

#### Simultaneous confidence interval

#### Bonferroni method

### $\alpha = P(X_c \ge 1|H_0), X_c \stackrel{H_0}{\sim} Bin(c, \alpha_c),$ where $X_c$ is the number of positive results among these c tests

 $P(X \ge 1|H_0) = 1 - (1 - \alpha_c)^c$ .

each individual test should be performed at the level  $\alpha/c$ .

 $\alpha = 1 - (1 - \alpha_c)^c \approx c\alpha_c$ . current setting of  $c = \binom{I}{2}$  pairwise differences  $(\mu_i - \mu_j)$ , studentised range distribution  $SR_{k_1,k_2}$ 

 $100(1-\alpha)\%$  simultaneous confidence interval: number of independent samples  $k_1$ 

 $B_{\mu_i-\mu_j} = \bar{y}_{i.} - \bar{y}_{j.} \pm t_{I(n-1)}(\frac{\alpha}{I(I-1)}) \cdot s_p \sqrt{\frac{2}{n}}, \quad | \text{Tukey's } 100(1-\alpha)\% \text{ simultaneous confidence interval} |$ 

e confidence interval for a single difference  $T_{\mu_i - \mu_j} = \bar{y}_{i.} - \bar{y}_{j.} \pm q_{I,I(n-1)}(\alpha) \cdot \frac{s_p}{\sqrt{n}}$ 

$$I_{\mu_i - \mu_j} = y_i - y_j \pm t_{I(n-1)}(\frac{\alpha}{2}) \cdot s_p \sqrt{\frac{\alpha}{n}}, \quad 1 \le \alpha \to \frac{\alpha}{c} = \frac{2\alpha}{I(I-1)}.$$

pairwise differences  $\mu_i - \mu_j$  are not independent as required by Bonferroni 1 9.1 Chi-squared test of homogeneity Bonferroni method gives slightly wider intervals compared to the Tukey

#### 9.2 Chi-squared test of independence

contrast to the previous setting of $J$ independent set		$b_1$	$b_2$		$b_J$	total
s, the total counts $(n_1, n_2, \dots, n_J)$ are random outcomes	$a_1$	$c_{11}$	$c_{12}$		$c_{1J}$	$c_1$
$(C_{11},, C_{IJ}) \sim \text{Mn}(n_{}; \pi_{11},, \pi_{IJ}),$	$a_2$	$c_{21}$	$c_{22}$	• • •	$c_{1J}$ $c_{2J}$	C2
characterised by $IJ-1$ degrees of freedom.	$a_I$	$c_{I1}$	$c_{I2}$		$c_{IJ}$	$c_I$
	total	$n_1$	$n_2$		$n_J$	n
U for all pains (i, i)						

 $H_0: \pi_{ij} = \pi_i.\pi_{\cdot j}$  for all pairs (i, j). maximum likelihood estimates of  $\pi_i$  and  $\pi_{ij}$ null hypothesis of independence the expected cell counts

$$e_{ij} = n\hat{\pi}_{ij} = n\hat{\pi}_i, \hat{\pi}_{\cdot j} = \frac{c_i n_j}{n} \text{ df} = (IJ-1) - (I-1+J-1) = (I-1)(J-1) \\ \text{maximum likelihood estimates of } \pi_i \in \frac{c_i n_j}{n} = c_i / n, \quad i = 1, \dots, I. \text{ estimates consumes } (I-1) \\ \text{The chi-squared tests of homogeneity and independence have the same test rejection rule.}$$

#### 10.1 Simple linear regression model

$$Y = \beta_0 + \beta_1 x + \sigma Z, \qquad Z \sim \mathcal{N}(0,1),$$

$$y_i = \beta_0 + \beta_1 x_i + e_i, \quad i = 1, \dots, n, \text{ where } C = (2\pi)^{-n/2}.$$

$$L(\theta) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}\right\} = C\sigma^{-n} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^{n} e_i^2\right\},$$

$$L(\theta) = \sum_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}\right\} = \frac{1}{2\sigma^2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^{n} e_i^2\right\},$$

$$\frac{1}{2\sigma^2} \exp\left\{-\frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}\right\} = \frac{1}{2\sigma^2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^{n} e_i^2\right\},$$

$$\begin{array}{l} & & \\ l(\theta) = \ln C - (n/2) \ln \sigma^2 - \frac{\sum_{i=1}^n e_i^2}{2\sigma^2} \quad \overline{x} = \frac{x_1 + \ldots + x_n}{n}, \quad \bar{y} = \frac{y_1 + \ldots + y_n}{n}, \\ e_i = y_i - \beta_0 - \beta_1 x_i, \quad \overline{x^2} = \frac{x_1^2 + \ldots + x_n^2}{n}, \quad \overline{y^2} = \frac{y_1^2 + \ldots + y_n^2}{n}, \end{array}$$

$$n^{-1} \sum_{i=1}^n e_i^2 = \beta_0^2 + 2\beta_0 \beta_1 \bar{x} - 2\beta_0 \bar{y} - 2\beta_1 \overline{xy} + \beta_1^2 \overline{x^2} + \overline{y^2},$$

$$\overline{xy} = \frac{x_1y_1 + \dots + x_ny_n}{n} \quad \frac{\partial l}{\partial \beta_0} = -\frac{n}{\sigma^2} (\beta_0 + \beta_1 \overline{x} - \overline{y}),$$

$$\frac{\partial l}{\partial \beta_1} = -\frac{n}{\sigma^2} (\beta_0 \overline{x} - \overline{xy} + \beta_1 \overline{x^2}),$$

$$\begin{aligned} &\hat{e}_i = y_i - b_0 - b_1 x_i, \\ &\mathbf{s}_{\mathrm{E}} = \hat{e}_1^2 + \ldots + \hat{e}_n^2, \end{aligned} \quad \frac{\partial l}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n e_i^2,$$

$$b_0 + b_1 \overline{x} = \overline{y}, \quad b_0 \overline{x} + b_1 \overline{x^2} = \overline{xy}, \quad \hat{\sigma}^2 = \frac{\mathrm{ss_E}}{n},$$

$$b_1 = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}, \qquad b_0 = \bar{y} - b_1\bar{x}.$$

$$\operatorname{ssp} = \min \left\{ \sum (y - \beta_0 - \beta_1 x_1)^2 \right\}$$

$$ss_{E} = \min_{\beta_{0},\beta_{1}} \left\{ \sum_{i=1} (y_{i} - \beta_{0} - \beta_{1}x_{i})^{2} \right\},$$
The maximum likelihood estimate of  $\hat{z}_{i}$ 

