

UPPSALA UNIVERSITET

FÖRELÄSNINGSANTECKNINGAR

Inferensteori

Rami Abou Zahra

Inlämningsdatum
November 2, 2022

CONTENTS

1. TODO	2
2. Data Analysis (K6)	3
2.1. Location Measures	3
2.2. Dispersion measures	3
2.3. Graphical illustration	4
2.4. Data materials in several dimensions	5
3. Statistical Inference	7
4. Estimation	7
4.1. Properties of estimates	8
4.2. Asymptotic properties	9

1. TODO

- Experiment in r (QQ-plot of exp vs $n(0,1)$ data)
- Understand .dat files
- Add from slides
- Add proof from book of theorem 4.9
- Problems 7.2.2 in the book

2. DATA ANALYSIS (K6)

Vi kommer undersöka statistisk säkerställd skillnad (Opinion polls example), hypotestestning (räknar sannolikheten att hypotesen är sann).

Anmärkning:

Vanligtvis antar vi att datan är normalfördelad, men inte i alla fall (såsom stickprov av lön)

2.1. Location Measures.

A data set is given by x_1, \dots, x_n

Definition/Sats 2.1: Sample mean

Sample mean is given by:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Definition/Sats 2.2: Median

The "middle value" of the sorted data. Different from the mean.

If n is even, the median is defined as the mean of the two middle values

Definition/Sats 2.3: Mode

This doesn't work if it's continuous data but it can be made discrete (such as age/time)

Mode is the most common data value

Example:

Let our data points be:

32 34 41 44 45 50 50 54 55 57 58 60 63

Find mean, median mode:

Mean: 13 data sets $\Rightarrow n = 13$:

$$\frac{1}{13}(32 + 34 + 41 + 44 + 45 + 50 + 50 + 54 + 55 + 57 + 58 + 60 + 63) \approx 49.46$$

Median: The middle value is 50

Mode: 50 is the only data value appearing more than once.

Anmärkning:

In this example, the median = mode. This is not always the case!

2.2. Dispersion measures.

Describes the "spread" of the data, such as the variance. We have the following:

Definition/Sats 2.4: Sample variance

The sample variance is given by:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Definition/Sats 2.5: Sample standard variance

Is given by:

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Definition/Sats 2.6: Range

Variationsbredden. The difference between the largest and the smallest values of the data

Definition/Sats 2.7: Inter quartile range

Kvartilavståndet is the difference between the upper and lower quartiles.
If we have an odd amount of data it is including the median!

Definition/Sats 2.8: Mid-range

The mean between the biggest and smallest value in the sample

Definition/Sats 2.9: Lower/Upper quartile

The *lower quartile* is the median of the lower half of the data material including the median if n is odd

The *upper quartile* is the median of the upper half of the data material including the median if n is odd

Example:

0 0 1 1 2 2

Here, the mean is given by $\frac{(1+1+2+2)}{6} = 1$.

Therefore, the sample variance is given by $\frac{4}{5}$ and the sample standard deviation $\sqrt{\frac{4}{5}}$

We can find the inter quartile range by looking at the half, like this:

$[0 \underbrace{0}_{\Delta} 1] [1 \underbrace{2}_{\Delta} 2]$

Therefore, the inter quartile range here is $2 - 0 = 2$

2.3. Graphical illustration.**Stem av leafplots:**

```
u = c(32,34,...)
stem(u)
```

Boxplots:

Uses quartiles, max min, and median. Useful if you want a quick look at the dispersion of data.

Bar chart:

Good for illustrating the frequency of each data point, but for large data points the data is hard to read

Histogram:

Attempts to fix the readability issues with the bar chart and is easier to compare with probability density functions.

