

Mathematical Statistics

Anna Janicka

Lecture III, 05.03.2018

INTRODUCTION TO STATISTICS

Plan for today

1. Introduction to Mathematical Statistics
 - the statistical model
2. Statistics and their distributions
 - the normal model
3. Estimation – introduction



MATHEMATICAL STATISTICS



WARSAW UNIVERSITY
Faculty of Economic Sciences

Assumptions

Empirical data reflect the functioning of a random mechanism

Therefore: we are dealing with random variables defined over some probabilistic space; the realizations of these random variables are the collected data.

Problem: we do not know the distribution of these random variables...



Difference between Probability Calculus and Mathematical Statistics

1. PC, example:

- Phrasing: in a production process each produced unit may be defective. This happens with probability 10%. The defects of different units are independent.
- Problems: What is the chance that in a batch of 50 items, exactly 6 will be defective? What is the average number of defective elements? What is the most probable number of defective elements?
- Solution: we build a probabilistic model. Here: a Bernoulli Scheme with $n=50$, $p=0,1$.

Alternatively, if we are interested in questions dealing with order (e.g. what is the chance that the first 5 items are defective?): a different model



Difference between Probability Calculus and Mathematical Statistics – cont.

2. MS, example:

- Formulation: An inspector verified a batch of 50 items, with the following results (1– item defective, 0 – OK):

0 1 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1

- Problems: what is the probability that an item is defective (assessment)? Is the producer's declaration that defectiveness is equal to 10% credible?
- Solution: we build a statistical model, i.e. a probabilistic model with unknown distribution parameter(s).



Statistical Model

in PC:

Statistical Model: (X, F_X, P) (Ω, F, P)

where:

X – the space of values for the observed random variable X (often n -dimensional, if we have an n -dimensional sample X_1, \dots, X_n)

F_X – σ -algebra on X

P – a family of probability distributions P_θ , indexed by a parameter $\theta \in \Theta$

~~In a less formal setting we usually provide: X, P, Θ~~



Statistical model – example

$X = \{0, 1\}^n$ – sample space

Joint probability distribution:

$$\begin{aligned} P_{\theta}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) &= \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} \\ &= \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i} \end{aligned}$$

for $\theta \in [0, 1]$

(we have $n=50$, $X_2 = X_{10} = X_{15} = X_{32} = X_{42} =$
 $X_{50} = 1$, other $X_i = 0$)



Statistical model – example cont.

Alternative formulation (if we only record the *number* of defective items in a sample):

$X = \{0, 1, 2, \dots, n\}$ – sample space

Joint probability distribution:

$$P_{\theta}(X = x) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$$

for $\theta \in [0, 1]$

(we have $n=50$ and $X=6$)



Statistical model – example cont. (2)

Possible questions

Based on the sample:

□ What is the value of θ ?

■ we are interested in a precise value

■ we are interested in an interval (*confidence*)

→ *estimation*

□ Verification of the hypothesis that $\theta = 0.1$

→ *hypothesis testing*

Statistical Model: example 2

Growths on the market

An analyst studies the length of periods of growth on the stock market. He is interested in times of growth (until the first fall), in days. Assume the times of growth, X_1, X_2, \dots, X_n are a sample from an exponential distribution $\text{Exp}(\lambda)$, where:

λ – unknown parameter

$X = (0, \infty)^n$ – sample space

Joint probability distribution:

$$P_\lambda(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n) = \prod_{i=1}^n (1 - e^{-\lambda x_i})$$

$$f_\lambda(x_1, x_2, \dots, x_n) = \lambda^n e^{-\lambda \sum x_i}$$

for $\lambda > 0$



Statistical Model: example 3

Measurements with error

We repeat measuring μ , the results of measurements are independent random variables X_1, X_2, \dots, X_n , (our machine is not perfect). Each measurement is normally distributed $N(\mu, \sigma^2)$.

μ, σ^2 – unknown parameters (so $\theta = (\mu, \sigma)$)

$X = \mathbb{R}^n$ – sample space

Joint probability distribution:

$$P_{\mu, \sigma}(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n) = \prod_{i=1}^n \Phi\left(\frac{x_i - \mu}{\sigma}\right) \text{ or}$$

$$f_{\mu, \sigma}(x_1, x_2, \dots, x_n) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

for $\mu \in \mathbb{R}, \sigma > 0$



STATISTICS

(objects)



Statistics

Parameter estimation (both point and interval) as well as hypothesis testing are conducted based on *statistics*

Statistic = a function of observations, i.e. any random variable

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

The distribution of a statistic T depends on the distribution of X , but the statistic as such cannot depend on parameter θ , e.g.

~~$X_1 + X_2 = \theta$~~

Statistics – examples

$$T_1 = \sum_{i=1}^n X_i, \quad T_2 = \frac{1}{n} \sum_{i=1}^n X_i, \quad T_3 = \frac{1}{n} \sum_{i=1}^n X_i - 0.1$$

are statistics for a sample size of n ;

$$T_1 = X, \quad T_2 = \frac{X}{n}, \quad T_3 = \frac{X}{n} - 0.1$$

are statistics for a single observation

The choice of a statistic depends on the question we want to answer.



