# Principles of Biostatistics and Informatics

4th Lecture: Elements of Probability Calculus

2nd October 2014

Dániel VERES

In this lecture we are discussing on probability.

First I will show you two experiments to feel the strange word of probabilities.

Then we define the probability as a quantity based on the law of large numbers.

Thereafter we discuss on the probability of events (notation, and/or relation between events, mutually exclusive and independent events,Kolmogorov axioms and conditional probability).

After that we discuss on how to estimate/calculate probabilities using probability calculus and special theoretical distributions. We define the expected value and the theoretical variance. We describe the uniform, binomial, poisson, Gaussian, lognormal, exponential distributions with examples.

At the end I give you two example for how our mind works and how it have to be....

# The Beginnings... Let's Play a Game

Coin tossing game:
- The pot starts at 2 dollars and it is *doubled* every time a head appears.
- The first time a tail appears, the game ends and Peter wins whatever is in the pot:
  - Peter wins $2 if a *tail appears on the first toss*
  - Peter wins $4 if a head appears on the first toss and a *tail on the second*
  - 8$ if a head appears on the first two tosses and a *tail on the third*

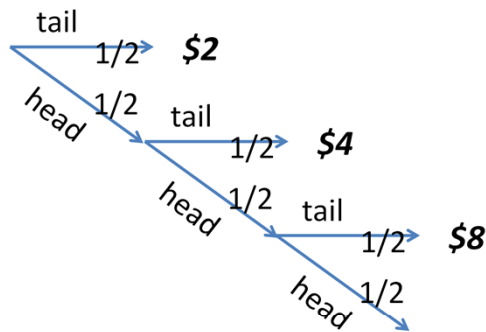Q: What would be a „fair" price for taking part in this game?

A: It would be the expected prize for one game.

Based on some opinion the beginnings of probability calculus was the Saint Petersburg paradox. It was a theoretical game publicated in 1715. A theoretical person – *Peter* – plays the game. The rules are the next.

1. The pot starts at 2 dollars and it is *doubled* every time a head appears.

2. The first time a tail appears, the game ends and Peter wins whatever is in the pot:
- Peter wins $2 if a *tail appears on the first toss*
- Peter wins $4 if a head appears on the first toss and a *tail on the second*
- 8$ if a head appears on the first two tosses and a *tail on the third*
- and so on

The question appers what would be a fair price for taking part in this game? It would be the expected prize (that Peter win) for one game.

## The Beginnings...

```
tail
    1/2 → $2
    1/2   tail
head      1/2 → $4
          1/2   tail
    head      1/2 → $8
              1/2
        head
```

The theoretical „fair" price:     **expected („mean") infinite $ in one game!**

$$\frac{1}{2} \cdot 2 + \frac{1}{2^2} \cdot 4 + ... + \frac{1}{2^n} \cdot 2^n$$

**In practice:** we never win an infinite sum...

   *Buffon*: Made 2048 tosses and won $9.82 on average (mean of the prizes). One million tossing: $10.94

What will be the prize Peter win?

In the half of the cases the first result is tail – so Peter win ½*2 dollars in average.

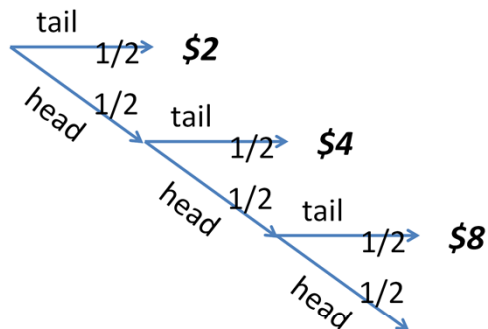In the half of the cases when the first toss is head (that is in the half of the cases) the second is tail – so we get ½*1/2*4 dollars in average and so on.

We could calculate the expected („average", or „mean") prize in one single game based on the equation in the slide. Based on that the expected price is infinite in one game.

But in practice we  realize the sum in one single game is never infinite! For example: *Buffon* (a famous mathematician) made 2048 tosses and won $9.82 on average (the mean of the prizes). With 1 million tossing we get $10.94 on average.

# The Beginnings...



The theoretical fair prize: **expected („mean") infinite \$ in one game!**

$$\frac{1}{2} \cdot 2 + \frac{1}{2^2} \cdot 4 + ... + \frac{1}{2^n} \cdot 2^n$$

**In practice:** we never win an infinite sum...

But the more we play the more the mean prize for one game increases and tends to the theoretical infinite prize.

**How many times do we play?**

We could experience the mean of *the prize won in one single game is increasing and tends to the theoretical infinite if we play more and more games.*

# An Other Experiment...

We have a quick test for a disease:
   white: healthy
   green: ill
We want to figure out whether there is an epidemic in a certain area based on the proportion of ill people. What we know is:
- In non-affected („healthy") areas:
   1 is green out of 10 people
- In affected areas:
   9 are green out of 10 people

Is there an epidemic in the unkown area in question?

In this other experiment we are modelling an epidemic investigation. We have a quick test for a disease. White coloured result indicate healthy, green one indicate ill people. We want to figure out whether there is an epidemic in a certain area based on the proportion of ill people.
What we know is:
        in a non-affected (healthy) area there is 1 ill out of 10 and
        in the affected area there is 9 ill out of 10 people.
Now we are begining to test the people in an unknown area where we would like to know wheather there is a disease – and may be an epidemic.