Easier to manipulate data for readability (use bigger/smaller intervals) (one can ask what the optimal width for a histogram would be)

Very often you can ask if the data follow a normal distribution, which can be hard by just looking at the histogram (because the width varies)

Thoughts:

Dynamically widths on histograms, the more sparse data the greater the width and the more dense, the less the width

QQ-plot:

Is the data normally distributed? You order your data and construct a table with your data and compare it with if it was normally distributed:

$$\Phi(z) = \frac{i - 0.5}{n}$$

If data was perfectly normal on both axis, x_i would be a linear function of z , ie. normally distributed $N(0, 1)$

We plot z on the x -axis and x_i on the y -axis

The name comes from quantile-quantile-plot (QQ-plot). It is a graphical way of comparing two probability distributions (sannolikhetsfördelning)

2.4. Data materials in several dimensions.

We can calculate correlation through sample covariance:

Definition/Sats 2.10: Sample covariance

$$c_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Not scale invariant (if you measure x in meters and go to cm then it is not the same). Therefore we need to norm it with something, which is where the correlation comes in:

Definition/Sats 2.11: Sample correlation coefficient

$$r_{xy} = \frac{c_{xy}}{s_x s_y}$$

Where s_x and s_y are the sample standard deviations for x and y

Definition/Sats 2.12: Sample correlation satisfies

The sample correlation coefficient satisfies

$$-1 \leq r_{xy} \leq 1$$

If it is 1, then there is a strong positive correlation (the linear regression has a line with positive derivative), similarly for negative.

When it is 0 there is no *linear* relation. There might be other, for example quadratic relation.

Bevis 2.1: Sample correlation satisfaction

$$\begin{aligned}
0 &\leq \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} - \frac{y_i - \bar{y}}{s_y} \right)^2 \\
&= \frac{1}{s_x^2} \frac{1}{n-1} \underbrace{\sum_i (x_i - \bar{x})^2}_{s_x^2} + \frac{1}{s_y^2} \frac{1}{n-1} \underbrace{\sum_i (y_i - \bar{y})^2}_{s_y^2} - 2 \frac{1}{s_x s_y} \frac{1}{n-1} \underbrace{\sum_i (x_i - \bar{x})(y_i - \bar{y})}_{c_{xy}} \\
&= 2 - 2r_{xy} \Rightarrow r_{xy} \leq 1 \\
0 &\leq \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} + \frac{y_i - \bar{y}}{s_y} \right)^2 = 2 + 2r_{xy} \\
&\Rightarrow -1 \leq r_{xy}
\end{aligned}$$

□

3. STATISTICAL INFERENCE

Anmärkning:

If $X \sim Hyp$, then for a large population $X \sim Bin$ (because the chance of picking the same one in a large population is so small)

Example:

See slide 1

The biggest difference between

Definition/Sats 3.1: Sample (stickprov)

x_1, x_2, \dots, x_n is a *sample* from the random variable X with distribution F_X

If $X = (X_1, \dots, X_n)$ are independent, we have a *random sample* from X

We can purposely choose our random variable in such way that makes it easier for us to analyse. It also allows us to compare these observations.

Example:

See slide 2

4. ESTIMATION

Suppose we have one unknown parameter θ

We can write the following for our data:

$$x = (x_1, x_2, \dots, x_n)$$

$$X = (X_1, X_2, \dots, X_n)$$

Definition/Sats 4.1: Estimate (skattning)

An *estimate* $\theta^* = \theta^*(x)$ is a function of the sample x

The estimate is an observation of the estimator $\theta^*(X)$

Example:

See slide 4

The greater the stickprov the better the estimate (because less and less variance)

Anmärkning:

The estimator is not distributed with the same distribution, since the estimator is not always an integer.

If we have different estimates, we need to make a reasonable choice such that our error is as little as possible (this is why we introduce estimators)

4.1. Properties of estimates.

We can take $\theta^* - \theta = E(\theta^*(X)) - \theta + (\theta^* - E(\theta^*(X)))$

It turns out, this is equal to the systematic error + random error

Definition/Sats 4.2: Unbiased (väntevärdesriktigt)

An estimate θ^* is said to be *unbiased* if it satisfies $E(\theta^*(X)) = \theta$

This is the same as saying it has no systematic error (therefore, we only have the random error left)

Example:

We show this by:

$$E(\mu^*(X)) = E(\bar{X}) = \mu$$

Let

$$p^*(X) = \frac{X}{1000} \quad X \sim \text{Bin}(1000, p)$$

Is $p^*(x)$ an unbiased estimate of p ? (slide 5)