The maximum likelihood estimate of  $\hat{\sigma}^2 = \frac{ss_E}{r}$  is : biased but asymptotically unbiased estimate of  $\sigma^2$ unbiased estimate of  $\sigma^2$   $s^2 = \frac{ss_E}{s^2}$ 

#### 10.3 Multiple linear regression

$$y_1 = \beta_0 + \beta_1 x_{1,1} + \ldots + \beta_{p-1} x_{1,p-1} + e_1,$$

$$y_n = \beta_0 + \beta_1 x_{n,1} + \ldots + \beta_{p-1} x_{n,p-1} + e_n,$$

$$\mathbf{y} = (y_1, \dots, y_n)^\mathsf{T}, \quad \boldsymbol{\beta} = (\beta_0, \dots, \beta_{p-1})^\mathsf{T}, \quad \boldsymbol{e} = (e_1, \dots, e_n)^\mathsf{T},$$

$$\mathbb{X} = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,p-1} \\ \dots & \dots & \dots & \dots \\ 1 & x_{n,1} & \dots & x_{n,p-1} \end{pmatrix}.$$

$$\boldsymbol{b} = (\mathbb{X}^{\mathsf{T}}\mathbb{X})^{-1}\mathbb{X}^{\mathsf{T}}\mathbf{y},$$

:ted responses  $\hat{\mathbf{y}} = \mathbb{X}\mathbf{b}$ :

$$\hat{\boldsymbol{y}} = \mathbb{P}\boldsymbol{y}$$
, where  $\mathbb{P} = \mathbb{X}(\mathbb{X}^{\mathsf{T}}\mathbb{X})^{-1}\mathbb{X}^{\mathsf{T}}$   
 $p \times p$  matrix with elements  $Cov(B_i, B_j)$ .

$$E\{(\boldsymbol{B}-\boldsymbol{\beta})(\boldsymbol{B}-\boldsymbol{\beta})^{\mathsf{T}}\} = \sigma^{2}(\mathbb{X}^{\mathsf{T}}\mathbb{X})^{-1}.$$

$$\hat{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbb{I} - \mathbb{P})\mathbf{y}$$

a zero mean vector and a covariance matrix  $\sigma^2(\mathbb{I} - \mathbb{P})$ 

#### Denote by $d_0^2,\dots,d_{p-1}^2$ the diagonal elements of the matrix $(\mathbb{X}^\intercal\mathbb{X})^{-1}$

#### 8.3 Kruskal-Wallis test

non-parametric I independent random samples

 $X_i = \bar{Y}_i - \mu_i \sim \text{N}(0, \frac{\sigma}{\sqrt{n}}), \quad i = 1, \dots, I, \\ \text{pooled sample of size } N = I \cdot n. \sum_{k} \sum_{k} r_{ik} = 1 + 2 + \dots + N = \frac{N(N+1)}{2} \ \bar{r}. = \frac{N(N+1)}{2N} = \frac{N+1}{2}$  are independent normal variables.

Let  $r_{ik}$  be the pooled ranks  $_{
m measures}$  the discrepancy between the sample means of the ranks  $\frac{\max\{X_1, \dots, X_I\} - \min\{X_1, \dots, X_I\}}{S_p/\sqrt{n}} \sim SR_{k_1, k_2}, \quad k_1 = I, \ k_2 = I(n-1).$ 

we the pooled ranks measures the discrepancy between the sample means of the ranks 
$$w = \frac{12n}{N(N+1)} \sum_{i=1}^{I} (\bar{r}_{i.} - \frac{N+1}{2})^2 \quad \bar{r}_{i.} = \frac{r_{i1} + \dots + r_{in}}{n}, \quad i = 1, \dots, I.$$

A large value of w would indicate a deviation from the null distribution. degrees of freedom  $k_2$  used in the variance estimate  $s_p^2$ . For  $I=3,\,n\geq 5$  or  $I\geq 4,\,n\geq 4$ , one can use the approximate null distribution

## $W \stackrel{H_0}{pprox} \chi_{I-1}^2$ . Categorical data analysis

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{j}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{i}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{i}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{i}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{i}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{j}} = \bar{y}_{i}. - \bar{y}_{i}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{i}} = \bar{y}_{i}. \pm t_{I(n-1)}(\underline{\alpha}) + \frac{s_{p}}{\sqrt{n}}.$$

$$I_{\mu_{i}-\mu_{i}} = \bar$$

overall mean rank is

 $H_0: \pi_{i|j} = \pi_i$  for all pairs (i, j).

Tukey takes into account the dependences

$$\begin{aligned} &: \pi_{ij} = \pi_{i}.\pi_{.j} & \text{ for all pairs } (i,j). & & \dots & \dots & \dots & \dots \\ &: \pi_{i|j} = \mathrm{P}(A = a_i|B = b_j) = \frac{\pi_{ij}}{\pi_{.j}} & \frac{a_i}{\cot ||\pi_{.1}||} & \frac{\pi_{11}}{\pi_{.2}} & \dots & \frac{\pi_{JJ}}{\pi_{.J}} & \frac{\pi_{L}}{\pi_{.J}} \\ &: \pi_{.j} = \mathrm{P}(A = a_i|B = b_j) = \frac{\pi_{ij}}{\pi_{.j}} & \frac{a_i}{\cot ||\pi_{.1}||} & \frac{\pi_{21}}{\pi_{.2}} & \dots & \frac{\pi_{JJ}}{\pi_{.J}} & \frac{\pi_{L}}{\pi_{.J}} \end{aligned}$$

,		sample 1	sample 2	 sample $J$	total counts
	category $a_1$	$c_{11}$	$c_{12}$	 $c_{1J}$	$c_1$
	category $a_2$	$c_{21}$	$c_{22}$	 $c_{2J}$	$c_2$
-				 	
	category $a_I$	$c_{I1}$	$c_{I2}$	 $c_{IJ}$	$c_I$
	sample sizes	$n_1$	$n_2$	 $n_J$	n

#### J multinomial distributions

Tukev method

normality assumption with equal variances,

are independent normal variables.

