

# Föreläsning 15: Försöksplanering och repetition

## Matematisk statistik

David Bolin  
Chalmers University of Technology  
Oktober 19, 2015



# Utfall och utfallsrum

## Slumpmässigt försök

Man brukar säga att ett slumpmässigt försök är ett försök som kan upprepas under väsentligen identiska förhållanden där utfallet av försöket inte exakt kan förutsägas i det enskilda fallet.

## Utfall och utfallsrum

Resultatet av ett försök kallas för ett *utfall*  $\omega$ , och mängden av alla möjliga utfall kallas för *utfallsrummet*  $\Omega$ .

## Händelser

Man kan också sätta samman olika utfall till *händelser*. En händelse  $A$  är en mängd av utfall, det vill säga en delmängd av utfallsrummet  $\Omega$ .

# Snitt, union och komplement

Låt  $A$  och  $B$  vara två händelser, vi definierar

## Komplement, $A^c$

Mängden av alla utfall som inte finns i  $A$ .  $A^c = \Omega \setminus A$ .

## Union, $A \cup B$

Mängden av alla utfall i  $A$  eller  $B$ .

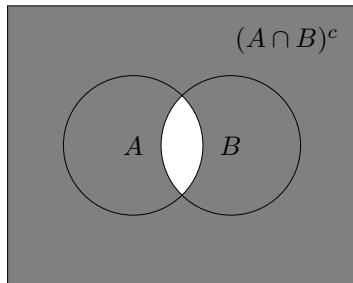
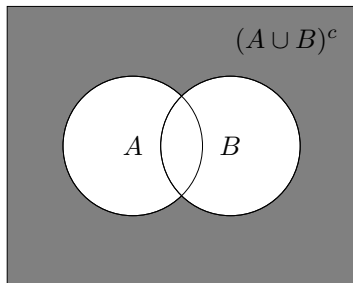
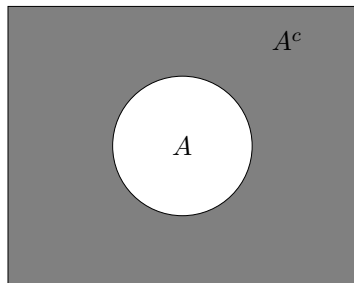
## Snitt $A \cap B$

Mängden av alla utfall som finns i  $A$  och  $B$ .

## Den tomma mängden $\emptyset$

Brukas kallas för en omöjlig händelse. Om  $A \cap B = \emptyset$  så säger vi att  $A$  och  $B$  är oförenliga, disjunkta, eller utesluter varandra.

## Venndiagram



# Kolmogorovs axiomsystem

## Kolmogorovs axiomsystem

Låt  $\Omega$  vara ett utfallsrum. En sannolikhet, eller sannolikhetsmått  $P : \Omega \rightarrow \mathbb{R}$  är en reell funktion som tar händelser i  $\Omega$  som argument och returnerar ett reellt tal, och som uppfyller

- 1  $0 \leq P(A) \leq 1$ .
- 2  $P(\Omega) = 1$ .
- 3 Om  $A_1, A_2, \dots$  är en följd av disjunkta händelser så gäller att

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

Speciellt så gäller att om händelserna  $A$  och  $B$  är disjunkta så är

$$P(A \cup B) = P(A) + P(B).$$

# Egenskaper hos sannolikhetsmättet

Ur axiomen kan vi härleda följande egenskaper.

## Egenskaper

För ett sannolikhetsmått  $P$  gäller att

- 1  $P(\emptyset) = 0$ .
- 2  $P(A^c) = 1 - P(A)$ .
- 3  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

Att dessa är uppfyllda inses enkelt genom att rita Venndiagram!

## Betingade sannolikheter

Vi vill ibland beräkna sannolikheten för en händelse  $A$  givet att vi vet att en annan händelse  $B$  har inträffat. Vi vill då veta den så kallade betingade sannolikheten för  $A$  givet  $B$ .

### Betingad sannolikhet

Antag att  $P(B) > 0$ . Den betingade sannolikheten för  $A$  givet  $B$  definieras som

$$P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

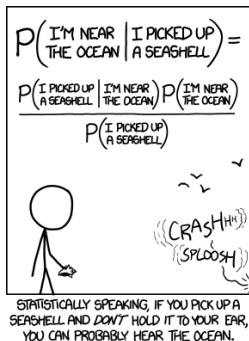
### Multiplikationssatsen för sannolikheter

Om  $A$  och  $B$  är händelser så gäller att

$$P(A \cap B) = P(B|A)P(A) = P(A|B)P(B).$$

Multiplikationssatsen är speciellt användbar för att beräkna sannolikheten för successiva händelser som påverkar varandra.

## Bayes sats



## Bayes sats

För två händelser  $A$  och  $B$  gäller att

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A^c)P(A^c)}$$



# Slumpvariabler

En slumpvariabel beskriver utfallet av ett slumpmässigt försök med ett numeriskt värde:

## Slumpvariabel

En stokastisk variabel (slumpvariabel)  $X$  är en funktion som tar element från  $\Omega$  som argument och som returnerar ett reellt tal.

Vi betecknar ofta slumpvariabler med stora bokstäver,  $X$ ,  $Y$  och  $Z$ . Utfall betecknas ofta med respektive små bokstäver,  $x$ ,  $y$  och  $z$ .