Distribution of statistics

In many cases statistical models refer to a common set of assumptions → similar models are applied.

Similar questions are posed → similar statistics are calculated.

The most commonly used is the normal model



The normal model

X_1, X_2, \dots, X_n are a sample from $N(\mu, \sigma^2)$.

The most important statistics (*in general, not only for this model*):

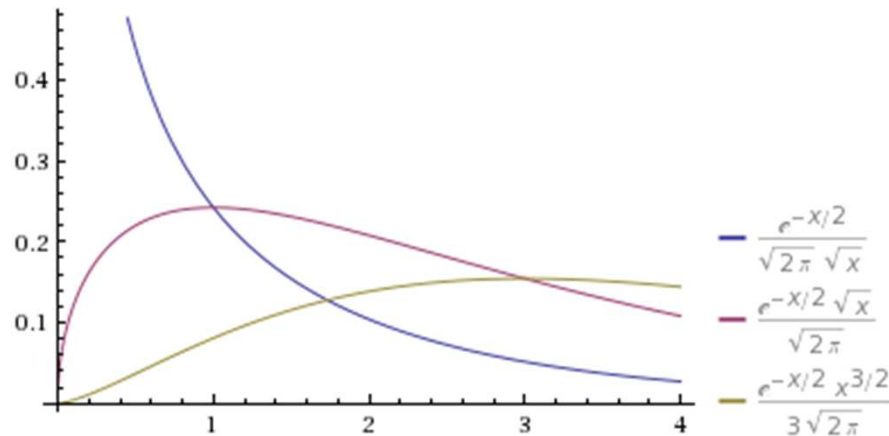
Mean: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

sample variance: $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2,$

standard deviation: $S = \sqrt{S^2}$

what are their
distributions?

Chi-squared Distribution $\chi^2(n)$



A special case of the gamma distribution.

The sum of squares of n IIN random variables (independent identically $N(0,1)$ distributed) has a $\chi^2(n)$ distribution

$$\mathbb{E}X = n, \quad \text{Var}X = 2n$$

The normal model – cont. (1)

Theorem: In the normal model, the \bar{X} and S^2 statistics are independent random variables such that

$$\begin{aligned}\bar{X} &\sim N(\mu, \sigma^2/n) & \frac{(\bar{X} - \mu)}{\sigma} \sqrt{n} &\sim N(0,1) \\ \frac{n-1}{\sigma^2} S^2 &\sim \chi^2(n-1)\end{aligned}$$

in particular:

$$E_{\mu, \sigma} S^2 = \sigma^2, \text{ and } \text{Var} S^2 = \frac{2\sigma^4}{(n-1)}$$



The normal model – cont. (2)

In the normal model, *the variable*

$$T = \frac{\sqrt{n}(\bar{X} - \mu)}{S}$$

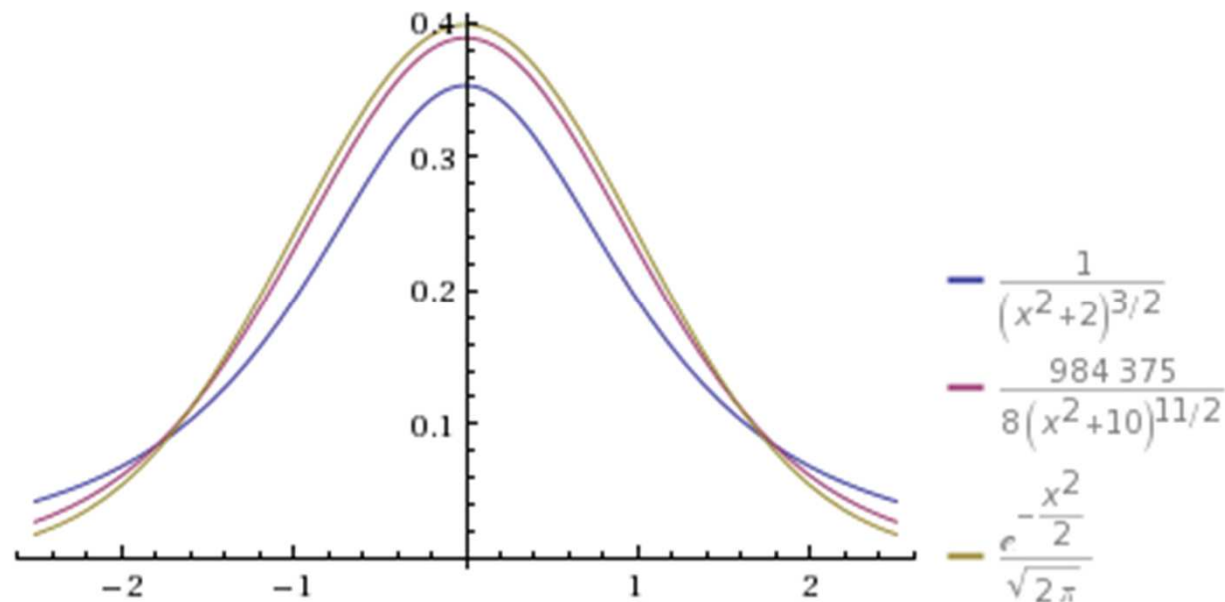
has a t-Student distribution with $n - 1$ degrees of freedom, $T \sim t(n - 1)$