 $(C_{1j}, \dots, C_{Ij}) \sim \text{Mn}(n_j; \pi_{1|j}, \dots, \pi_{I|j}), \quad j = 1, \dots, J.$ degrees of freedom for J independent samples from I-dimensional multinomial s J(I-1).  $H_0: \pi_{i|j} = \pi_i \text{ for all } (i,j), \text{ odds}(A) =$ 

 $\hat{\pi}_i = c_i/n, \quad i = 1, \dots, I.$  estimates consumes (I-1) degr expected cell counts  $e_{ij} = n_j \cdot \hat{\pi}_i = c_i n_j / n$ 

$$\mathbf{x}^2 = \sum_{i=1}^{I} \sum_{j=1}^{J} \frac{(c_{ij} - e_{ij})^2}{e_{ij}} = \sum_{i=1}^{I} \sum_{j=1}^{J} \frac{(c_{ij} - c_{i}n_{j}/n)^2}{c_{i}n_{j}/n}.$$

reject  $H_0$  for larger values of  $\mathbf{x}^2 \mathbf{X}^2 \approx \chi_{\mathrm{df}}^2$ 

# $\frac{\mathrm{P}(A)}{\mathrm{P}(A^c)} = \frac{\mathrm{P}(A)}{1 - \mathrm{P}(A)} \, \mathrm{P}(A) = \frac{\mathrm{odds}(A)}{1 + \mathrm{odds}(A)}, \, \mathrm{odds}(A|B) = \frac{\mathrm{P}(A|B)}{\mathrm{P}(A^c|B)} = \frac{\mathrm{P}(AB)}{\mathrm{P}(A^cB)}$

null hypothesis of homogeneity states the equality of J population distributions

9.3 Matched-pairs designs 9.4 Odds ratios

odds ratio for a pair of events 
$$(A, B)$$
  $\Delta_{AB} = \frac{\text{odds}(A|B)}{\text{odds}(A|B^c)},$ 

$$\Delta_{AB} = \frac{P(AB)P(A^cB)^c}{P(A^cB)P(AB^c)}, \quad \Delta_{AB} = \Delta_{BA}, \quad \Delta_{AB^c} = \frac{1}{\Delta_{AB}}.$$

if  $\Delta_{AB}=1$ , then events A and B are independent, if  $\Delta_{AB}>1$ , then  $\mathrm{P}(A|B)>\mathrm{P}(A|B^c)$  so that B increases the probability of A, if  $\Delta_{AB}<1$ , then  $\mathrm{P}(A|B)<\mathrm{P}(A|B^c)$  so that B decreases the probability of A

 $ss_R = \sum (\hat{y}_i - \bar{y})^2 = (n-1)b_1^2 s_x^2 = (n-1)r^2 s_y^2$  regression sum of squares.

#### df = J(I-1) - (I-1) = (I-1)(J-1). Coefficient of determination

**Residuals** 
$$\hat{y}_i = b_0 + b_1 x_i$$
,  $ss_T = ss_R + ss_E$ ,  $ss_T = \sum_i (y_i - \bar{y})^2 = (n-1)s_y^2$  the total sum of squares,

$$\hat{e}_i = y_i - \hat{y}_i$$
,  $\hat{e}_1 + \dots + \hat{e}_n = 0$ ,  
 $x_1\hat{e}_1 + \dots + x_n\hat{e}_n = 0$ ,  $Cov(\hat{E}_i, \hat{E}_i) = -\sigma^2 \frac{c_{ij}}{c_{ij}}$ .

$$x_1\hat{e}_1 + \ldots + x_n\hat{e}_n = 0, \quad \text{Cov}(\hat{E}_i, \hat{E}_j) = -\sigma^2 \frac{c_{ij}}{n-1}, \\ \hat{y}_1\hat{e}_1 + \ldots + \hat{y}_n\hat{e}_n = 0.$$

$$Var(\hat{E}_{i}) = \sigma^{2} \left( 1 - \frac{c_{ii}}{n-1} \right), \quad c_{ij} = \frac{\sum_{k=1}^{n} (x_{i} - x_{k})(x_{j} - x_{k})}{ns_{x}^{2}}.$$

$$\begin{aligned} & \textbf{Sample correlation coefficient} & s^2 = \frac{n-1}{n-2} \, s_y^2 (1-r^2). \\ & s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \ s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) (y_i - \bar{y}). \\ & \textbf{10.2 Confidence intervals and hypothesis testing} \\ & s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2, \ \text{sample correlation coefficient} \\ & s_{xy} \quad \text{to } S_y \quad \text{to }$$

$$s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$
, sample correlation coefficier  $y = \bar{y} + r \frac{s_y}{s_x} (x - \bar{x})$ ,  $r = \frac{s_{xy}}{s_x s_y}$ .  $b_1 = \frac{r s_y}{s_x}$ ,  $s_y = \frac{r s_y}{s_x}$ ,

$$y = \bar{y} + r \frac{s_y}{s_x} (x - \bar{x}), \quad r - \frac{1}{s_x s_y}. \quad \sigma_1 - \left(\frac{y - \bar{y}}{s_y}\right) = r\left(\frac{x - \bar{x}}{s_x}\right).$$