## Diskreta slumpvariabler

En diskret slumpvariabel kan endast anta ett ändligt eller uppräknligt antal värden, i allmänhet någon delmängd av heltalen.

# Sannolikhetsfunktionen

## Sannolikhetsfunktion

Till en diskret stokastisk variabel  $X$  definierar vi sannolikhetsfunktionen  $p(k)$  genom  $p(k) = P(X = k)$ .

Eftersom  $p(k)$  beskriver sannolikheter måste vi ha att

- $p(k) \geq 0$  för alla  $k$ .
- $\sum_{k=-\infty}^{\infty} p(k) = 1$ .

Dessa två villkor är både nödvändiga och tillräckliga för att en funktion  $p(k)$  ska vara en sannolikhetsfunktion.

Vi kan också visa att

$$P(m \leq X \leq n) = \sum_{k=m}^n p(k)$$

om  $m$  och  $n$  är heltal.

# Kontinuerliga fördelningar

I praktiken stöter vi ofta på situationer då de möjliga värdena är alla värden i ett intervall, tex  $[0, 1]$ , eller är alla positiva reella tal.

## Kontinuerlig slumpvariabel

En kontinuerlig slumpvariabel kan anta alla värden i något (eller några) intervall av reella tal och sannolikheten för att den antar varje specifikt värde är noll.

Kontinuerliga slumpvariabler kan beskrivas med täthetsfunktionen.

## Täthetsfunktion

Låt  $X$  vara en kontinuerlig slumpvariabel, en funktion  $f(x)$  så att

$$f(x) \geq 0, \quad \int_{-\infty}^{\infty} f(x)dx = 1, \quad P(a \leq X \leq b) = \int_a^b f(x)dx,$$

kallas en täthetsfunktion.

# Fördelningsfunktioner

Ett annat vanligt sätt att beskriva en stokastisk variabel är att använda den så kallade fördelningsfunktionen.

## Fördelningsfunktionen

Låt  $X$  vara en slumpvariabel. Dess fördelningsfunktion ges då av  $F(x) = P(X \leq x)$ . Om  $X$  är en diskret variabel har vi

$$F(x) = \sum_{k \leq x} p(k),$$

och om  $X$  är kontinuerlig har vi

$$F(x) = \int_{-\infty}^x f(y) dy.$$

## Väntevärde

Väntevärdet för en slumpvariabel definieras som

$$E(X) = \begin{cases} \sum_{k=-\infty}^{\infty} kp(k) & \text{om } X \text{ är diskret,} \\ \int_{-\infty}^{\infty} xf(x)dx & \text{om } X \text{ är kontinuerlig.} \end{cases}$$

Väntevärdet är "tyngdpunkten" i fördelningen.

## Sats

Vi har att

$$E(g(X)) = \begin{cases} \sum_{k=-\infty}^{\infty} g(k)p(k), & \text{om } X \text{ är diskret,} \\ \int_{-\infty}^{\infty} g(x)f(x)dx, & \text{om } X \text{ är kontinuerlig.} \end{cases}$$

# Varians och standardavvikelse

## Varians

Variansen av en stokastisk variabel definieras som

$$V(X) = E[(X - \mu)^2], \text{ d\AAr } \mu \text{ \AAr v\AAntev\AArdet av } X.$$

Vi ser allts\AA att variansen definieras som v\AAntev\AArdet av den kvadratavvikelsen av  $X$  fr\AAn dess v\AAntev\AArde. Vi ber\AAknar variansen som

$$V(X) = \begin{cases} \sum_{k=-\infty}^{\infty} (k - \mu)^2 p_X(k), & \text{om } X \text{ \AAr diskret} \\ \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx, & \text{om } X \text{ \AAr kontinuerlig.} \end{cases}$$

Ett enklare s\AAtt \AAr ofta att ber\AAkna variansen som

$$V(X) = E(X^2) - \mu^2$$

Standardavvikelsen av en stokastisk variabel ges av  $\sigma = \sqrt{V(X)}$ .

# Normalfördelningen

## Normalfördelningen

En kontinuerlig slumpvariabel  $X$  är normalfördelad,  $N(\mu, \sigma^2)$ , med parametrar  $\mu \in \mathbb{R}$  och  $\sigma > 0$ , om den har täthetsfunktionen

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

Fördelningsfunktionen ges av

$$F(x) = \int_{-\infty}^x \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^2\right) dy$$

## Parametrar

Om  $X \sim N(\mu, \sigma^2)$  har vi att  $E(X) = \mu$  och  $V(X) = \sigma^2$ .

## Beräkning av sannolikheter

## Standardiserad normalfördelning

En slumpvariabel  $Z$  sägs ha en standardiserad normalfördelning om  $Z \sim N(0, 1)$ . Vi beteckningar fördelningsfunktion och täthetsfunktion för denna fördelning med  $\varphi(x)$  respektive  $\Phi(x)$ .

## Sats

Om  $X \sim N(\mu, \sigma^2)$  så gäller att  $aX + b \sim N(a\mu + b, a^2\sigma^2)$ .

Detta betyder att om  $X \sim N(\mu, \sigma^2)$

- kan vi skriva  $X = \mu + \sigma Z$  där  $Z \sim N(0, 1)$ .
- har vi att  $Z = (X - \mu)/\sigma \sim N(0, 1)$ .