# An Other Experiment...

We have a quick test for a disease:
   white: healthy
   green: ill
We want to figure out whether there is an epidemic in a certain area based on the proportion of ill people. What we know is:
- In non-affected („healthy") areas:
   1 is green out of 10 people
- In affected areas:
   9 are green out of 10 people

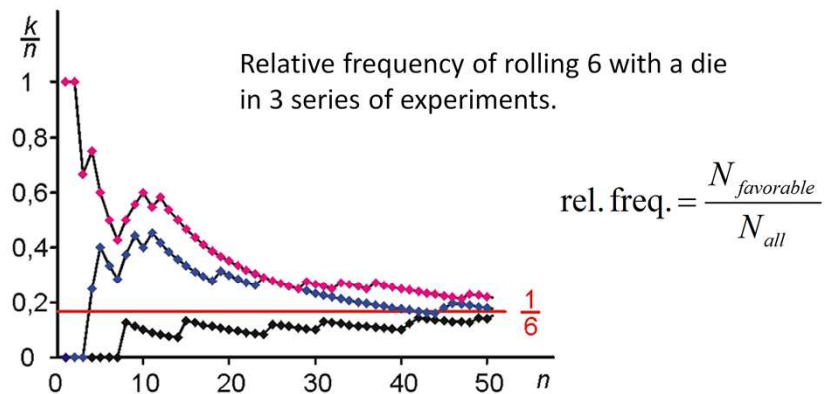Is there an epidemic in the unkown area in question?

***Increasing the number of measurements increase the „certanity".***

We have limited time and resources so we couldn't test every people, but we have to make a decision.
We found *that increasing the number of measurements increase the „certanity".*
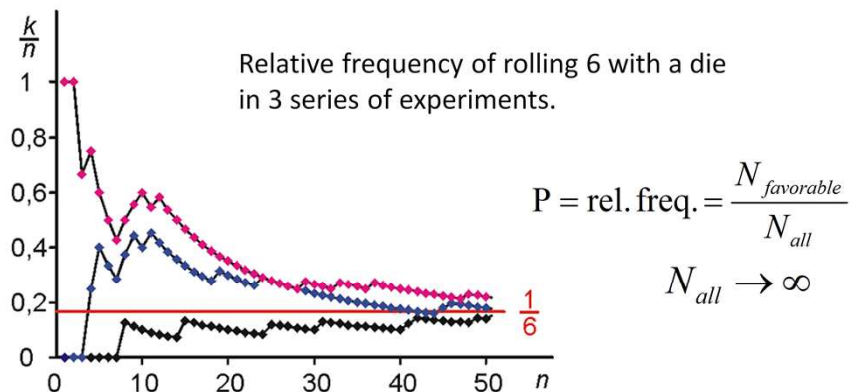
In an other experiment we roll a die 50 times and we count the relative frequency of rolling 6. We repeat this experiment 3 times. We expirience that the *relative frequencies* (frequency of favorable/all measurement) *tend to a certain value independetly from the actual series of experiments if we increase the number of the rolls.*

# Probability as a Quantity

Relative frequency of rolling 6 with a die in 3 series of experiments.

$$P = \text{rel. freq.} = \frac{N_{favorable}}{N_{all}}$$

$$N_{all} \rightarrow \infty$$

$\frac{k}{n}$ axis values: 1, 0,8, 0,6, 0,4, 0,2, 0

$n$ axis values: 0, 10, 20, 30, 40, 50

$\frac{1}{6}$

**Law of large numbers** (on relative frequencies): the relative frequency in an infinite sequence tends to a certain value.
We assign that **certain value** to an **event**: *1/6* to *rolling 6* with a die.
This value is called the ***probability of an event***.
This is an *empirical law – cannot be proven* by logical sequence.

---

*Law of large numbers (on relative frequencies): the relative frequency in an infinite sequence tends to a certain value.*
We assign that certain value to an event: 1/6 to rolling 6 with a die.
This value is called *probability of an event*.
The relative frequency is equal to the probability if the sequence is infinite.
This is an empirical law – can not be proven by logical sequence.

# Probability of Events I.

**Notation:**

Event: **A**
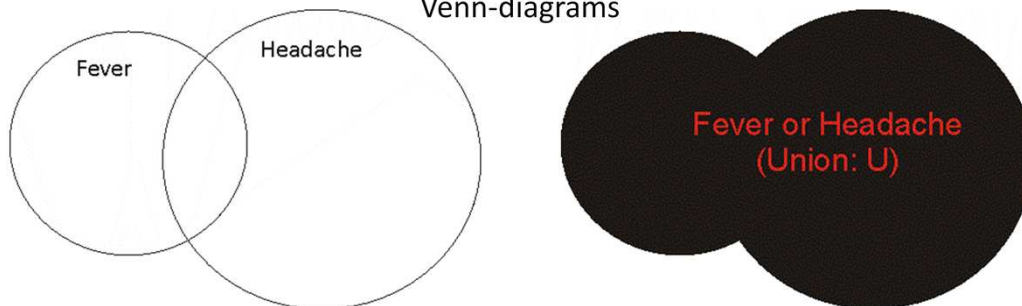
*(the patient has fever)*

Probability that event A occurs**: P(A)**

*(the probability that the patient has fever)*

Probability that event A **or** event B occur: **P(AorB)**, **P(A+B)**, **P(A∪B)**