Intervals for individual observations is if 
$$y_p = \beta_0 + \beta_1 x_p + \sigma z_p$$
,  $z_p$  is generated by the N(0,1)  $\mu_p = \beta_0 + \beta_1 x_p$ 

$$\begin{split} B_0 &\sim \mathcal{N}(\beta_0, \sigma_0), \qquad \sigma_0^2 = \frac{\sigma^2 \sum x_i^2}{n(n-1)s_x^2}, \qquad s_{b_0}^2 = \frac{s^2 \sum x_i^2}{n(n-1)s_x^2}, \qquad \frac{B_0 - \beta_0}{S_{B_0}} \sim t_{n-2}, \\ B_1 &\sim \mathcal{N}(\beta_1, \sigma_1), \qquad \sigma_1^2 = \frac{\sigma^2}{(n-1)s_x^2}, \qquad s_{b_1}^2 = \frac{s^2}{(n-1)s_x^2}, \qquad \frac{B_1 - \beta_1}{S_{B_1}} \sim t_{n-2}. \\ \mathrm{Cov}\left(B_0, B_1\right) = -\frac{\sigma^2 \bar{x}}{(n-1)s_x^2} & \text{which is negative, if } \bar{x} > 0, \text{ and positive, if } \bar{x} < 0. \\ H_0 : \quad \beta_i = \beta^*. \end{split}$$

 $\frac{\text{SSR}}{\text{SS}_{\text{T}}} = r^2, \qquad \frac{\text{SSE}}{\text{SS}_{\text{T}}} = 1 - r^2.$  the squared sample correlation coefficient  $r^2$  is called the coefficient of determination.

 $ss_E = ss_T(1 - r^2) = (n - 1)s_u^2(1 - r^2)$ , unbiased estimate of  $\sigma^2$ :

 $\hat{\mu}_p = b_0 + b_1 x_p$ . prediction interval of the response value  $y_p$ 

$$Var(B_0 + B_1 x_p) = Var(B_0) + x_p^2 Var(B_1) + 2x_p Cov(B_0, B_1) = \frac{\sigma^2}{n} + \frac{\sigma^2}{n-1} \cdot (\frac{x_p - \bar{x}}{s_p})^2$$
.

$$I_{\mu_p} = b_0 + b_1 x_p \pm t_{n-2} (\frac{\alpha}{2}) \cdot s \sqrt{\frac{1}{n} + \frac{1}{n-1} (\frac{x_p - \bar{x}}{s_x})^2}$$
. 10.4 Coefficients of multiple determination and model utility tests

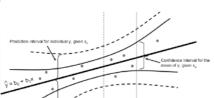
$$(\frac{-\bar{x}}{x})^2$$
. 10.4 Coefficients of multiple determination and model utility tests  $R^2 = 1 - \frac{\text{Ss}_E}{x}$ ,  $\text{Ss}_T = (n-1)s_x^2$ ,  $R^2 = 1 - \frac{n-1}{x}$ .

$$\operatorname{Var}(Y_p - \hat{\mu}_p) = \operatorname{Var}(\mu_p + \sigma Z_p - \hat{\mu}_p) =$$

$$R^2 = 1 - \frac{\mathrm{ss_E}}{\mathrm{ss_T}}, \qquad \mathrm{ss_T} = (n-1)s_y^2. \ \ R_a^2 = 1 - \frac{n-1}{n-p} \cdot \frac{\mathrm{ss_E}}{\mathrm{ss_T}}$$

$$= \sigma^2 + \text{Var}(\hat{\mu}_p) = \sigma^2 (1 + \frac{1}{n} + \frac{1}{n-1} \cdot (\frac{x_p - \bar{x}}{s_x})^2).$$
Prediction Interval vs. Confidence Interval

$$R_a^2 = 1 - \frac{s^2}{s_y^2},$$



confidence interval goes to zero, while the width of the 95% prediction interval converges to  $1.96\sigma$ 

standard error of 
$$b_i$$
  $s_{b_j} = sd_j$ ,  $\frac{B_j - \beta_j}{S_{B_j}} \sim t_{n-p}, \quad j = 0, 1, \dots, p-1$ .

## Parametric models

 $\mathbb{E}(\bar{X}) = \mu, \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}, \quad \mathbb{E}(S^2) = \sigma^2, \quad \text{Var}(S^2) = \frac{\sigma^4}{n} \left( \mathbb{E}(\frac{X - \mu}{\sigma})^4 - \frac{n - 3}{n - 1} \right). \quad \mathbf{Gamma, 1} \quad \mathbf{Gam}(\alpha, \lambda) \quad \mu = \frac{\alpha}{\lambda}, \quad \sigma^2 = \frac{\alpha}{\lambda^2}.$  $\mu = \mathrm{E}(X), \quad \sigma^2 = \mathrm{Var}(X). \quad \mathrm{Var}(X) = \mathrm{E}((X-\mu)^2). \quad \mathcal{F}(\mu,\sigma) = \mathrm{N}(\mu,\sigma), \quad \frac{\bar{X}-\mu}{S/\sqrt{n}} \sim t_{n-1}$ 

$$\begin{split} & E(X) = \int_{-\infty}^{\infty} x f(x) dx \quad \text{or} \quad E(X) = \sum_{i=1}^{\infty} x_i p_i, \\ & z\text{-score}, \quad \text{Cov}(X_1, X_2) = E((X_1 - \mu_1)(X_2 - \mu_2)). \end{split} \\ & z = \frac{X - \mu}{2} \quad \text{F}(X) = \frac{x}{2} \quad \text{Tr}(X_1 - \mu_1)(X_2 - \mu_2). \end{split}$$
gamma function  $\Gamma(a) = \int_{-\infty}^{\infty} x^{a-1}e^{-x}dx$ ,  $Z = \frac{X - \mu}{\sigma}$ ,  $E(X) = \mu$ , variance  $Var(X) = \sigma^2$ ,