Vi kan använda detta för att beräkna sannolikheter:

$$P(X < x) = P\left(\frac{X - \mu}{\sigma} < \frac{x - \mu}{\sigma}\right) = P\left(Z < \frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$



# Centrala gränsvärdessatsen

En av de viktigaste satserna inom statistikteorin, och en av anledningarna till varför normalfördelningen är så viktig.

## CGS

Låt  $X_1, \dots, X_n$  vara oberoende och likafördelade slumpvariabler med väntevärde  $\mu$  och varians  $\sigma^2 < \infty$ . Då gäller att

$$P\left(\frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} \leq x\right) \rightarrow \Phi(x), \quad \text{då } n \rightarrow \infty.$$

Om  $n$  är stort har vi enligt satsen att

- $\sum_{i=1}^n X_i$  är approximativt  $N(n\mu, n\sigma^2)$ -fördelad.
- $\bar{X} = n^{-1} \sum_{i=1}^n X_i$  är approximativt  $N(\mu, \sigma^2/n)$ -fördelad.

Hur stor  $n$  måste vara beror på fördelningen av  $X_i$ .

## Tvådimensionella fördelningar

## Definition

En två dimensionell slumpvariabel  $(X, Y)$  tillordnar två numeriska värden till varje utfall i  $\Omega$ .

För diskreta variabler har vi sannolikhetsfunktionen

$$p_{X,Y}(i, j) = P(X = i, Y = j).$$

För en sannolikhetsfunktion har vi att  $p_{X,Y}(i, j) \geq 0$  och  $\sum_{i,j} p_{X,Y}(i, j) = 1$ . För kontinuerliga variabler har vi en täthetsfunktion  $f_{X,Y}(x, y)$  som är sådan att

- 1  $f_{X,Y}(x, y) \geq 0$ ,
- 2  $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy = 1$ , och
- 3  $P(a \leq X \leq b \text{ och } c \leq Y \leq d) = \int_a^b \int_c^d f_{X,Y}(x, y) dx dy$ .

## Marginalfördelningar

Givet en diskret tvådimensionell slumpvariabel  $(X, Y)$  definieras marginalfördelningarna för  $X$  och  $Y$  som

$$p_X(i) = \sum_{j=-\infty}^{\infty} p_{X,Y}(i, j), \quad p_Y(j) = \sum_{i=-\infty}^{\infty} p_{X,Y}(i, j)$$

I det kontinuerliga fallet har vi på motsvarande sätt

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy, \quad f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx$$

## Väntevärden av funktioner

Väntevärdet av en funktion  $H(X, Y)$  definieras som

$$E(H(X, Y)) = \begin{cases} \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} H(i, j) p_{X,Y}(i, j), & \text{diskret} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} H(x, y) f_{X,Y}(x, y) dx dy, & \text{kontinuerlig} \end{cases}$$

## Beroende

## Oberoende händelser

Två händelser  $A$  och  $B$  är oberoende om  $P(A \cap B) = P(A)P(B)$ .

## Oberoende slumpvariabler

Två slumpvariabler  $X$  och  $Y$  är oberoende om deras täthetsfunktion (eller sannolikhetsfunktion) kan skrivas som produkten av marginalfördelningarna. Det vill säga, för det diskreta fallet

$$p_{X,Y}(i, j) = p_X(i)p_Y(j) \quad \text{diskreta variabler}$$

$$f_{X,Y}(x, y) = f_X(x)f_Y(y) \quad \text{kontinuerliga variabler}$$

## Kovarians

Kovariansen mellan  $X$  och  $Y$  definieras som

$$C(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] \text{ där } \mu_X = E(X) \text{ och } \mu_Y = E(Y).$$

## Sats

Kovariansen kan beräknas som

$$C(X, Y) = E(XY) - E(X)E(Y)$$

## Sats

Om  $X$  och  $Y$  är oberoende så är  $C(X, Y) = 0$  och  $E(XY) = E(X)E(Y)$ .

## Korrelation

Den så kallade korrelationskoefficienten definieras som

$$\rho(X, Y) = \frac{C(X, Y)}{\sqrt{V(X)V(Y)}}.$$

# Betingade fördelningar

Den betingade fördelningen för  $X$  givet  $Y = y$  definieras för det kontinuerliga fallet som

$$f_{X|y}(x) = \frac{f_{X,Y}(x, y)}{f_Y(y)},$$

givet att  $f_Y(y) > 0$ . På motsvarande sätt har vi den betingade fördelningen för  $Y$  givet  $X = x$ :

$$f_{Y|x}(y) = \frac{f_{X,Y}(x, y)}{f_X(x)},$$

givet att  $f_X(x) > 0$ .

Samma formler gäller för diskreta variabler om vi byter  $f$  mot  $p$ .

# Punktskattningar

## Stickprov

Ett stickprov av storlek  $n$  är  $n$  oberoende observationer av en slumpvariabel  $X$ . Vi kan alltså skriva stickprovet som observationer av  $n$  slumpvariabler  $X_1, \dots, X_n$  där alla  $X_i$  är oberoende och likafördelade.

## Skattning

En skattning av en parameter  $\theta$  är en funktion  $\hat{\theta}(X_1, \dots, X_n)$  av observationerna.