*(the probability that the patient has fever or headache)*

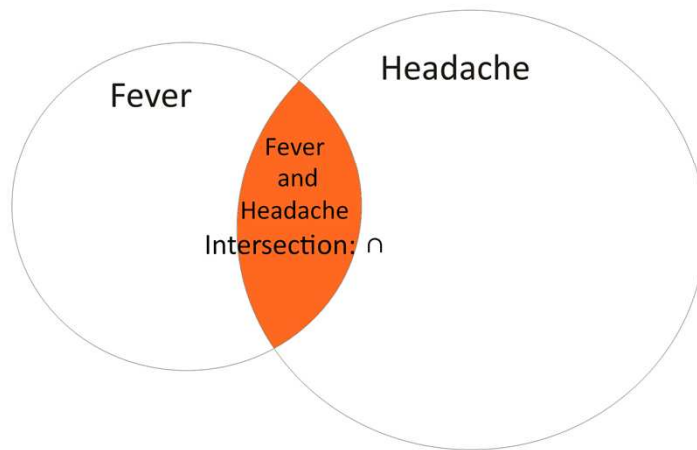Venn-diagrams

Headache

Fever

Fever or Headache
(Union: U)

Now we describe some properties of events' propability. First about the notation.
(Eaxamples are given in parentheses with italian format)
An event is notated with a capital letter – e.g. A (*the patient has fever*). Its probability
symbolized as P(A). Example P(A) = the probability that the patient has fever.
Probability that event A *or* event B occur (*the probability that the patient has a fever or a
headache*) could be notated in three way: P(AorB), P(A+B), P(AUB).  The last one refer to
a set definition: the *Union* – abbreviated U – is all value that belongs to at least one of
the sets. On Venn-diagram we can plot it as it shown in the slide.

# Probability of Events II.

Prob. that both events A **and** B occur: **P(AandB)**, **P(A\*B)**, **P(AB)**, **P(A∩B)**
*(the probability that the patient has both fever and headache)*

Fever

Headache

Fever
and
Headache
Intersection: ∩
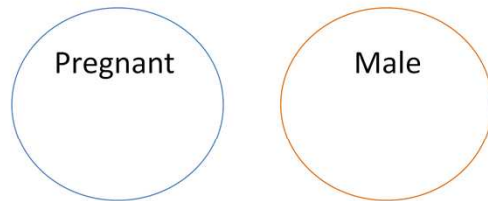
The probability that both events A and B occur could be notated as P(AandB), P(A*B), P(AB), P(A∩B). (*The robability that the patient has fever and headache.*) The ∩ is the symbol in set theory for the intersection. The Venn-diagram for intersection is shown in the slide.

# Probability of Events III.

**Mutually exclusive events:** A and B cannot occur at the same time.
*(the patient is both pregnant and male)*        *(A∩B)=0*



**Independent events**: occurrence of A does not affect the occurrence
  of B
*(our first patient is male and the second one is female)*

---

We have to highlight two event relations.
First the mutually exclusive events that mean two events (A,B) cannot occur at the same time. (*The patient is both pregnant and male.*) It means the intesection of this two event is an empty set. On Venn-diagram we see that the two set has no intersection.
In the case of indepent events the occurrence of an event doesn't affect the occurence of the second one. (*Our first patient is male and the second one is female.*)

# Probability of Events IV.

**Conditional** probability
Probability of A **given that** B has occurred: **P(A|B).**
*(the probability that a patient suffering from a viral infection has actually flu − and not some other type of viral infection)*

Before we go on we have to define the conditional probability. Conditional probability is the probability of an event given an other event has occured. (*The probability that a patient suffering from a viral infection has actually flu – and not some other type of viral infection.*) The notation for the conditional probability is P(A|B).

# Probability of Events V.

**Axioms on probability of events (Kolmogorov):**

*1*. **0 ≤ P(A) ≤ 1**

2. **P(*sure*) = 1** *(The patient will die sooner or later)*
   **P(*impossible*) = 0** *(I'm 310 cm tall)*

3. *Mutually exclusive* events (i.e. P(AandB)=0)
   **P(AorB)=P(A)+P(B)**
   *(probability of being pregnant or male)*

And a theorem:
+4. *Independent* events: **P(AandB)=P(A)*P(B)**
   *(probability that our first patient is male and the second one is female)*

To describe the probability of events we have axioms. Now we show the Kolmogorov axioms. (In a simplified way.)

1. The probability of an event is between 0 and 1.
2. The probability of a sure event (*the patient will die sooner or later* – we know that life is a sexually transmitted lethal disease☺) is 1. The probability of an impossible event is 0 (*I'm 310 cm tall*).
3. The probability of A or B events occur if A and B are mutually exclusive events is the sum of the probabilty of A and the probability of B events. (*The probability that being pregnant or male is the probability that being pregnant + the probability that being male.*)

A theorem based on the axioms:

+4 The probability of A and B events occur if A and B are independent events is the multiplication of the probability of A and the probability of B.

(*Probability that our first patient is male and the second one is female  is the probability that our first patient is male * the probability that our second patient is female.*)

These mentioned statements are true from other way round. For example if P(A)*(P(B)=P(AandB) then A and B are independent events.