 $X \sim N(\mu, \sigma) \quad f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^{-}}{2\sigma^{2}}}, \quad -\infty < x < \infty. \quad \Gamma(k) = (k-1)!, \quad k = 1, 2, \dots$   $\Phi(x) = P(Z \le x), \quad -\infty < x < \infty. \quad \text{if } Z, Z_{1}, \dots, Z_{k} \text{ are independent random variables with } N(0, 1) \text{-distribution},$   $P(X \le x) = P(\frac{X-\mu}{\sigma} \le \frac{x-\mu}{\sigma}) = P(Z \le \frac{x-\mu}{\sigma}) = \Phi(\frac{x-\mu}{\sigma}). \quad \frac{Z}{\sqrt{(Z_{1}^{2} + \dots + Z_{k}^{2})/k}} \sim t_{k}.$ 

 $\bar{X} = \frac{X_1 + \ldots + X_n}{n} \ \bar{X} \approx N(\mu, \frac{\sigma}{\sqrt{n}}), \quad 1.5$  Bernoulli, binomial, and multinomial distributions

$$\sigma^{2} = \sum_{j=1}^{k} w_{j} (\mu_{j} - \mu)^{2} + \sum_{j=1}^{k} w_{j} \sigma_{j}^{2}.$$

variation within the strata  $\sum_{j=1}^{k} w_j \sigma_i^2$ .

Bin
$$(n, p)$$
.  $P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, ..., n,$   
 $\mu = np, \quad \sigma^2 = np(1-p). \quad \text{Bin}(n, p) \approx N(np, \sqrt{np(1-p)}).$ 

$$P(X = x) = \frac{\mu^x}{r!}e^{-\mu}, \quad x = 0, 1, ...,$$
  
 $E(X) = \mu, \quad Var(X) = \mu.$ 

arule of thumb good enough it both 
$$np \ge 3$$
 and  $n(1-p) \ge 3$ , 
$$P(X \le x) \approx \Phi\left(\frac{x-np}{\sqrt{np(1-p)}}\right).$$
 Continuity correction suggests replacing  $x$  by  $x+\frac{1}{2}$  on the right hand side 
$$P(X \le x) = P(X < x+1), \quad P(X \le x) \approx \Phi\left(\frac{x-np}{\sqrt{np(1-p)}}\right), \quad \text{Hype}$$

Poisson distribution is obtained as an approximation for the Bin(n, p) distribution i

 $n \to \infty$ ,  $p \to 0$ , and  $np \to \mu$ .

describe the number of rear events (like accidents) observed during a given time interval.

#### Geometric distribution

sequence of independent Bernoulli trials with probability p of success.

the distribution of the number X of trials needed to get one success, the distribution of the number Y = X - 1 of failures before the first success.

$$X \sim \text{Geom}(p)$$
  $P(X=x) = (1-p)^{x-1}p, \quad x=1,2,\ldots, \quad \mu=\frac{1}{n}, \quad \sigma^2=\frac{1-p}{n^2}$  Beta distribution

Beta distribution  $\operatorname{Beta}(a,b)$  is determined by two parameters  $a>0,\,b>0$  which are called pseudo-counts. It is defined by the probability density function

$$g(p) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} p^{a-1} (1-p)^{b-1}, \quad 0$$

its mean and variance are given by

$$\mu=rac{a}{a+b},\quad \sigma^2=rac{\mu(1-\mu)}{a+b+1}.$$

#### 2.1 Point estimation

## mean square error $E((\widehat{\Theta} - \theta)^2) = Var(\widehat{\Theta}) + (E(\widehat{\Theta}) - \theta)^2$

 $\text{mean square error vanishes } \mathbb{E}((\hat{\theta} - \theta)^2) \rightarrow 0 \text{ as } n \rightarrow \infty, \underbrace{\text{the point}}_{} \text{ estimate } \hat{\theta} \text{ is called } \underbrace{consistent.}_{} \mathbb{E}(S^2) = \underbrace{\frac{n}{n-1}}_{} \mathbb{E}\left(\frac{\sum_{i=1}^n X_i^2}{-\bar{X}^2} - \bar{X}^2\right) = \frac{n}{n-1}(\mathbb{E}(X^2) - \mathbb{E}(\bar{X}^2))$ 

 $I_{\mu} \approx \bar{x} \pm z(\frac{\alpha}{2}) \cdot \frac{s}{\sqrt{n}}$ 

#### 2.6 Exact confidence intervals

$$I_{\mu} = \bar{x} \pm t_{n-1} \left(\frac{\alpha}{2}\right) \cdot s_{\bar{x}},$$

$$I_{\sigma^2} = \Big(\frac{(n-1)s^2}{x_{n-1}(\frac{\alpha}{2})}; \frac{(n-1)s^2}{x_{n-1}(1-\frac{\alpha}{2})}\Big)$$

$$f(x) = \frac{\Gamma(\frac{k+1}{2})}{\sqrt{k\pi}\Gamma(\frac{k}{2})} \left(1 + \frac{x^2}{k}\right)^{-\frac{k+1}{2}}, -\infty < x < \infty$$

gamma function 
$$\Gamma(a) = \int_0^\infty x^{a-1} e^{-x} dx$$

if 
$$Z, Z_1, ..., Z_k$$
 are independent random variables with N(0, 1)-distribution

if 
$$Z,Z_1,\dots,Z_k$$
 are independent random variables with N(0,1)-distribution,

 $\mbox{Mixtures of normal distribution} \quad \mbox{Bernoulli, } P(X=1) = p, \quad P(X=0) = 1-p.$ 

binomial, 
$$X \sim Bin(1, p)$$
  $\mu = p$ ,  $\sigma^2 = p(1 - p)$ .

