) være gældende for  $z$  i mængden

$$[\tilde{z}(1, j_1), \tilde{z}(1, j_1 + 1)] \times \cdots \times [\tilde{z}(n, j_n), \tilde{z}(n, j_n + 1)].$$

I integralet i (54) skal vi integrere over  $f$ . Hvis vi bruger (B2) er det nemmest at formulere integralet ved at justere  $a_i^\pm$ 'erne. Det gøres med

$$\tilde{a}_i^-(j) = \begin{cases} a_i^+(1) - (a_i^+(2) - a_i^-(2)) & \text{hvis } j = 1, a_i^-(j) = -\infty \\ a_i^-(j) & \text{ellers.} \end{cases}$$

og

$$\tilde{a}_i^+(j) = \begin{cases} a_i^+(k_j) + (a_i^+(k_j - 1) - a_i^-(k_j - 1)) & \text{hvis } j = k_j, a_i^+(j) = \infty \\ a_i^+(j) & \text{ellers.} \end{cases}$$

For at gøre notationen nemmere i kanterne af grid'et, sættes også

$$h^{j_1, \dots, j_n} = h^{(j_1 \wedge k_1 - 1) \vee 1, \dots, (j_n \wedge k_n - 1) \vee 1}, \quad (67)$$

$$a_J(j_1, \dots, j_n) = a_J((j_1 \wedge k_1 - 1) \vee 1, \dots, (j_n \wedge k_n - 1) \vee 1) \quad (68)$$

for  $j_i \in \{0, \dots, k_i\}, i = 1, \dots, n, J \subseteq \{1, \dots, n\}$

Bemærk at et  $f \in \mathcal{F}$  er kendetegnet ved  $j_1, \dots, j_n$ . Da kan vi skrive ligningssystemet (54) som  $k$  ligninger, der hver består af en sum af  $2^n$  integraler.

$$\begin{aligned} r(j_1, \dots, j_n) &= \int_{[\tilde{a}_1^-(j_1), z(1, j_1)] \times \cdots \times [\tilde{a}_n^-(j_n), z(n, j_n)]} h^{j_1-1, \dots, j_n-1}(z) \, dz \\ &\quad + \int_{[z(1, j_1), \tilde{a}_1^+(j_1)] \times [\tilde{a}_2^-(j_2), z(2, j_2)] \times \cdots \times [\tilde{a}_n^-(j_n), z(n, j_n)]} h^{j_1, j_2-1, \dots, j_n-1}(z) \, dz \\ &\quad \vdots \\ &\quad + \int_{[z(1, j_1), \tilde{a}_1^+(j_1)] \times \cdots \times [z(n, j_n), \tilde{a}_n^+(j_n)]} h^{j_1, j_2, \dots, j_n}(z) \, dz \end{aligned} \quad (69)$$

Ved at bruge  $H$  fra tidligere kan vi sætte

$$\begin{aligned} & \int_{[\tilde{a}_1^-(j_1), z(1, j_1)] \times \dots \times [\tilde{a}_n^-(j_n), z(n, j_n)]} h^{j_1-1, \dots, j_n-1}(z) \, dz \\ &= \sum_{J \subseteq \{1, \dots, n\}} \left( a_J(j_1-1, \dots, j_n-1) \right. \\ & \quad \left. \cdot H([\tilde{a}_1^-(j_1), z(1, j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n)], J) \right) \end{aligned} \quad (70)$$

Definer nu

$$\mathbf{V}(j_1, \dots, j_n) = \begin{pmatrix} \vec{v}(z(1, j_1), z(2, j_2), \dots, z(n, j_n)) \\ \vec{v}(z(1, j_1+1), z(2, j_2), \dots, z(n, j_n)) \\ \vec{v}(z(1, j_1), z(2, j_2+1), \dots, z(n, j_n)) \\ \vdots \\ \vec{v}(z(1, j_1+1), z(2, j_2+1), \dots, z(n, j_n+1)) \end{pmatrix}. \quad (71)$$

med

$$\mathbf{V}(j_1, \dots, j_n) = \mathbf{V}((j_1 \wedge [k_1 - 1]) \vee 1, \dots, (j_n \wedge [k_n - 1]) \vee 1) \quad (72)$$

for  $j_i \in \{0, \dots, k_i\}$ ,  $i = 1, \dots, n$ , samt

$$\mathbf{y}(j_1, \dots, j_n) = \begin{pmatrix} y(j_1, j_2, \dots, j_n) \\ y(j_1+1, j_2, \dots, j_n) \\ \vdots \\ y(j_1+1, \dots, j_n+1) \end{pmatrix} \quad (73)$$

for  $j_i \in \{1, \dots, k_i - 1\}$ ,  $i = 1, \dots, n$ . Vi definerer

$$\mathbf{y}(j_1, \dots, j_n) = \mathbf{y}((j_1 \wedge [k_1 - 1]) \vee 1, \dots, (j_n \wedge [k_n - 1]) \vee 1) \quad (74)$$

for generelle  $j_i \in \{0, \dots, k_i\}$ ,  $i = 1, \dots, n$ .