# Probability of Events VI.

*Conditional* events: **P(A|B)\*P(B)=P(A)**
Example:

*the probability that a patient suffering from a viral infection has actually flu is* 14% *= P(A|B).*

*The probability that a patient coming to our office has viral infection is 8% =P(B)*

*The probability of occurrence of flu infections at our office:*
   *P(A)* = 0,14 * 0,08 = 0,0112 = 1,12%

There is an other important calculation that you have to know on conditional events. I try to explain it based on an example.
*The probability that a patient suffering from a viral infection has actually flu is* 14% *=P(A|B)* - that is the conditional probability
*The probability that a patient coming to our office has viral infection is 8% =P(B)* – that is the probability of the condition.
*The probability of occurrence of flu infections at our office* – that is the probability of our event in the given sample.
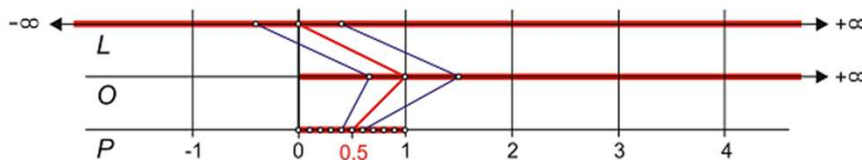      *P(A)* = 0,14 * 0,08 = 0,0112 = 1,12%

# Odds

***Odds (O)***: the ratio of the probability that a given event occurs and the probability that it does not occur. (how much larger is the probability of an event occurring than of not occurring)

$$O = \frac{P(A)}{1 - P(A)}$$

***Logit (L)***: natural logarithm of odds

Logit, Odds, Probability



There is an other „probability like" parameter in probability calculus that used very often. This parameter is the odds. Calculated as a ratio of the probability that a given event occurs and the probability that it does not occur. (The meaning is how much larger is the probability of an evet occuring than of not occuring.)
The logit is a rarely used parameter. It is the natural logarithm of odds.
The connection between the value of the logit, odds and probability is shown in the slide.
We can notice for example that the probability is between 0 and 1. The odds between 0 and infinite. The odds is over 1 if the probability is over 0.5. The logit is negative if the probability is smaller than 0.5.

# Probability Calculus

Permutations,
Variations,
Combinations

Probability calculation and statistics are based on the permutations, variations and combinations.

## *Probability Calculus Example*

During last year's flu epidemic 402 out of the total 2989 patients who turned up at a doctor's office required vaccination. Based on last year's data what is the probability that 4 vaccines will be sufficient (exactly, i.e. no vaccines will be left), if we are expecting a total number of 25 patients?

$$P = \binom{n}{k} \cdot (p)^k \cdot (1-p)^{(n-k)} = \binom{25}{4} \cdot \left(\frac{402}{2989}\right)^4 \cdot \left(1 - \frac{402}{2989}\right)^{(25-4)} \approx 0{,}2$$

How to calculate? How to read out from a graph, table?
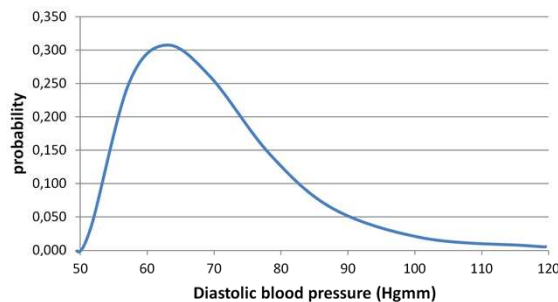Which equation, table should we use?

An example why and how we use the probability calculus.
During last year's flu epidemic 402 out of the total 2989 patients who turned up at a doctor's office required vaccination. Based on last year's data what is the probability that 4 vaccines will be sufficient (exactly, i.e. no vaccines left) in a certain day, if we are expecting a total number of 25 patients?
For answering the question we use the Bernoulli distribution's (see later) equation. I show this equation to highlight that Bernoulli distribution is based on probability calculus.
To answer question like that, our main questions will be: How to calculate? How do we know the equations? Which equation, table we should use?
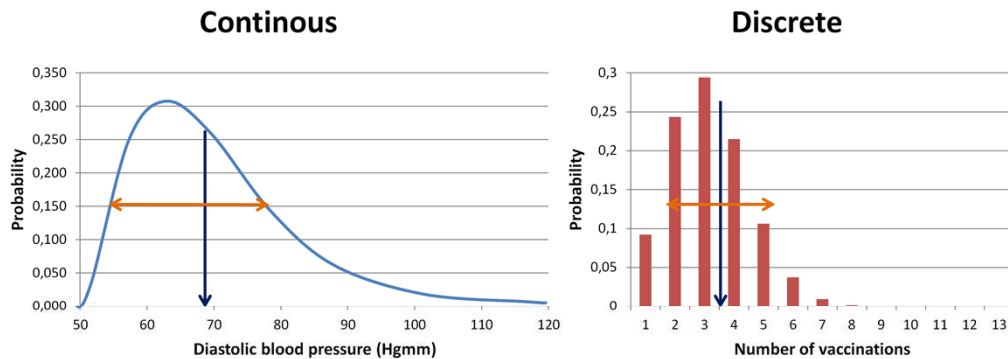
# Theoretical Distributions



I know the probability for all value based on experiments. (very rare)

**I can calculate (or estimate) the probability based on a few parameters using *special distributions.***

***What are the parameters and which distribution should I use?***

For answering the questions about probabilities we can use theoretical distributions. Theoretical distributions shows the probability for a given value. In a very rare cases we know the probabilities for every outcomes based on a large number of experiments. But usually we can calculate or more often estimate the probabilities based on a few value (parameters) using *special distributions*. So the question is what are these parameters and which special distribution I should use for the given problem.

