

NMAI059 Probability and statistics 1

Class 13

Robert Šámal

Statistics – What have we learnt

- ▶ basic setup: we consider random sample X_1, \dots, X_n from distribution F_ϑ — describes the measurement process, all ways how could it go
 - ▶ we measure data – particular numbers x_1, \dots, x_n , so called realization of random sample, — what did we really measure
1. point estimation: determine best possible number, estimate for the parameter ϑ , or some function of it, $g(\vartheta)$.
 2. interval estimation: determine an interval, that contains the unknown parameter ϑ with a large probability
 3. hypothesis testing — *best fixed*

Overview

Hypothesis testing

Goodness of fit test

Linear regression

Hypothesis testing – illustration

- ▶ We want to test, if a coin is fair.
- ▶ H_0 : it is fair
- ▶ H_1 : not fair (“Scientists discovered, that casino XY uses loaded coin.”)
- ▶ Results: Reject H_0 /don't reject H_0
- ▶ Type I error: false rejection. We reject H_0 , even if it is true. Embarrassing.
- ▶ Type II error: false non-rejection. We don't reject H_0 , even if it is false. Unused opportunity.
- ▶ Need to find k such that we will reject H_0 if $|S - n/2| > k$.

foss coin n-tosses

$$n = \underline{1000}$$

$$S := \underline{\# \text{ of heads}}$$

$$S = \underline{200} \quad p = 0$$

Hypothesis testing – general approach

$X_i \rightarrow X_i \sim \text{Ber}(p), \text{ i.i.d}$
 $X = X_1, \dots, X_n \sim \text{Bern}(n, p)$

- ▶ We choose an appropriate statistical model.
- ▶ We choose *significance level* α : prob. of false rejection of H_0 . Typically $\alpha = 0.05$ (medicine/psychology – much less in high-energy physics).
- ▶ We determine *test statistics* $T = h(X_1, \dots, X_n)$, that we will determine from the measured data.
 $\Rightarrow \sum X_i \sim \text{Bern}(n, p)$
- ▶ We determine *rejection region* – set W .
 $W = \{n: \frac{n}{n} = 1 > k\}$
- ▶ We measure x_1, \dots, x_n – so-called realizations of X_1, \dots, X_n .
- ▶ Decision rule: we reject H_0 iff $h(x_1, \dots, x_n) \in W$.
- ▶ $\alpha = P(h(X) \in W; H_0)$ $P(\text{Type I error})$
- ▶ $\beta = P(h(X) \notin W; H_1)$ α . strength of the test $1 - \beta$
 $P(\text{Type II error})$
- ▶ often we do not choose α in advance but compute so-called *p-value*: minimal α , for which we would reject H_0 .

Hypothesis testing – an example

- ▶ X_1, \dots, X_n random sample from $N(\vartheta, \sigma^2)$
- ▶ σ^2 known
- ▶ $H_0 : \vartheta = 0$ $H_1 : \vartheta \neq 0$

mean temp.

mean = false temp.

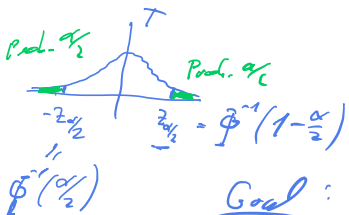
$$\alpha = 0.05$$

Assume $H_0 : \vartheta = 0$

$$Z := \bar{X}_n \sim N(\vartheta, \sigma^2/n)$$

$$T := \frac{\bar{X}_n}{\sigma/\sqrt{n}} \sim N(0, 1)$$

the disto. of test stat., assume H_0 , should be known (at least approx.)



$$W = (-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, +\infty)$$

If $T \in W$, we reject H_0

Goal: $P(T \in W; H_0) = \alpha$

Hypothesis testing – an example

two-sample test

- ind. $\left\{ \begin{array}{l} \blacktriangleright X_1, \dots, X_{n_1} \text{ random sample from } \underline{Ber}(\vartheta_X) \\ \blacktriangleright Y_1, \dots, Y_{n_2} \text{ random sample from } \underline{Ber}(\vartheta_Y) \\ \blacktriangleright H_0 : \vartheta_X = \vartheta_Y, \quad H_1 : \vartheta_X \neq \vartheta_Y \end{array} \right.$

$X_i = [i\text{-th male patient recovered}]$

$Y_j = [j\text{-th female recovered}]$

$$\hat{\vartheta}_X = \bar{X}_{n_1} = \frac{X_1 + \dots + X_{n_1}}{n_1} \quad \hat{\vartheta}_Y = \bar{Y}_{n_2} = \frac{Y_1 + \dots + Y_{n_2}}{n_2}$$

approx. normal by CLT

we estimate $\vartheta_X - \vartheta_Y = \hat{\vartheta}_X - \hat{\vartheta}_Y =: Z$ if this is too large, we reject H_0