Två egenskaper vi vill att vår skattare ska ha är att

- Den ska vara *väntevärdesriktig* (unbiased på engelska), vilket betyder att vi vill att  $E(\hat{\theta}(X_1, \dots, X_n)) = \theta$ .
- Den ska ha låg varians om  $n$  är stort, helst vill vi att  $V(\hat{\theta}(X_1, \dots, X_n)) \rightarrow 0$  då  $n \rightarrow \infty$ .

## Skattning av väntevärde och varians

Låt  $x_1, \dots, x_n$  vara observationer av oberoende och likafördelade s.v.  $X_1, \dots, X_n$  med väntevärde  $\mu$  och standardavvikelse  $\sigma$ .

Väntevärdesriktiga skattningar av  $\mu$  och  $\sigma^2$  är då

- $\mu^* = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$
- $(\sigma^2)^* = S^2 = \frac{Q}{n-1} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ .

Om  $X_i \in N(\mu, \sigma^2)$  så har vi  $\mu^* \sim N(\mu, \frac{\sigma^2}{n})$  och  $\frac{Q}{\sigma^2} \in \chi^2(n-1)$

Låt  $x_{i1}, \dots, x_{in_i}$  vara ober. obs. från  $N(\mu_i, \sigma)$  då  $i = 1, \dots, k$ . Då är en poolad variansskattning

- $(\sigma^2)^* = S_p^2 = \frac{Q}{f} = \frac{(n_1-1)S_1^2 + \dots + (n_k-1)S_k^2}{(n_1-1) + \dots + (n_k-1)}$

Eftersom  $X_{ij} \sim N(\mu_i, \sigma^2)$  så har vi  $\frac{Q}{\sigma^2} \sim \chi^2(f)$



# Maximum likelihood-metoden

Tanken bakom maximum likelihood metoden är att vi vill hitta det parametervärde som mest troligast producerade det stickprov vi har. Metoden är baserad på den så kallade likelihood-funktionen, som för ett stickprov  $x_1 \dots, x_n$  är

$$L(\theta) = \prod_{i=1}^n f(x_i)$$

där  $f(x)$  är täthetsfunktionen för fördelningen och  $\theta$  är parametrarna vi vill skatta. I det diskreta fallet byter vi täthetsfunktionen mot sannolikhetsfunktionen.

## ML skattare

Maximum likelihood-skattaren av en parameter  $\theta$  ges av

$$\hat{\theta}_{ML} = \arg \max_{\theta} L(\theta)$$

# Konfidensintervall

## Konfidensintervall

Låt  $X_1, \dots, X_n$  vara slumpvariabler med en fördelning som har  $\theta$  som en parameter med  $\theta_0$  som sant okänt värde. Ett  $100(1 - \alpha)\%$  konfidensintervall för  $\theta$  med konfidensgraden  $1 - \alpha$  är ett intervall  $[a(X_1, \dots, X_n), b(X_1, \dots, X_n)]$  sådant att

$$P(a \leq \theta_0 \leq b) = 1 - \alpha.$$

- $(a, b)$  är ett slumpmässigt intervall, eftersom  $a$  och  $b$  är slumpvariabler som beror av  $X_1, \dots, X_n$ .
- Konfidensintervallet ska alltså tolkas som att om vi gör upprepade mätningar av  $X_1, \dots, X_n$  och bildar konfidensintervallet för alla dessa mätningar så kommer  $100(1 - \alpha)\%$  av dessa intervall täcka det sanna värdet  $\theta_0$ .

## Konfidensintervall

Om  $\theta^*(X_1, \dots, X_n) \in N(\theta, V(\theta^*))$  så har vi

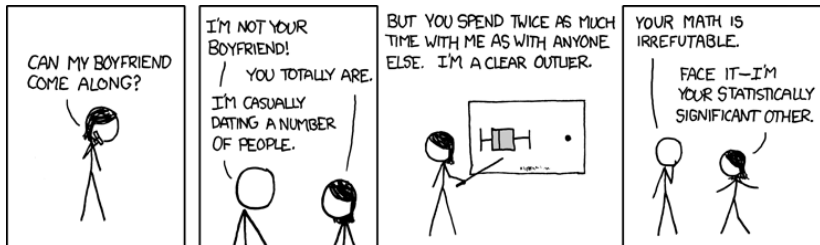
- $I_\theta = (\theta^* \pm \lambda_{\alpha/2} \cdot \sqrt{V(\theta^*)})$  om  $V(\theta^*)$  är känd
- $I_\theta = (\theta^* \pm t_{\alpha/2}(f) \cdot s \cdot c)$  om  $V(\theta^*) = c^2 \cdot \sigma^2$  där  $\sigma^2 = V(X_i)$ ,  
 $c$  är en konstant och  $(\sigma^2)^* = S^2 = \frac{Q}{f}$  med  $\frac{Q}{\sigma^2} \sim \chi^2(f)$

Om  $\theta^*(X_1, \dots, X_n) \approx N(\theta, V(\theta^*))$  enligt CGS (el. dyl.) så kan vi använda normalbaserade konfidensintervall approximativt.