***t*-Student Distribution $t(n)$, $n=1,2,\dots$**

defined as the distribution of the random variable

$\frac{\sqrt{n}X}{\sqrt{Y}}$ for independent X and Y , $X \sim N(0,1)$, $Y \sim \chi^2(n)$



$$\mathbb{E}X = 0 \quad n > 1 \quad \text{Var}X = \frac{n}{n-2} \quad n > 2$$

ESTIMATION



Point Estimation

- The choice, on the base of the data, of *the best* parameter θ , from the set of parameters which may describe P_θ
- An **Estimator** of parameter θ is any statistic $T = T(X_1, X_2, \dots, X_n)$ with values in Θ (we interpret it as an approximation of θ). Usually denoted by $\hat{\theta}$
- Sometimes we estimate $g(\theta)$ rather than θ .



Estimation: an example

Empirical frequency

Quality control example:

0 1 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 1

□ Model: $X = \{0, 1, 2, \dots, n\}$ (here $n=50$),

$$P_{\theta}(X = x) = \binom{n}{x} \theta^x (1 - \theta)^{n-x} \text{ for } \theta \in [0, 1]$$

□ parameter θ : probability of faulty element

□ an obvious estimator: $\hat{\theta} = X/n = 6/50$

n – sample size

X – number of faulty elements in sample



Problems with (frequency) estimators...

Example: three genotypes in a population, with frequencies $\theta^2 : 2\theta(1-\theta) : (1-\theta)^2$

In a population of size n , N_1 and N_2 and N_3 individuals of particular genotypes were observed.

Should we take $\hat{\theta} = \sqrt{N_1/n}$? or rather $\hat{\theta} = 1 - \sqrt{N_3/n}$? How about $\hat{\theta} = N_1/n + \frac{1}{2} N_2/n$?

Maybe something else?

→ How do we choose the best one?



Estimation – sample characteristics

Sample characteristics:

estimators based on the empirical distribution (empirical CDF)



Empirical CDF

- Let X_1, X_2, \dots, X_n be a sample from a distribution given by F (modeled by $\{P_F\}$)

(n -th) empirical CDF

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(-\infty, t]}(X_i) = \frac{\text{number of observations } X_i : X_i \leq t}{n}$$

- For a given realization $\{X_i\}$ it is a function of t , the CDF of the empirical distribution (uniform over x_1, x_2, \dots, x_n). For a given t it is a statistic with a distribution

$$P(\hat{F}(t) = \frac{k}{n}) = \binom{n}{k} F(t)^k (1 - F(t))^{n-k}, \quad k = 0, 1, \dots, n$$



Empirical CDF: properties

1. $E_F \hat{F}_n(t) = F(t)$

2. $\text{Var} \hat{F}_n(t) = \frac{1}{n} F(t)(1 - F(t))$

3. from CLT: $\frac{\hat{F}_n(t) - F(t)}{\sqrt{F(t)(1 - F(t))}} \sqrt{n} \xrightarrow{n \rightarrow \infty} N(0,1)$

i.e., for any z : $P\left(\frac{\hat{F}_n(t) - F(t)}{\sqrt{F(t)(1 - F(t))}} \sqrt{n} \leq z\right) \rightarrow \Phi(z)$

4. Glivenko-Cantelli Theorem

$$\sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)| \xrightarrow{a.s.} 0$$

for $n \rightarrow \infty$

if sample size increases, we will approximate the unknown distribution with any given level of precision



Order statistics

- Let X_1, X_2, \dots, X_n be a sample from a distribution with CDF F . If we organize the observations in ascending order:

$X_{1:n}, X_{2:n}, \dots, X_{n:n} \leftarrow$ **order statistics**

($X_{1:n} = \min, X_{n:n} = \max$)

- An empirical CDF is a stair-like function, constant over intervals $[X_{i:n}, X_{i+1:n})$



Distribution of order statistics

- Let X_1, X_2, \dots, X_n be independent random variables from a distribution with CDF F . Then $X_{k:n}$ has a CDF equal to

$$F_{k:n}(x) = P(X_{k:n} \leq x) = \sum_{i=k}^n \binom{n}{i} (F(x))^i (1 - F(x))^{n-i}$$

- If additionally the distribution is continuous with density f , then $X_{k:n}$ has density

$$f_{k:n}(x) = n \binom{n-1}{k-1} f(x) (F(x))^{k-1} (1 - F(x))^{n-k}$$





WARSAW UNIVERSITY
Faculty of Economic Sciences