$EZ = E(\hat{\vartheta}_X) - E(\hat{\vartheta}_Y) = \vartheta_X - \vartheta_Y = 0$ ASSUMING H_0 $\vartheta = \vartheta_X = \vartheta_Y$

$\text{var}(Z) = \text{var}(\hat{\vartheta}_X - \hat{\vartheta}_Y) = \text{var}(\hat{\vartheta}_X) + \text{var}(\hat{\vartheta}_Y) = \frac{\text{var}(X_i) \cdot n_1}{n_1^2} + \frac{\text{var}(Y_j) \cdot n_2}{n_2^2}$

$\overset{\sigma^2}{\text{ind.}}$

we estimate ϑ by $\hat{\vartheta} = \frac{\sum X_i + \sum Y_j}{n_1 + n_2} = \frac{\vartheta_X \cdot (1 - \vartheta_X)}{n_1} + \frac{\vartheta_Y \cdot (1 - \vartheta_Y)}{n_2} = \frac{\vartheta \cdot (1 - \vartheta)}{(\frac{1}{n_1} + \frac{1}{n_2})}$

$\hat{\sigma}^2 = \hat{\vartheta}(1 - \hat{\vartheta}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$ $T = \frac{\hat{\vartheta}_X - \hat{\vartheta}_Y}{\sqrt{\hat{\vartheta}(1 - \hat{\vartheta}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim N(0, 1)$

$\hat{\vartheta}_X, \hat{\vartheta}_Y$ - approx. normal $\rightarrow \hat{\vartheta}_X - \hat{\vartheta}_Y$ approx. normal

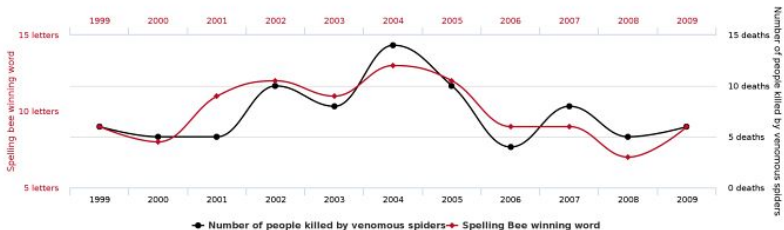
$P(T \leq \alpha) = \alpha$

- ▶ we first gain data, then look for interesting stuff
- ▶ – given enough data, there will be random coincidences
- ▶ even worse, we may test, until we get the desired outcome
- ▶ reproducibility – after exploratory analysis of the data we make an independent measurement and a confirmatory analysis.
- ▶ or we split the data in advance to a part for hypothesis formation and part for verification . . . simple example of cross validation

Letters in winning word of Scripps National Spelling Bee

correlates with

Number of people killed by venomous spiders



Overview

Hypothesis testing

Goodness of fit test

Linear regression

χ_k^2 – chi-square distribution $E(Q) = E(Z_1^2) + \dots + E(Z_k^2)$ (easy)
 $\sim 1 + \dots + 1 = k$

Definition

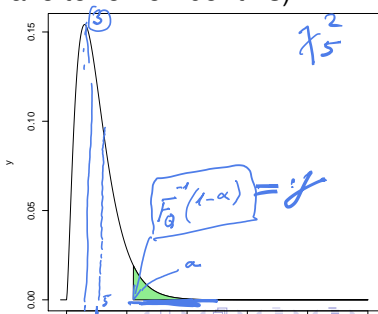
$Z_1, \dots, Z_k \sim N(0, 1)$ i.i.d. The distribution of r.v.

$E(Z_i) = 0$
 $var(Z_i) = 1 = E(Z_i^2) - (E(Z_i))^2$

$Q = Z_1^2 + \dots + Z_k^2$

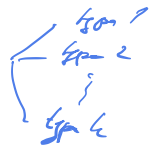
is called chi-square(d) with k degrees of freedom (really $k!$), and denoted χ_k^2 .

- ▶ $E(Q) = k$ (easy)
- ▶ $var(Q) = 2k$ (fyi, you don't have to remember this)
- ▶ density can be written by a reasonable formula
- ▶ $Q \doteq N(k, 2k)$ for large k (CLT)



Multinomial and categorical distribution

expt. with k poss. outcomes



Definition

Given $p_1, \dots, p_k \geq 0$ so, that $p_1 + p_2 + \dots + p_k = 1$.

we repeat n -times an experiment with k possible outcomes, where the i -th has probability p_i

$X_i :=$ how many times we got the i -th outcome (X_1, \dots, X_k) has multinomial distribution with parameters $n, (p_1, \dots, p_k)$.