Om  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  med  $(\sigma^2)^* = S^2 = \frac{Q}{f}$  och  $\frac{Q}{\sigma^2} \in \chi^2(f)$  så

$$I_{\sigma^2} = \left( \frac{f \cdot s^2}{\chi_{\alpha/2}^2(f)}, \frac{f \cdot s^2}{\chi_{1-\alpha/2}^2(f)} \right)$$

## Hypotesprövning



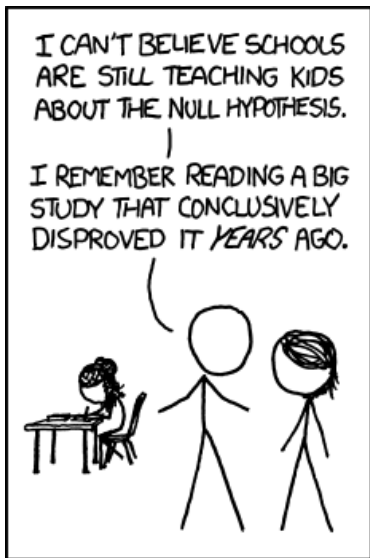
Vi har ett stickprov med någon fördelning med okänd parameter  $\theta$ .  
Vi vill testa *nollhypotesen*

$$H_0 : \theta = \theta_0$$

mot *mothypotesen*

$$H_1 : \theta \neq \theta_0.$$

Om testet bekräftar att  $\theta \neq \theta_0$  säger vi att  $H_0$  förkastas till förmån för  $H_1$ . Man brukar då också säga att  $\theta$  är signifikant skild från  $\theta_0$ .



# Konfidensintervallmetoden

Om vi bildar ett 95% konfidensintervall  $I_\theta$  för  $\theta$  kan vi direkt använda detta för att göra hypotestestet.

- Vi förkastar  $H_0$  om intervallet inte innehåller  $\theta_0$ .
- Om intervallet däremot innehåller  $\theta_0$  kan vi inte förkasta  $H_0$ .

Konfidensgraden vi använder när vi beräknar konfidensintervallet är också konfidensgraden för hypotestestet vi gör.

## Definition

Felrisken eller signifikansnivån definieras som

$$\alpha = P(H_0 \text{ förkastas} \mid H_0 \text{ sann}).$$

Översatt till konfidensintervall är signifikansnivån

$$\alpha = P(\theta \notin I_\theta \text{ om } \theta = \theta_0) = P(\theta_0 \notin I_\theta).$$

## Typ 2 fel

Felet vi gör om vi förkastar  $H_0$  trots att den är sann brukar kallas för ett "Typ 1 fel".

Den andra sortens fel vi kan göra är "Typ 2 fel":

## Definition

Under hypotesprövning säger vi att vi gör ett typ 2 fel om vi inte förkastar  $H_0$  trots att  $H_1$  är sann. Sannolikheten för att göra ett typ 2 fel brukar betecknas med  $\beta$ :

$$\beta = P(H_0 \text{ förkastas inte} \mid H_1 \text{ sann}).$$

Sannolikheten för att göra ett sådant fel betecknas med  $\beta$ . Sannolikheten att *inte* göra ett typ 2 fel kallas för testets *styrka*, och ges alltså av  $1 - \beta$ .

# Teststorheter och p-värden

## Teststorhet

En teststorhet  $T = T(X_1, \dots, X_n)$  är en funktion av observationerna, och alltså en slumpvariabel.  $T_{obs} = T(x_1, \dots, x_n)$  är ett observerat värde av teststorheten för givna observationer.

## p-värde

p-värdet eller signifikanssannolikheten definieras som sannolikheten under nollhypotesen att vi får ett värde  $|T|$  som är lika stort eller större än det observerade värdet  $|T_{obs}|$ .

## Kritiskt område

Givet en signifikansnivå  $\alpha$  definierar vi det kritiska området  $C_\alpha$  som de värden på teststorheten  $T$  som leder till att man förkastar  $H_0$  på nivån  $\alpha$ .



# Oberoende stickprov

Vi antar att vi har två oberoende stickprov

- $n_1$  observationer  $X_{11}, X_{12}, \dots, X_{1n_1}$  från  $N(\mu_1, \sigma_1^2)$ .
- $n_2$  observationer  $X_{21}, X_{22}, \dots, X_{2n_2}$  från  $N(\mu_2, \sigma_2^2)$ .

Den intressanta parametern är nu  $\theta = \mu_1 - \mu_2$  som vi skattar med  $\theta^* = \bar{X}_1 - \bar{X}_2$  och vi vill nu bilda ett konfidensintervall för  $\theta$  samt testa hypotesen

$$H_0 : \theta = 0,$$

vilket är samma sak som att testa  $\mu_1 = \mu_2$ , mot

$$H_1 : \theta \neq 0$$

eller någon ensidig hypotes,  $\theta > 0$  (som motsvarar  $\mu_1 > \mu_2$ ) eller  $\theta < 0$  (som motsvarar  $\mu_1 < \mu_2$ ). Vi skiljer på tre fall:

- 1  $\sigma_1$  och  $\sigma_2$  är kända.
- 2  $\sigma_1 = \sigma_2 = \sigma$  där  $\sigma$  är okänd.
- 3  $\sigma_1$  och  $\sigma_2$  är okända och ej säkert lika.