**Poisson,**  $X \sim \text{Pois}(\mu)$  is a discrete distribution rule of thumb good enough if both  $np \geq 5$  and  $n(1-p) \geq 5$ ,

 $P(X \le x) \approx \Phi\left(\frac{x + \frac{1}{2} - np}{\sqrt{np(1-p)}}\right),$ 

$$P(X < x) \approx \Phi\left(\frac{x - \frac{1}{2} - np}{\sqrt{np(1-p)}}\right).$$

$$P(X < x) \approx \Phi\left(\frac{x - \frac{1}{2} - np}{\sqrt{np(1-p)}}\right).$$

Simple random sample

$$X$$
 is the number of successes in  $n$  Bernoulli trials which depend on each other. 
$$\mathrm{P}(X=x) = \frac{\binom{B}{x}\binom{W}{n-x}}{\binom{N}{n}}, \quad \mu = np, \quad \sigma^2 = np(1-p)\frac{N-n}{N-1}.$$
 Compared to the variance of the  $\mathrm{Bin}(n,p)$  distribution, the last formula contains the factor

B = Np balls are black and

W = N(1 - p) balls are white.

 $\frac{N-n}{N-1} = 1 - \frac{n-1}{N-1}$ , finite population correction factor.

 $Hg(N, n, p) \approx N(\mu, \sigma), \quad \mu = np, \quad \sigma = \sqrt{np(1-p)}\sqrt{1 - \frac{n-1}{N-1}},$ 

applied provided  $np \ge 5$  and  $n(1-p) \ge 5$ .

2.4 Dichotomous data

The estimated standard error for the sample proportion  $\hat{p}$  is  $s_{\hat{p}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1}}$ Simple random sample  $s_{\hat{x}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \sqrt{1-\frac{n}{n}}$ 

shape parameter  $\, \alpha > 0 \, ;$  inverse scale (or the rate) parameter  $\, \lambda > 0 \, .$ 

 $\alpha \gg 1$  chi-squared

if  $\alpha = 1$ ,  $Gam(1, \lambda) = Exp(\lambda)$ . exponential,

chi-squared distribution with k degrees of freedom is the gamma distribution with  $\alpha = \frac{k}{2}, \lambda = \frac{1}{2}$ .

n  $\omega_1,\dots,\omega_k$  are independent random variables with N(0,1)-distribution  $Z_1^2+\dots+Z_k^2\sim\chi_k^2$ 

if  $(X_1,\ldots,X_n)$  are independent and random variables each having the  $N(\mu,\sigma)$ 

multinomial distribution  $(X_1, \ldots, X_r) \sim \operatorname{Mn}(n; p_1, \ldots, p_r)$ 

$$\begin{split} \mathbf{P}(X_1 = x_1, \dots, X_r = x_r) &= \binom{n}{x_1, \dots, x_r} p_1^{x_1} \dots p_r^{x_r}, x_i = 0, \dots, n, \quad i = 1, \dots, r, \\ p_1 + \dots + p_r &= 1. \quad \text{marginal distribution of } X_i \text{ is binomial } X_i \sim \text{Bin}(n, p_i) \end{split}$$

if  $X_i \sim \text{Exp}(\lambda)$ , i = 1, ..., k are independent,

 $X_1 + \ldots + X_k \sim \operatorname{Gam}(k, \lambda), \quad k = 1, 2, \ldots$ 

restricted to positive values,

large values of the shape parameter  $Gam(\alpha, \lambda) \approx N(\frac{\alpha}{\lambda}, \frac{\sqrt{\alpha}}{\lambda}),$ 

if  $Z_1, \ldots, Z_k$  are independent

 $Cov(X_i, X_j) = -np_i p_j, \quad i \neq j.$ 

Hypergeometric distribution  $X \sim \text{Hg}(N, n, p)$ 

replacement from a box with N balls, of which

number x of black balls among n balls drawn without

 $\frac{\sum_{i=1}^{n} (X_i - \bar{X})^2}{\sigma^2} \sim \chi_{n-1}^2$ .

An approximate  $100(1-\alpha)\%$  two-sided confidence interval for p is given by  $I_p \approx \hat{p} \pm z(\frac{\alpha}{2}) \cdot s_{\hat{p}}$ 

- sampling with replacement produces what we call a random sample consisting of independent and identicall valid both for sampling with replacement and for sampling distributed observations, the bias size  $E(\widehat{\Theta}) - \theta$ , measuring the lack of accuracy (systematic error), sampling without replacement produces a so called simple random sample having identically distributed but dependent observations.

$$= \frac{\frac{n}{n-1} E\left(\frac{D(n-1)}{n} - X^2\right) = \frac{n}{n-1} (E(X^2) - E(X^2))}{= \frac{n}{n-1} (\sigma^2 + \mu^2 - \frac{\sigma^2}{n} (1 - \frac{n-1}{N-1}) - \mu^2) = \frac{n}{n-1} (\sigma^2 - \frac{\sigma^2}{n} (1 - \frac{n-1}{N-1})) = \sigma^2 \frac{N}{N-1}, \sigma^2 = \overline{\sigma^2} + \sum_{\substack{i=1 \ N=1}}^{k} w_j (\mu_j - \mu)^2.$$

 $s_{\bar{x}}^2=\frac{s^2}{n}\frac{N-1}{N}(1-\frac{n-1}{N-1})=\frac{s^2}{n}(1-\frac{n}{N}).$  unbiased estimate of  $\mathrm{Var}(\bar{X})$ :

sampling without replacement, the formula for the estimated standard error of the sample mean 
$$\bar{x}$$

$$s_{\bar{x}} = \frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}. \ I_{\mu} \approx \bar{x} \pm z (\alpha/2) s_{\bar{x}},$$

$$Var(\bar{X}_s) = w_1^2 Var(\bar{X}_1) + ... + w_k^2 Var(\bar{X}_k) = \frac{w_1^2 \sigma_1^2}{n_1} + ... + \frac{w_k^2 \sigma_k^2}{n_k}.$$

$$s_{\bar{x}}^2 = w_1^2 s_{\bar{x}_1}^2 + ... + w_k^2 s_{\bar{x}_k}^2 = \frac{w_1^2 s_1^2}{n_k} + ... + \frac{w_k^2 s_k^2}{n_k}.$$

A confidence interval for the mean based on a stratified sample:  $I_{\mu} \approx \bar{x}_s \pm z(\frac{\alpha}{2}) \cdot s_{\bar{x}_c}$ 

$$\bar{\sigma} = w_1 \sigma_1 + \ldots + w_k \sigma_k$$

The stratified sample mean  $\bar{X}_{so}$  with the optimal allocation  $n_i = n \frac{w_j \sigma_j}{\pi}$  has the smallest variance  $Var(\bar{X}_{so}) = \frac{\bar{\sigma}^2}{n}$  among all allocations of n observations.