- ▶ trivial example: $X_i =$ number of die rolls that equaled i
- ▶ important example: $X_i =$ number of occurrences of i -th letter,
- ▶ $P(X_1 = x_1, \dots, X_k = x_k) = \binom{n}{x_1, \dots, x_k} p_1^{x_1} \dots p_k^{x_k}$

gives $x_1 + \dots + x_k = n$

$$X_1 + \dots + X_k = n$$

Pearson χ^2 statistics

- ▶ (X_1, \dots, X_k) – multinomial distribution with parameters $n, (p_1, \dots, p_k)$ as above
- ▶ $E_i := \mathbb{E}(X_i) = np_i$
- ▶ Pearson χ^2 statistics is the function

$$T_n = (Z_n)^2 \quad Z_n \xrightarrow{d} N(0,1)$$

$$Z_n^2 \xrightarrow{d} \chi^2_k$$

$$\chi^2 = T := \sum_{i=1}^k \frac{(X_i - E_i)^2}{E_i}$$

k fixed; as $n \rightarrow \infty$

▶ **Theorem** $T \xrightarrow{d} \chi_{k-1}^2$

sketch of proof for $k=2$

$$X_1 + X_2 = n$$

$$p_1 + p_2 = 1$$

$$E_1 = np_1$$

$$T = \frac{(X_1 - np_1)^2}{np_1} + \frac{(X_2 - np_2)^2}{np_2} = \frac{(X_1 - np_1)^2}{np_1} + \frac{(n - X_1 - n(p_2))^2}{np_2}$$

$$\frac{(X_1 - np_1)^2}{n} \cdot \left(\frac{1}{p_1} + \frac{1}{p_2} \right) = \frac{(X_1 - np_1)^2}{n p_1 p_2} = \left(\frac{X_1 - np_1}{\sqrt{np_1 p_2}} \right)^2$$

Z_i – observed

$$X_1 \sim \text{Bin}(n, p_1)$$

$$\frac{X_1 - np_1}{\sqrt{np_1(1-p_1)}} \approx \text{approx.} \sim N(0,1) \text{ CCT}$$

(the limit does not depend on $p_1 \dots p_k$)

Goodness of fit test

- ▶ (X_1, \dots, X_k) – multinomial distribution with parameters $n, (\vartheta_1, \dots, \vartheta_k)$ as above
- ▶ n known, ϑ unknown
- ▶ Null hypothesis $H_0: \vartheta = \vartheta^*$ for some given ϑ^*
- ▶ $E_i := n\vartheta_i^*$ for all i
- ▶ We use the statistics $\chi^2 = T := \sum_{i=1}^k \frac{(X_i - E_i)^2}{E_i}$ approx. $T \sim \chi_{k-1}^2$
- ▶ We reject H_0 iff $T > \gamma$
- ▶ $\gamma := F_Q^{-1}(1 - \alpha)$, where $Q \sim \chi_{k-1}^2$
- ▶ $P(\text{l type error}) = P(T > \gamma; H_0) \rightarrow P(Q > \gamma) = \alpha$

$1 - F_T(\gamma) \xrightarrow{n \rightarrow \infty} 1 - F_Q(\gamma)$

Goodness of fit test – example

- ▶ We roll a die repeatedly (600 times). The numbers 1 upto 6 came up with frequencies 92, 120, 88, 98, 95, 107.
- ▶ Is the die fair? $x_1 \quad x_2 \quad \dots \quad x_6$

$$n = 600 \rightarrow * \left(\frac{1}{6}, \dots, \frac{1}{6} \right)$$

$$E_i = n \theta_i^* = 100$$

$$T = \sum_{i=1}^6 \frac{(x_i - E_i)^2}{E_i} = \sum \frac{(x_i - 100)^2}{100} = \frac{8^2}{100} + \frac{20^2}{100} + \frac{12^2}{100} + \frac{2^2}{100} + \frac{5^2}{100} + \frac{2^2}{100}$$

$$= \underline{6.86} < 11.1 \Rightarrow \text{do not reject } H_0$$

(we trust die is fair)

$$Q \sim \chi_5^2$$

$$F_{0.95}^{-1}(0.95) = 11.1$$

$$1 - F_Q(6.86)$$

$$p\text{-value} = 0.23$$

Extensions

$$X_i = \mathbb{1}_{\{Y \in B_i\}}$$

- ▶ To study a distribution of an arbitrary r.v. Y we can pick “bins” B_1, \dots, B_k (a partition of \mathbb{R}) and look how often $Y \in B_i$ (this will be measured by r.v. X_i).
- ▶ Similar test for independence of discrete random variables.