# Parameters of Theoretical Distributions

**Continous**

Probability — Diastolic blood pressure (Hgmm)

**Discrete**

Probability — Number of vaccinations

- **Expected value(E)**

$$E(\xi) = \int_{-\infty}^{\infty} p_i \cdot x_i$$

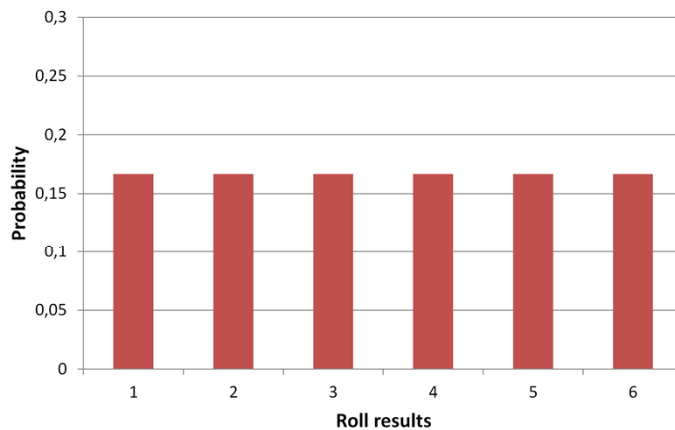$$E(\xi) = \sum_{i=1}^{m} p_i \cdot x_i$$

- Theoretical variance (Var, D$^2$)

$$Var(\xi) = E\left[\left(\xi - E(\xi)\right)^2\right]$$

Theoretical distributions have similar parameters as mentioned in descriptive statistics. There is a parameter that describe the center of the distribution and an other one that describe the width of the distribution.

The first one called *expected value* (abbreviated with E), the second is the *theoretical variance* (Var). In the equation *x* is the given value and *p* is the probability of that value. The expected value calculated slightly differently for continous and discrete variables. As I showed in the lecture the expected value is equal with the *mean of the population*. For countinous variable we use infinite small binwidth for summarisation – that is the integral (∫).

This two indicator (the expected value and the variance) defines exactly the distribution that means knowing this indicators we could calculate the probability for all value.

# Uniform Distribution

$$E(\xi) = \frac{1}{2}(a+b)$$

$$Var(\xi) = \frac{1}{12}(b-a)^2$$

Distribution of a perfect die *(eg. probability of rolling 4)*
Ideal workload distribution throughout the day
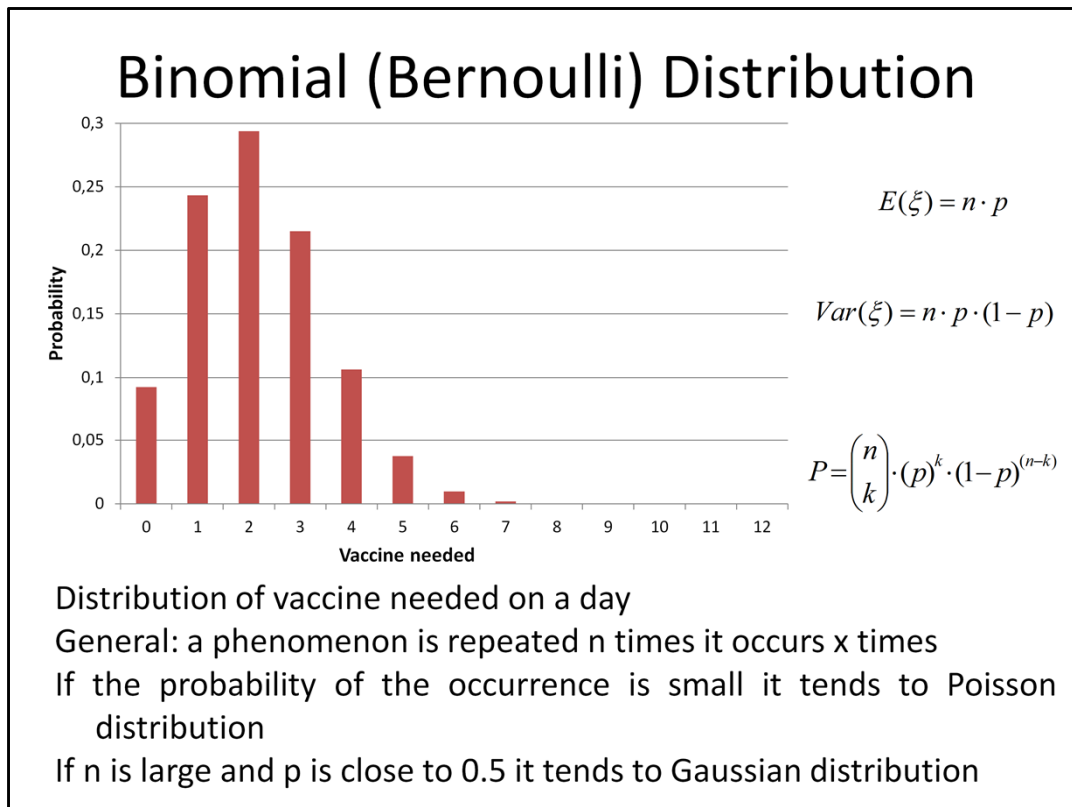Temperature distribution in an empty lecture hall

Let's see first the *uniform distribution*.
We have uniform distribution for example if we rolling a dice or we talk on the ideal workload or temperature distribution in an empty space.
For example using the uniform distribution we can calculate the probability of rolling 4 with a die.
The formula of the expected value and the variance is available in the formula collection. Here *a* and *b* are the smallest and largest outcomes. In the case of a six sided die the expected value is 0.5*(1+6)=3.5, and the variance is 1/12*(6-1)^2≈2.08.