The stratified sample mean  $\bar{X}_{sp}$  for the proportional allocation  $n_j = nw_j$  has the variance  $\text{Var}(\bar{X}_{sp}) = \frac{\overline{\sigma^2}}{n}$ .

$$Var(\bar{X}_{so}) \le Var(\bar{X}_{sp}) \le Var(\bar{X}),$$

$$\frac{(\bar{\sigma})^2}{n} \leq \frac{\overline{\sigma^2}}{n} \leq \frac{\sigma^2}{n}$$

The statistical tools introduced in this course so far are based on the so called frequentist approach. In the parametric case, the frequentist treats the data x as randomly generated by a distribution  $f(x|\theta)$  involving the unknown true population parameter value  $\theta$ , which may be estimated using the method of maximum likelihood. This section presents basic concepts of the Bayesian approach relying on the following model for the observed data x:

A  
priori distribution 
$$\xrightarrow{g(\theta)}$$
 generates a value  $\theta \xrightarrow{f(x|\theta)}$  data  
  $x$ .

The model assumes that before the data is collected the parameter of interest  $\theta$  is randomly generated by a prior distribution  $g(\theta)$ . The computational power of the Bayesian approach stems from the possibility to treat  $\theta$  as a realisation of a random variable  $\Theta$ .

The prior distribution  $g(\theta)$  brings into the statistical model our knowledge (or lack of knowledge) on  $\theta$  before the data x is generated using a conditional distribution  $f(x|\theta)$ , which in this section is called the likelihood function. After the data x is generated by such a two-step procedure involving the pair  $g(\theta)$  and  $f(x|\theta)$ , we may update our knowledge on  $\theta$  and compute a posterior distribution  $h(\theta|x)$  using the Bayes formula

$$h(\theta|x) = \frac{f(x|\theta)g(\theta)}{\phi(x)}$$
.

The denominator, depending on whether the distribution is continuous or discrete,

$$\phi(x) = \int f(x|\theta)g(\theta)d\theta$$
 or  $\phi(x) = \sum_{\alpha} f(x|\theta)g(\theta)$ 

gives the marginal distribution of the random data X. For a fixed realization x, treating the denominator  $\phi(x)$  as a constant which does not explicitly involve  $\theta$ , the Bayes formula can be summarized as

posterior 
$$\propto$$
 likelihood  $\times$  prior

where the sign  $\propto$  means proportional.

If we have no prior knowledge on  $\theta$ , the prior distribution is often modelled by the uniform distribution. In this case of uninformative prior, with  $g(\theta)$  being a constant over a certain interval, we have  $h(\theta|x) \propto f(x|\theta)$ , implying that the posterior knowledge comes solely from the likelihood function.

#### Multinomial-Dirichlet model

The multinomial-Dirichlet model is a multivariate version of the binomial-beta model. For both the binomial-beta and multinomial-Dirichlet models, the updating rule has the form

the posterior pseudo-counts = the prior pseudo-counts plus the sample counts

#### Dirichlet distribution

The Dirichlet distribution  $Dir(\alpha_1, \dots, \alpha_r)$  is a multivariate extension of the beta distribution. It is a probability distribution over the vectors  $(p_1,\dots,p_r)$  with non-negative components such that

$$p_1 + \ldots + p_r = 1$$
.

The positive parameters  $\alpha_1, \dots, \alpha_r$  of the Dirichlet distribution are often called the pseudo-counts. The probability density function of  $Dir(\alpha_1, ..., \alpha_r)$  is given by

$$g(p_1,\ldots,p_r) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)\ldots\Gamma(\alpha_r)}p_1^{\alpha_1-1}\ldots p_r^{\alpha_r-1}, \quad \alpha_0 = \alpha_1+\ldots+\alpha_r.$$

The marginal distributions of the random vector  $(X_1,\ldots,X_r)\sim \mathrm{Dir}(\alpha_1,\ldots,\alpha_r)$  are the beta distributions

$$X_j \sim \text{Beta}(\alpha_j, \alpha_0 - \alpha_j), \quad j = 1, \dots, r.$$

Different components of the vector have negative covariances

$$Cov(X_i, X_j) = -\frac{\alpha_i \alpha_j}{\alpha_0^2(\alpha_0 + 1)}$$
 for  $i \neq j$ .

The figure below illustrates four examples of  $\mathrm{Dir}(\alpha_1,\alpha_2,\alpha_3)$  distribution. Each triangle contains n=300 points generated using different sets of parameters  $(\alpha_1, \alpha_2, \alpha_3)$ :

upper left (0.3, 0.3, 0.1), upper right (13, 16, 15), lower left (1, 1, 1), lower right (3, 0.1, 1).

A dot in a triangle gives a realisation  $(x_1, x_2, x_3)$  of the vector  $(X_1, X_2, X_3) \sim \text{Dir}(\alpha_1, \alpha_2, \alpha_3)$  as the distances to the bottom edge of the triangle  $(x_1)$ , to the right edge of the triangle  $(x_2)$ , and to the left edge of the triangle  $(x_3)$ .