Vi fortolker  $y(j_1, \dots, j_n)$  som interpolationens  $y$ -værdi i  $(z(1, j_1), \dots, z(n, j_n))$ .  
Altså er  $\mathbf{y}$  vektoren af  $y$ -værdier, der indgår i  $h^{j_1, \dots, j_n}$ . Indsat i (70) giver det

$$\sum_{J \subseteq \{1, \dots, n\}} \left( \mathbf{e}_J^T \mathbf{V}(j_1-1, \dots, j_n-1)^{-1} \mathbf{y}(j_1-1, \dots, j_n-1) \right. \\ \left. \cdot H([\tilde{a}_1^-(j_1), z(1, j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n)], J) \right) \quad (75)$$

hvor  $\mathbf{e}_J$  er en vektor af 0'er på nær et 1-tal i den indgang der svarer til den koordinat der svarer til  $J$ . Hvis vi definerer vektoren af længde  $2^n$

$$\mathbf{H}([\tilde{a}_1^-(j_1), z(1, j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n)]) \quad (76)$$

$$= \begin{pmatrix} H([\tilde{a}_1^-(j_1), z(1, j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n)], \emptyset) \\ H([\tilde{a}_1^-(j_1), z(1, j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n)], \{1\}) \\ \vdots \\ H([\tilde{a}_1^-(j_1), z(1, j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n)], \{1, \dots, n\}) \end{pmatrix} \quad (77)$$

kan vi skrive (75) som

$$\begin{aligned} & \left[ \mathbf{H}([\tilde{a}_1^-(j_1), z(1, j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n))] \right]^T \\ & \cdot \mathbf{V}(j_1 - 1, \dots, j_n - 1)^{-1} \mathbf{y}(j_1 - 1, \dots, j_n - 1) \end{aligned} \quad (78)$$

Vi kan derfor skrive ligningssystemet fra (69) som

$$\begin{aligned} r(j_1, \dots, j_n) = & \left\{ \begin{aligned} & \left[ \mathbf{H}([\tilde{a}_1^-(j_1), z(1, j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n))] \right]^T \\ & \cdot \mathbf{V}(j_1 - 1, \dots, j_n - 1)^{-1} \mathbf{y}(j_1 - 1, \dots, j_n - 1) \\ & + \left[ \mathbf{H}([z(1, j_1), \tilde{a}_1^+(j_1)), \dots, (\tilde{a}_n^-(j_n), z(n, j_n))] \right]^T \\ & \cdot \mathbf{V}(j_1, j_2 - 1, \dots, j_n - 1)^{-1} \mathbf{y}(j_1, j_2 - 1, \dots, j_n - 1) + \\ & \vdots \\ & + \left[ \mathbf{H}([z(1, j_1), \tilde{a}_1^+(j_1)), \dots, (z(n, j_n), \tilde{a}_n^+(j_n))] \right]^T \\ & \cdot \mathbf{V}(j_1, \dots, j_n)^{-1} \mathbf{y}(j_1, \dots, j_n) \end{aligned} \right\} \end{aligned} \quad (79)$$

Ud fra (79) kræver det et stykke (trivielt) omrøkeringsarbejde at få det på formen

$$\{r(j_1, \dots, j_n)\}_{j_i \in \{1, \dots, k_i\}, i=1, \dots, n} = B \{y(j_1, \dots, j_n)\}_{j_i \in \{1, \dots, k_i\}, i=1, \dots, n} \quad (80)$$

som løses ved at invertere  $k \times k$  matricen  $B$ . Når de  $k$  y-værdier er fundet regnes matrix produkterne

$$\mathbf{V}(j_1, \dots, j_n)^{-1} \mathbf{y}(j_1, \dots, j_n) \quad (81)$$

og koefficienterne sættes ind i risk ratio filerne.