# Testet i de olika fallen

Vi bildar teststorheten

$$T = \frac{\theta^*}{\sqrt{V(\theta^*)}}$$

- Om  $\sigma_1$  och  $\sigma_2$  är kända gäller att  $V(\theta^*) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$  och  $T$  är  $N(0, 1)$ -fördelad.
- Om  $\sigma_1 = \sigma_2 = \sigma$  där  $\sigma$  är okänd använder vi en poolad variansskattning och skattar  $V(\theta^*)$  med  $s_p^2(1/n_1 + 1/n_2)$ .  $T$  är då  $t(n_1 + n_2 - 2)$ -fördelad.
- Om  $\sigma_1$  och  $\sigma_2$  är okända och  $\sigma_1 \neq \sigma_2$  skattar vi  $V(\theta^*)$  med  $s_1^2/n_1 + s_2^2/n_2$ .  $T$  är då approximativt  $t(f)$ -fördelad med

$$f = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}}$$

# Jämförelse av variansen för oberoende stickprov

För att jämföra de två varianserna inför vi teststorheten

$$T = \frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2}$$

$T$  är  $F(n_1 - 1, n_2 - 1)$ -fördelad. Vi låter  $F_\alpha(f_1, f_2)$  beteckna  $\alpha$ -kvantilen i  $F$ -fördelningen och får att ett konfidensintervall för  $\sigma_1^2/\sigma_2^2$  ges av

$$I_{\sigma_1^2/\sigma_2^2} = \left[ \frac{s_1^2/s_2^2}{F_{\alpha/2}(n_1 - 1, n_2 - 1)}, \frac{s_1^2/s_2^2}{F_{1-\alpha/2}(n_1 - 1, n_2 - 1)} \right]$$

För hypotestest av  $\sigma_1^2/\sigma_2^2 = 1$  bildar vi teststorheten

$$T = s_1^2/s_2^2$$

som under  $H_0$  är  $F(n_1 - 1, n_2 - 1)$ -fördelad.

## Stickprov i par

En annan vanlig situation är att mätningarna uppkommer i par, till exempel om

- Man vill studera hur mycket rökare går upp i vikt när de slutar röka. Man mäter då vikten före och efter för varje person före och efter den slutar röka och jämförelsen sker för varje person.
- Man vill studera systematiska skillnader mellan två mätmetoder och använder varje metod på var och ett av ett antal prover och jämför de två metoderna för varje prov.

Modellen vi nu ansätter är att vi har  $n$  observationer i två stickprov

$$X_1, X_2, \dots, X_n \qquad Y_1, Y_2, \dots, Y_n$$

För varje mätning bildar vi differensen som antas vara normalfördelad:

$$D_i = X_i - Y_i \sim N(\Delta, \sigma^2)$$

Vi vill nu testa om  $\Delta = 0$ , vilket görs som vanligt för normalfördelade mätningar.

# Binomialfördelningen

Låt  $X \sim \text{Bin}(n, p)$ .  $p^* = x/n$  är en väntevärdesriktig skattning av  $p$  med varians  $V(p^*) = p(1 - p)/n$ .

Enligt CGS är  $X$  approximativt  $N(np, np(1 - p))$ -fördelad om  $n$  är stor ( $np(1 - p) > 10$ ). Vi har då att

$$T = \frac{p^* - p}{\sqrt{p^*(1 - p^*)/n}}$$

är approximativt  $N(0, 1)$ -fördelad.

Om vi vill testa

$$H_0 : p = p_0$$

$$H_1 : p \neq p_0$$

eller en ensidig hypotes  $H_1 : p > p_0$  eller  $H_1 : p < p_0$  så är  $X \sim \text{Bin}(n, p_0)$  under  $H_0$ , och vi kan direkt beräkna p-värdet

## Jämförelser i binomialfördelningen

Antag att  $X_1 \sim \text{Bin}(n_1, p_1)$  och  $X_2 \sim \text{Bin}(n_2, p_2)$  och vi vill testa om det är rimligt att anta att  $p_1 = p_2$ .

$p_1^* - p_2^*$  är en väntevärdesriktig skattning av  $p_1 - p_2$ .

Om  $n_1 p_1 (1 - p_1)$  och  $n_2 p_2 (1 - p_2)$  båda är stora (säg större än 10) så är  $p_1^* - p_2^*$  approximativt normalfördelad. Vi har

$$V(p_1^* - p_2^*) = \frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}$$

ett konfidensintervall för  $p_1 - p_2$  fås som

$$I_{p_1 - p_2} = \left( p_1^* - p_2^* \pm z_{\alpha/2} \sqrt{\frac{p_1^*(1 - p_1^*)}{n_1} + \frac{p_2^*(1 - p_2^*)}{n_2}} \right)$$

# Hypotestest

Om vi vill testa  $H_0 : p_1 = p_2$  kan vi använda en poolad variansskattning eftersom vi har att under  $H_0$  så gäller

$$V_{H_0}(p_1^* - p_2^*) = p(1 - p)\left(\frac{1}{n_1} - \frac{1}{n_2}\right)$$

Här är  $p$  det gemensamma värdet under  $H_0$ , vilket skattas som

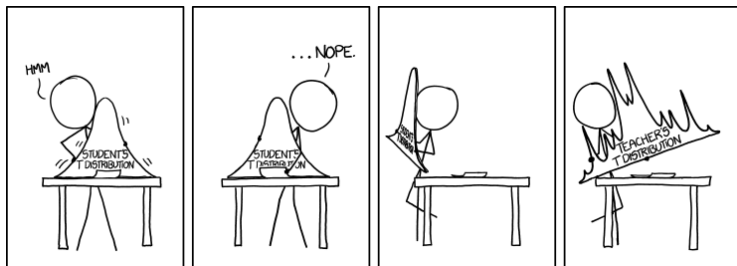
$$p^* = \frac{x_1 + x_2}{n_1 + n_2}$$

Vi bilar då teststorheten

$$T = \frac{p_1^* - p_2^*}{\sqrt{p^*(1 - p^*)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

som under  $H_0$  är approximativt  $N(0, 1)$  fördelad.