#### 5.3 Credibility interval

Given data x coming from a parametric model with the likelihood function  $f(x|\theta)$ , a  $100(1-\alpha)\%$  confidence interval for the parameter  $\theta$ ,

$$I_{\theta} = (a_1(x), a_2(x)),$$

48

is viewed as a realisation of a random interval  $(a_1(X), a_2(X))$  such that

$$P(a_1(X) < \theta < a_2(X)) = 1 - \alpha.$$

This frequentist interpretation of the confidence level  $100(1-\alpha)\%$  is rather cumbersome as it requires mentioning other samples and potential confidence intervals which as a group cover the true unknown value of  $\theta$  with probability

In the framework of Bayesian inference we can refer to  $\theta$  as a realisation of a random variable  $\Theta$  with a certain posterior distribution  $h(\theta|x)$ . This allows us to define a  $100(1-\alpha)\%$  credibility interval (or credible interval)

$$J_\theta = (b_1(x),b_2(x))$$

by the relation based on the posterior distribution

$$P(b_1(x) < \Theta < b_2(x)|x) = 1 - \alpha.$$

The interpretation of the credibility interval is more intuitive as it does not refer to some potential, never observed

#### 5.1 Conjugate priors

Suppose the data x is generated by a parametric model having the likelihood function  $f(x|\theta)$ . Consider a parametric family of the prior distributions G.

$$\mathcal G$$
 is called a family of conjugate priors for the likelihood function  $f(x|\theta)$  if for any prior  $g(\theta)\in\mathcal G,$  the corresponding posterior distribution  $h(\theta|x)\in\mathcal G$ 

The next table presents five Bayesian models involving conjugate priors. The details of the first three models come next. Notice that the posterior variance is always smaller than the prior variance. This list also illustrates that the contribution of the prior distribution to the posterior distribution becomes smaller as the sample size n increase.

Parametric model for the data	Unknown $\theta$	Prior	Posterior distribution
$X_1, \dots, X_n \sim N(\mu, \sigma)$	$\theta = \mu$	$N(\mu_0, \sigma_0)$	$N(\gamma_n \mu_0 + (1 - \gamma_n)\bar{x}; \sigma_0 \sqrt{\gamma_n})$
$X \sim Bin(n, p)$	$\theta = p$	Beta(a, b)	Beta(a + x, b + n - x)
$(X_1,, X_r) \sim Mn(n; p_1,, p_r)$	$\theta = (p_1,, p_r)$	$Dir(\alpha_1,, \alpha_r)$	$Dir(\alpha_1 + x_1,, \alpha_r + x_r)$
$X_1, \dots, X_n \sim \text{Geom}(p)$	$\theta = p$	Beta(a, b)	$Beta(a + n, b + n\bar{x} - n)$
$X_1, \dots, X_n \sim Pois(\mu)$	$\theta = \mu$	$Gam(\alpha_0, \lambda_0)$	$Gam(\alpha_0 + n\bar{x}, \lambda_0 + n)$
$X_1, \dots, X_n \sim \text{Gam}(\alpha, \lambda)$	$\theta = \lambda$	$Gam(\alpha_0, \lambda_0)$	$Gam(\alpha_0 + \alpha n, \lambda_0 + n\bar{x})$

#### Normal-normal model

Suppose a random sample  $(x_1, \dots, x_n)$  is drawn from the normal distribution  $N(\mu, \sigma)$  with a known standard deviation  $\sigma$  and the unknown mean  $\theta = \mu$ . Taking the normal prior  $\Theta \sim N(\mu_0, \sigma_0)$  with known  $(\mu_0, \sigma_0)$  results in the normal posterior  $N(\mu_1, \sigma_1)$  with

$$\mu_1 = \gamma_n \mu_0 + (1 - \gamma_n)\bar{x}, \quad \sigma_1^2 = \sigma_0^2 \gamma_n,$$

$$\gamma_n = \frac{\sigma^2}{\sigma^2 + n\sigma_0^2} = \frac{\frac{\sigma^2}{n}}{\frac{\sigma^2}{n} + \sigma_0^2}$$

posterior mean  $\mu_1$  gets close to the maximum likelihood estimate and the input  $\gamma_{n}\mu_0$  involving the prior mean becomes negligible. is the shrinkage factor which becomes smaller for the larger sample sizes n. As a result for the large samples, the

#### Binomial-beta model

Next, we introduce the beta distribution which serves as a convenient family of conjugate priors for Bayesian inference for p, in the case when the data x is generated by the Bin(n, p).

#### 5.2 Bayesian estimation

In the language of decision theory, finding a point estimate a for the unknown population parameter  $\theta$  is an action of assigning the value a to the unknown parameter  $\theta$ . In the frequentist setting, the optimal a is found by maximising the likelihood function. In the Bayesian setting, the optimal choice of a is determined by the so-called loss function  $l(\theta,a)$ . The so-called Bayes action minimises the posterior risk

$$R(a|x) = E(l(\Theta, a)|x),$$

computed using the posterior distribution

$$R(a|x) = \int l(\theta,a)h(\theta|x)d\theta \quad \text{ or } \quad R(a|x) = \sum_{\alpha} l(\theta,a)h(\theta|x).$$

We consider two loss functions leading to two different Bayesian estimators. These two loss functions called the zero-one loss and the squared error loss

$$\boxed{\text{Zero-one loss function: } l(\theta,a) = 1_{\{\theta \neq a\}}} \\ \boxed{\text{Squared error loss: } l(\theta,a) = (\theta-a)^2}$$

are schematically depicted on the figure below.

#### Zero-one loss function and maximum a posteriori probability

With the zero-one loss function, the posterior risk is equal to the probability of misclassification

$$R(a|x) = \sum_{\theta \neq a} h(\theta|x) = 1 - h(a|x).$$

In this case, to minimise the risk we have to maximise the posterior probability h(a|x). We define In this case, to minimise the risk we have to minimise the risk we have of  $\hat{\theta}$  that maximises  $h(\hat{\theta}|x)$ . Observe that with the uninformative prior,  $\hat{\theta}_{map} = \hat{\theta}_{mle}$ .

define  $\hat{\theta}_{map}$  as the value

#### Squared error loss function and posterior mean estimate

Using the squared error loss function we find that the posterior risk is a sum of two components

$$R(a|x) = E((\Theta - a)^{2}|x) = Var(\Theta|x) + (E(\Theta|x) - a)^{2}.$$

Since the first component is independent of a, we minimise the posterior risk by putting

$$\hat{\theta}_{pme} = E(\Theta|x),$$

resulting in the posterior mean value as the Bayesian point estimate of  $\theta.$