Overview

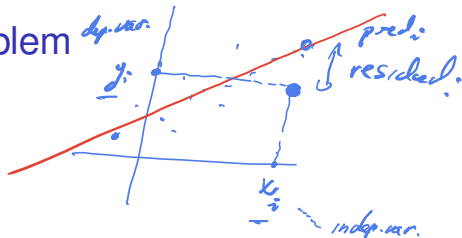
Hypothesis testing

Goodness of fit test

Linear regression

Linear regression – the problem

- ▶ data: (x_i, y_i) for $i = 1, \dots, n$
- ▶ goal: $y = \vartheta_0 + \vartheta_1 x$



- ▶ we measure how good fit we have by the quadratic error:

$$\sum_{i=1}^n \underbrace{(y_i - \underbrace{(\vartheta_0 + \vartheta_1 x_i)}_{\text{pred}_i})}_{\text{residual}_i}^2$$

Linear regression – solution

- ▶ To minimize

$$\sum_{i=1}^n (y_i - (\vartheta_0 + \vartheta_1 x_i))^2$$

we know from LA

- ▶ the optimal parameters are

$$\hat{\vartheta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{\vartheta}_0 = \bar{y} - \vartheta_1 \bar{x},$$

where $\bar{x} := (x_1 + \dots + x_n)/n$, $\bar{y} := (y_1 + \dots + y_n)/n$.

$\frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})$ — sample covariance

$\frac{1}{n-1} \sum (x_i - \bar{x})^2$ — sample variance

Linear regression – why sum of squares?

- ▶ We assume that x_1, \dots, x_n are fixed, y_i is a realization of a r.v.

$$Y_i = \vartheta_0 + \vartheta_1 x_i + W_i$$

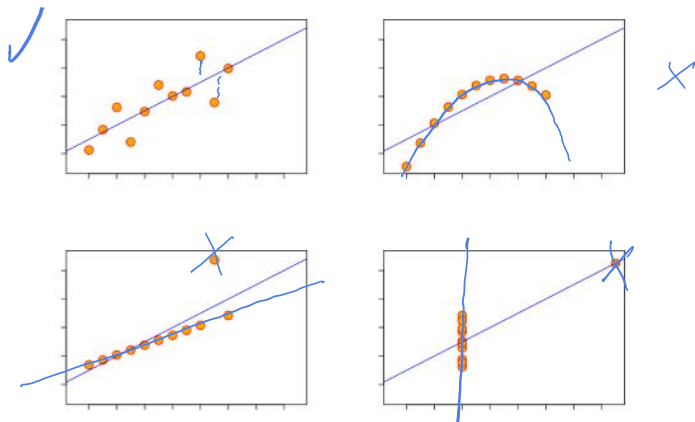
- ▶ $W_i \sim N(0, \sigma^2)$ for all i ; W_1, \dots, W_k iid
- ▶ maximal likelihood:

*to take prod.
by physical
eg'*

$$L(y; \vartheta) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i - \vartheta_0 - \vartheta_1 x_i)^2}{2\sigma^2}}$$

- ▶ $\ell(y; \vartheta) = \log L(y; \vartheta) = a + b \sum_{i=1}^n (y_i - \vartheta_0 - \vartheta_1 x_i)^2$

Limits of regression



(data: Francis Anscombe 1973, image: wikieditor Schutz)

- ▶ nonlinear regression
- ▶ logistic regression

box w. balls 1, ..., N

$N = ?$



→ estimate N

→ pick 5 balls x_1, \dots, x_5 with repeat.

moment method: 1-st moment = $\frac{N+1}{2} = \frac{x_1 + \dots + x_5}{5}$

$$\Rightarrow \hat{N} = 2 \frac{x_1 + \dots + x_5}{5} - 1$$

$$\uparrow \mathbb{E}(\hat{N}) = N$$

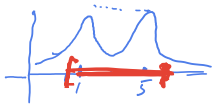
MSE larger than other methods

1, 1, 1, 1, 11 → 3, but $N \geq 11$
↳ $\sum_{i=1}^n x_i = 5$

max. likelihood

$$\hat{N} = \max \{x_1, \dots, x_5\}$$

$$\hat{N} = \max(x_1, \dots, x_5), \text{ const}$$



Simpson's paradox

Treatment Stone size	Treatment A	Treatment B
Small stones	Group 1 93% (81/87)	Group 2 87% (234/270)
Large stones	Group 3 73% (192/263)	Group 4 69% (55/80)
Both	78% (273/350)	83% (289/350)