## Icke-parametriska metoder



Om vi inte hittar någon fördelning som passar när vi vill göra ett test kan vi använda någon icke-parametrisk metod.

Givet ett stickprov av storlek  $n$  vill vi testa  $H_0 : M = M_0$  mot en ensidig eller tvåsidig mothypotes. Vi kan då använda ett

- Teckentest eller Wilcoxons ranktest.

Givet två oberoende stickprov  $X_1, \dots, X_m$  och  $Y_1, \dots, Y_n$  från några fördelningar  $X$  respektive  $Y$  vill vi testa  $H_0 : M_X = M_Y$ .

- Vi kan då använda Wilcoxons ranksummetest.



# Enkel linjär regression

Vi har talpar  $(x_i, Y_i)$  där  $x_i$  är ett fixt värde och  $Y_i$  är en slumpvariabel. Modellen vi ansätter för  $Y_i$  är

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

där  $\varepsilon_i$  är oberoende  $N(0, \sigma^2)$  slumpvariabler som beskriver mätfelet och  $\beta_0$  och  $\beta_1$  är okända parametrar som beskriver det linjära sambandet.

Vi har alltså att väntevärdet av  $Y$  ges av funktionen

$$E(Y_i) = f(\beta_0, \beta_1, x_i)$$

Vi skattar parametrarna enligt minsta-kvadratmetoden genom att minimera kvadratfelet för den datan vi har

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - f(\beta_0, \beta_1, x_i))^2$$

# Multipel linjär regression

Vi har gjort mätningar av en responsvariabel  $Y$  för fixerade värden på förklarande variabler  $x_1, \dots, x_p$ , som kan väljas utan fel.

Vi antar en linjär modell för  $(Y_i, x_{1,i}, \dots, x_{p,i}), i = 1, \dots, n$ :

$$Y_i = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_p x_{p,i} + \varepsilon_i \quad (1)$$

där  $\varepsilon_i$  är oberoende  $N(0, \sigma^2)$  slumpvariabler som beskriver mätfelet och  $\beta_i$  är parametrar som beskriver beroendet.

Vi kan formulera modellen på matrisform som

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{E}$$

där

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & \cdots & x_{p,1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1,n} & \cdots & x_{p,n} \end{pmatrix} \quad \mathbf{E} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix}$$

## Parameterskattning

## Sats

Givet att  $\mathbf{X}^T \mathbf{X}$  är inverterbar fås minstakvadrat-skattningen av  $\beta$  som  $\beta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ . Skattningarna är väntevärdesriktiga och vi har

$$V(\beta_i^*) = \sigma^2 (\text{diagonalelement nummer } i+1 \text{ i } (\mathbf{X}^T \mathbf{X})^{-1}).$$

Vidare är  $s^2 = Q_0 / (n - p - 1)$  där

$$Q_0 = \sum_{i=1}^n (y_i - \beta_0^* - \beta_1^* x_{1,i} - \dots - \beta_p^* x_{p,i})^2 = \mathbf{Y}^T \mathbf{Y} - \beta^T \mathbf{X}^T \mathbf{Y}.$$

## Sats

Med  $\mu_Y^*(\mathbf{x}_0) = \beta_0^* + \beta_1^* x_{0,1} + \dots + \beta_p^* x_{0,p}$  är

$$V(\mu_Y^*(\mathbf{x}_0)) = \sigma^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0$$

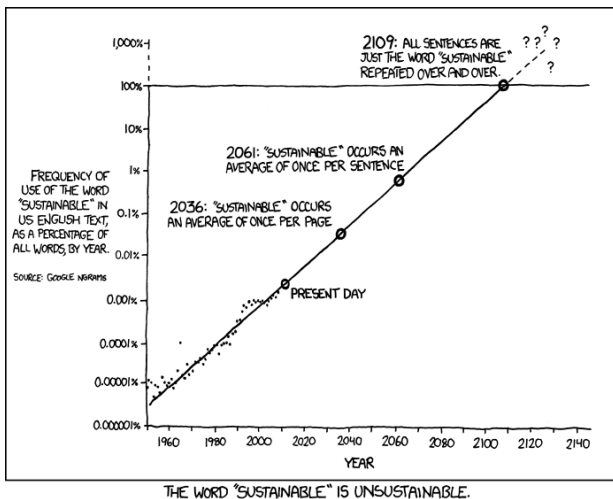
## Egenskaper hos skattningarna

- Om  $\varepsilon_i$  är normalfördelade är också  $\beta_i^*$  normalfördelade.
- $(n - p - 1)s^2/\sigma^2 \sim \chi^2(n - p - 1)$ .
- med  $\theta = \beta_i$  eller  $\mu_Y(\mathbf{x}_0)$  gäller  $(\theta^* - \theta)/d(\theta^*) \sim t(n - p - 1)$ .
- Ett konfidensintervall för  $\theta$  är