# Binomial (Bernoulli) Distribution

$$E(\xi) = n \cdot p$$

$$Var(\xi) = n \cdot p \cdot (1-p)$$

$$P = \binom{n}{k} \cdot (p)^k \cdot (1-p)^{(n-k)}$$

Distribution of vaccine needed on a day

General: a phenomenon is repeated n times it occurs x times

If the probability of the occurrence is small it tends to Poisson distribution

If n is large and p is close to 0.5 it tends to Gaussian distribution

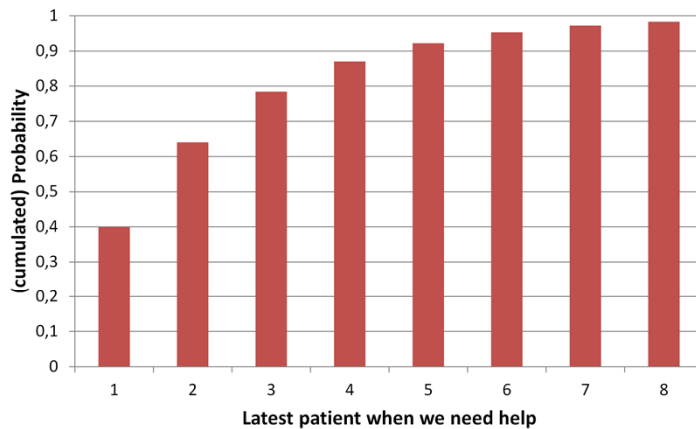The *binomial (or Bernoulli) distribution* is used in general if a phenomenon is repeated n times it occurs x times.

An exmple was the vaccination I mentioned before: During last year's flu epidemic 402 out of the total 2989 patients who turned up at a doctor's office required vaccination. Based on last year's data what is the probability that 4 vaccines will be sufficient (exactly, i.e. no vaccines left) in a certain day, if we are expecting a total number of 25 patients?

For the calculation we need the expected value and the variance – or the parameters that describe it. E=n*p and Var=n*p*(1-p) so we need *n* (in our case the expected number of paient: 25) and *p (*we could estimate it using last years data: p=402/2989*)*. Based on this values we can calculate the probability of *k* (in our case 4) using the equation of the distribution.

If the probability of occurance (p) is small the binomial distribution tends to a Poisson distribution.

If repetition (n) is large and the probability is close to 0.5 it tends to Gaussian distribution.

# Geometric Distribution

$$E(\xi) = \frac{1}{p}$$

$$Var(\xi) = \sqrt{\frac{1-p}{p^2}}$$

$$P = p \cdot (1-p)^{(n-1)}$$

Independent sequence of Bernoulli trials
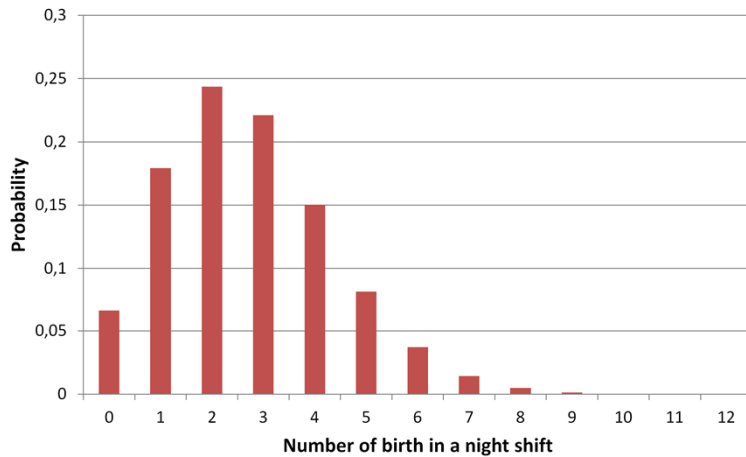When the first head occurs (St. Petersburg paradox)
The first patient when we need a nurse's to help

A geometric distribution is a special Binomial distribution. In this case we make independent sequence of Bernoulli trials.
A medical example: What is the probability that we can examine the first patient without calling the nurse to help us? Or the probability that we couldn't examine the xth people without any help before. In this graph I show the probability of this situation – this is a kind of a cumulative frequency distribution (we cumulate the probability that we need help in the first patient + we don't need help in the first, but in the second we need +we don't need help in the first two, only in the third patient....)
In the play of St. Petersburg paradox the prize of a single game follow geometric distribution too (see 2nd slide in this lecture).