If we have more than one unbiased estimate, which is the best one? Well, in that case we need to start looking at the random error. We can study this by looking at the variance

Definition/Sats 4.3: Efficiency comparison of estimates

If θ_1^* and θ_2^* are unbiased estimates of θ and

$$V(\theta_1^*(X)) \leq V(\theta_2^*(X))$$

For all θ with strict inequality for some. We say that θ_1^* is more *efficient* than θ_2^* (less random error)

Example:

See slide 6 & 7

Example: (stratification)

See slide 8

Here we assume that the number of men that take the plane is Binomially distributed $\text{Bin}(500, p + a)$ while for women $\text{Bin}(500, p - a)$

$$\begin{aligned} p_1^*(X) &= \frac{1}{1000}X \Rightarrow E(p_1^*) = p \\ p_2^*(y, z) &= \frac{1}{1000}y + \frac{1}{1000}z \\ &\Rightarrow \frac{1}{1000}E(y) + \frac{1}{1000}E(z) \\ &= \frac{1}{1000}500(p + a) + \frac{1}{1000}500(p - a) = p \\ V(p_1^*) &= \frac{p(1-p)}{1000} \geq V(p_2^*) = \frac{p(1-p)}{1000} - \frac{a^2}{1000} \end{aligned}$$

Definition/Sats 4.4: Standard error (medelfelet)

We want to assign a numerical value to the dispersion of an estimate.

Therefore, we define the *standard error* of the estimate θ^* is an estimate of the standard deviation $D(\theta^*)$

Denoted by $d(\theta^*(x)) = d(\theta^*)$

Recall that the standard deviation is given by \sqrt{Var}

Example:

See slide 9

4.2. Asymptotic properties.

The accuracy of an estimate should improve as the sample size increases, seems reasonable

Definition/Sats 4.5: Bias

The *bias* (väntevärdesfelet) for the estimate θ^* is defined as

$$B(\theta^*) = E(\theta^*) - \theta$$

Anmärkning:

An unbiased estimate has bias 0

Definition/Sats 4.6: Asymptotically unbiased

If the bias $B(\theta_n^*)$ tends to zero as $n \rightarrow \infty$ for all θ , the estimate θ_n^* is said to be *asymptotically unbiased*

Example:

Let x_1, \dots, x_n be a random sample from $N(\mu, \sigma^2)$ where μ is unknown. We want to estimate σ^2

The estimate $\sigma_n^{2*} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ is biased, but it is asymptotically unbiased

The estimate s_n^2 is unbiased for σ^2 , why?

$$\begin{aligned} s^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) \\ E(s^2) &= \frac{1}{n-1} \left(\sum_{i=1}^n E(x_i^2) - nE(\bar{x}^2) \right) \\ E(x_i^2) &= V(x_i) + (E(x_i))^2 = \sigma^2 + \mu^2 \\ E(\bar{x}^2) &= V(\bar{x}) + (E(\bar{x}))^2 = \frac{\sigma^2}{n} + \mu^2 \\ E(s^2) &= \frac{1}{n-1} \left(\sum_{i=1}^n (\sigma^2 + \mu^2) - n \left(\frac{\sigma^2}{n} + \mu^2 \right) \right) \\ &= \frac{1}{n-1} (n(\sigma^2 + \mu^2) - \sigma^2 - n\mu^2) = \frac{1}{n-1} (n-1)\sigma^2 = \sigma^2 \end{aligned}$$

Definition/Sats 4.7: Consistent estimate

The estimate θ_n^* is said to be *consistent* for θ if the corresponding estimator converges to θ in probability for all θ

Definition/Sats 4.8: Convergence in probability

The estimator θ_n^* converges to θ in probability if $\forall \varepsilon > 0$:

$$\lim_{n \rightarrow \infty} P(|\theta_n^* - \theta| > \varepsilon) = 0$$

Definition/Sats 4.9

If the estimate θ_n^* is asymptotically unbiased and

$$\lim_{n \rightarrow \infty} V(\theta_n^*) = 0$$

Then it is consistent

Example:

See slide 11

If an estimate is not consistent, then it does not matter if the sample size is increased, it wont yield better results.