$$I_\theta = [\theta^* \pm t_{\alpha/2}(n - p - 1)d(\theta^*)]$$

- För att testa  $H_0 : \theta = \theta_0$  mot  $H_1 : \theta \neq \theta_0$  används teststorheten  $T = (\theta^* - \theta_0)/d(\theta^*)$  som har kritiskt område  $C_\alpha = \{T : |T| > t_{\alpha/2}(n - p - 1)\}$ .
- Ett prediktionsintervall för en framtida observation  $Y(\mathbf{x}_0)$  ges av

$$I_{Y(\mathbf{x}_0)} = \left[ \mu_Y^*(\mathbf{x}_0) \pm t_{\alpha/2}(n - p - 1) \cdot s \sqrt{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0} \right]$$



Det viktigt att validera att modellantagandena är uppfyllda genom att studera residualerna. Det är också viktigt att inse att vi inte kan vara säkra på att modellen stämmer för  $x$  utanför mätområdet.

## Faktor försök

- Ett försök med  $k$  faktorer, där varje faktor testas med två nivåer kallas ett  $2^k$ -faktor försök.
- Till exempel för ett  $2^3$ -försök är modellen

$$Y_{ijkl} = \mu + Ax_i + Bx_j + Cx_k + ABx_ix_j + ACx_ix_k + BCx_jx_k + ABCx_ix_jx_k + \varepsilon_{ijkl}$$

med  $\varepsilon_{ijkl} \sim N(0, \sigma^2)$ , där  $x_i = -1$  om faktor  $A$  har låg nivå och  $x_i = 1$  om faktor  $A$  har hög nivå, och  $x_j$  och  $x_k$  är definierade på motsvarande sätt för faktor  $B$  och  $C$ .

- Försöket sägs ha  $n$  replikat om varje faktorkombination mäts  $n$  gånger.
- Huvudeffekter och samspelseffekter kan skattas med hjälp av Yates schema eller ett teckenschema.
- För en given effekt  $\theta$  har vi  $d(\hat{\theta}) = s/\sqrt{2^k n}$ .
- $I_\theta = [\hat{\theta} \pm t_{\alpha/2}(2^k(n-1))d(\hat{\theta})]$ .

Teckenschema för ett  $2^3$ -försök

| Försök | Medel           | $\mu$ | A | B | C | AB | AC | BC | ABC |
|--------|-----------------|-------|---|---|---|----|----|----|-----|
| (1)    | $\bar{y}_{111}$ | +     | - | - | - | +  | +  | +  | -   |
| a      | $\bar{y}_{211}$ | +     | + | - | - | -  | -  | +  | +   |
| b      | $\bar{y}_{121}$ | +     | - | + | - | -  | +  | -  | +   |
| ab     | $\bar{y}_{221}$ | +     | + | + | - | +  | -  | -  | -   |
| c      | $\bar{y}_{112}$ | +     | - | - | + | +  | -  | -  | +   |
| ac     | $\bar{y}_{212}$ | +     | + | - | + | -  | +  | -  | -   |
| bc     | $\bar{y}_{122}$ | +     | - | + | + | -  | -  | +  | -   |
| abc    | $\bar{y}_{222}$ | +     | + | + | + | +  | +  | +  | +   |

# Variansanalystabell

- Om  $\theta = 0$  är  $E(2^k n \hat{\theta}^2) = \sigma^2$ .
- Effekter vi vet är noll kan användas för skattning av  $\sigma$
- Med  $s_{\theta}^2 = 2^k n \hat{\theta}^2$  är  $T = s_{\theta}^2 / s^2 \sim F(1, 2^k(n-1))$  om  $\theta = 0$ .
- $Q = S_{yy}$  kan delas upp i olika kvadratsummor baserat på de olika effekterna i modellen.
- Vi kan sammanfatta informationen om hur vi testar effekterna i ett faktorförsök i en så kallad variansanalystabell. Till exempel för ett  $2^2$ -försök:

| Variation | Kvadratsumma                                  | Frihetsgrader  | Medelkvadratsumma            | Teststorhet      |
|-----------|---|----------------|------------------------------|------------------|
| Faktor A  | $Q_A = 4n\hat{A}^2$                           | $f_A = 1$      | $s_A^2 = Q_A / f_A$          | $s_A^2 / s^2$    |
| Faktor B  | $Q_B = 4n\hat{B}^2$                           | $f_B = 1$      | $s_B^2 = Q_B / f_B$          | $s_B^2 / s^2$    |
| Faktor AB | $Q_{AB} = 4n\hat{AB}$                         | $f_{AB} = 1$   | $s_{AB}^2 = Q_{AB} / f_{AB}$ | $s_{AB}^2 / s^2$ |
| Residual  | $Q_0 = \sum_{ijk} (y_{ijk} - \bar{y}_{ij})^2$ | $f_R = 4(n-1)$ | $s^2 = Q_0 / f_R$            |                  |
| Total     | $Q = \sum_{ijk} (y_{ijk} - \bar{y})^2$        | $4n - 1$       |                              |                  |