## Poisson Distribution

Probability

0,3
0,25
0,2
0,15
0,1
0,05
0

Number of birth in a night shift
0 1 2 3 4 5 6 7 8 9 10 11 12

$$E(\xi) = \lambda$$

$$Var(\xi) = \lambda$$

$$P = \frac{\lambda^x}{x!} \cdot e^{-\lambda}$$

Number of births during night shift
Number of white blood cells in the field of view
Number of decayed atoms in a radioactive substance during a
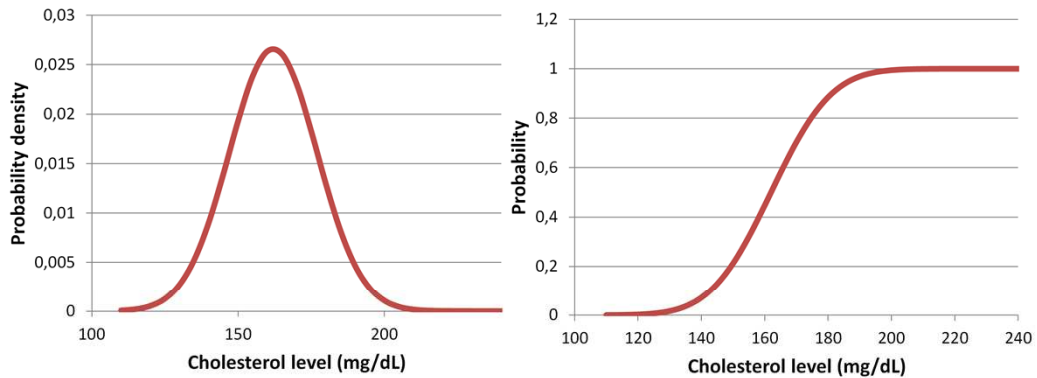given time interval

The Poisson distribution has a special attribution: the expected value and the variance are the same – so we need only 1 parameter to describe this distribution.

For example based on this distribution we can calculate the probability that we have 3 birth during our night shift. Other examples that follows Poisson distribution are: Number of white blood cells in the field of view, number of decayed atoms in a radioactive substance during a given time interval.

In general: number of elements in a given time interval or volume…, if the *probability of the occurrence is small*.

The expected value ($\lambda$) derived from *n*p* (number of „repetition" * probability).

# Normal (Gaussian) Distribution I.

Cholesterol level, glucose level.....
Height, BMI...
Diastolic blood pressure of adults
......

$$E(\xi) = \mu$$

$$Var(\xi) = \sigma^2$$

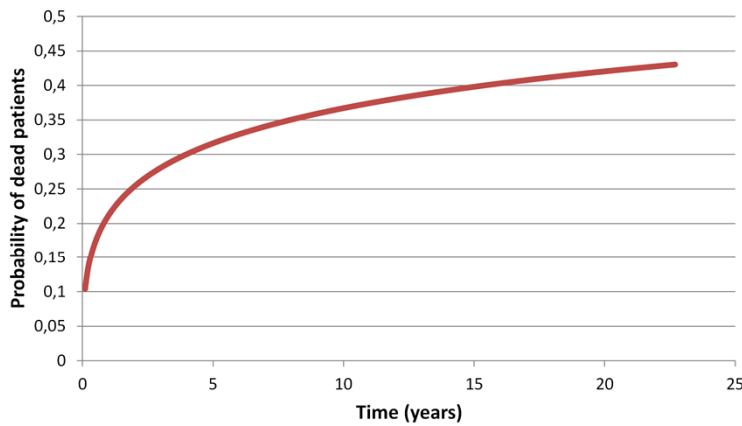$$P = \frac{1}{\sigma \cdot \sqrt{2 \cdot \pi}} \cdot e^{\frac{-(x-\mu)^2}{2 \cdot \sigma^2}}$$

The normal or Gaussian distribution is the most common in medical practice.
In this slide I ploted both the relative distribution and cumulative frequency functions, because this is the most important distribution for us. As you see against the other mentioned distribution this is symmetric one.
The most of the variables in medical practice follows normal (Gaussian) distribution – eg. enzime levels, height, body mass index (BMI), blood pressures...
Why?

# Gaussian Distribution II.

***Central limit theorem***: for given conditions, adding a large number of independent variables yields a normally distributed variable.

The reason of why we have normal distribution in most of the variables in medical practice described by the central limit theorem. It says that summarizing large number of independent variables resulted a normally distributed variable. In medical practice most of the measure values affected by several factor: gens from father, gens from mother, nutrition, way of life...

# Lognormal (Galton) Distribution
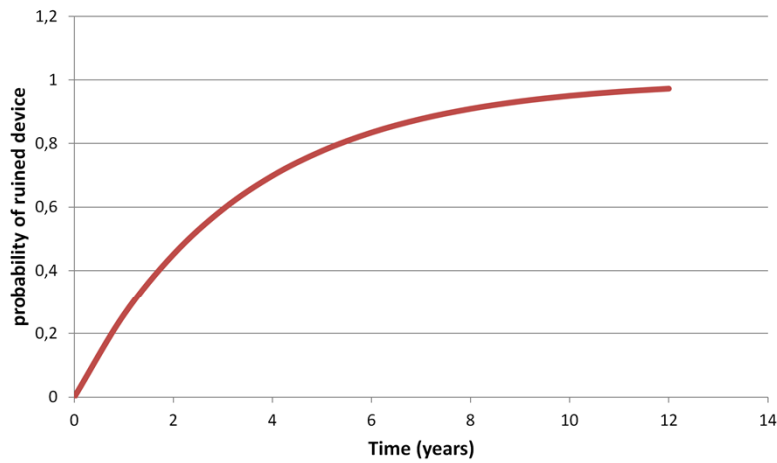
$$E(\xi) = e^{\mu + \sigma^2/2}$$

$$Var(\xi) = (e^{\sigma^2} - 1) \cdot e^{2\mu + \sigma^2}$$

$$P = \frac{1}{\sigma \cdot x \cdot \sqrt{2 \cdot \pi}} \cdot e^{\frac{-(\ln x - \mu)^2}{2 \cdot \sigma^2}}$$

Body mass and height in childhood
Survival time of a cancer

A common distribution in medical practice is the lognormal distribution. For example the body parameters in childhood, survival time of a cancer.
In general if the values of the variable are close to 0 and couldn't be negative instead of a normal distribution we get a lognormal distribution.

## Exponential Distribution

$$E(\xi) = \frac{1}{\lambda}$$

$$Var(\xi) = \frac{1}{\lambda^2}$$

$$P = \lambda \cdot e^{-\lambda \cdot x}$$

Anaesthetic equipment operating time (before the first error).
Lifetime of the individual atoms in the course of radioactive decay.

The exponential distribution is well know in biophysics and has some appearance in medical practice too. I give you two example: anaesthetic equipment operating time (before the first error) and the lifetime of the individual atoms in the course of radioactive decay.

# Human thinking and probability...

Tom is a quiet, shy, modest, hard-working guy who is happy to help others. Which is more probable?

a) Tom is a librarian
b) Tom is a blue-collar worker

Perhaps you think based on the description that Tomi is a librarian more likely but if you think it over the frequency of male librarians and male blue-collar worker you should realize that Tomi is probably a blue-collar worker and not a librarian.

# Human thinking and probability...

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

Which is more probable?
- a) Linda is a teacher in a secondary school
- b) Linda works in bookstore and participates in yoga courses
- c) Linda is a member of the league of women voters
- d) Linda is a bank teller.
- e) Linda is an insurance agent
- f) Linda is a bank teller and is active in the feminist movement.

I'd like to highlight two statement: *d* and *f*. I hope everybody found out that co-occurrence is less probable than occurrence of a given event. The intersection of sets is always equal or smaller than the sets. So it is less probable that Linda is a bank teller and active feminist than she is a bank teller.

# Test Questions #1

- Give the definition of fprobability based on relative frequences.
- What is the law of large numbers?
- How tends the relative frequencies to the probability? [fluctuations, infinite seuence]
- How we can prove the law of large numbers?.
- What is the union of two sets?
- How we can notate the probability that events A or B occur?
- How we can notate the probability that both event A and B occur at the same time?
- What is the intersection of two event?
- What does it mean mutually exclusive events?
- Give an example for mutually excusive events.
- What is the the value of intersection of two mutually exclusive events?
- What does independent events mean?
- Give an example for independent events.
- What is the conditional probability?.

- Give an example for conditional probability.
- How we could notate conditional probability?
- How to calculate P(A) if P(A|B) and P(B) is given?
- What are the Kolmogorov's axioms?
- What is the relation betweenA and B events, if P(AorB)=P(A)+P(B) is true?
- What is the relation between A and B events, if P(AB)=P(A)*P(B) is true?
- What is the probability of sure event?
- What is the probabilty of an impossible event?
- Give an example for sure and impossible events.
- What could the value of an event's probability be?
- Define the odds.
- Define the logit.
- Calculate the logit if the probability of an event is 0,12.
- Calculate the odds if the probability is 0,4.
- Calculate the probability if the odds is 3.
- Calculate the probability if the logit is – 32.

The following questions may be answered using lecture material, consultation with practice teacher, or your own investigation (on the library or the internet). These test questions are examples for multiple choice items that may occur in the midterm and exam tests.

# Test Questions #2

- How you can calculate the expected value of a continous distribution?
- How you can calculate the expected value of a discrete distribution?
- Which central tendency equal with the expected value in case of a population?
- Define the theoretical variance.
- What are the two indicators that define exactly a special distribution?
- How does the frequency distribution of a uniform distribution looks like?
- How does the frequency distribution of a Poisson distribution looks like?
- How does the frequency distribution of a Bernoulli distribution looks like?
- How does the frequency distribution of a Geometric distribution looks like?
- How does the frequency distribution of a Gaussian distribution looks like?
- How does the cumulative frequency distribution of a Gaussian distribution looks like?
- How does the frequency distribution of a exponential distribution looks like?
- How does the frequency distribution of a lognormal distribution looks like?

- Give two example for uniform distribution.
- Give two example for binomial distribution.
- Give two example for poisson distribution.
- Give two example for normal distribution.
- Give two example for lognormal distribution.
- Give two example for geometric distribution.
- Give two example for exponential distribution.
- How you can calculate the expected value of a uniform distribution?
- How you can calculate the expected value of a binomial distribution?
- How you can calculate the expected value of a lognormal distribution?
- How you can calculate the expected value of a exponential distribution?
- How you can calculate the expected value of a Poisson distribution?
- How you can calculate the expected value of a Gaussian distribution?
- What is the central limit theorem?
- Why are the most of the medical variables normally distributed?

# Test Questions #3

- Give a general describtion when we get a binomial distribution.
- Give a general describtion when we get a Poisson distribution.
- When we get lognormal distribution instead of normal distribution?
- How to convert lognormal distribution to normal distribution?
- Could be the co-occurance larger than the occurrence of one of the given events?