

Bioinformatics and Biostatistics BB2440: Biostatistics

Lecture 2: Bayes' Probability & Random Variables & Some Probability Distributions

Timo Koski

TK

04.09.2013



Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)



KTH Matematik

Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events



Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)

Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables



Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables
- Probability Distributions

Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables
- Probability Distributions
 - Binomial Distribution

Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables
- Probability Distributions
 - Binomial Distribution
 - Geometric Distribution



Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables
- Probability Distributions
 - Binomial Distribution
 - Geometric Distribution
 - Poisson Distribution

Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables
- Probability Distributions
 - Binomial Distribution
 - Geometric Distribution
 - Poisson Distribution
 - Geometric Distribution

Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables
- Probability Distributions
 - Binomial Distribution
 - Geometric Distribution
 - Poisson Distribution
 - Geometric Distribution
 - Mean and Variance



Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables
- Probability Distributions
 - Binomial Distribution
 - Geometric Distribution
 - Poisson Distribution
 - Geometric Distribution
 - Mean and Variance
 - Normal Distribution

Outline of Lecture 2.

- A Quick Summary of the Formal Rules of Probability (from Lect 1.)
- Conditionally Probable Events
- Bayes' Rule of Probability (Inverse Conditional Probability)
- Random Variables
- Probability Distributions
 - Binomial Distribution
 - Geometric Distribution
 - Poisson Distribution
 - Geometric Distribution
 - Mean and Variance
 - Normal Distribution
 - Exponential Distribution



Probability Summary

Pr denotes probability

A, B, and C denote specific events.

Pr(A) denotes the probability of event A occurring

The probability of an event A , $Pr(A)$, has the following properties:

- (a) $0 \leq Pr(A) \leq 1$;
- (b) $Pr(S) = 1$, if S is an event that is certain to happen (e.g., the sun rises).



Probability Summary

We say that two events A and B are **mutually exclusive**, if the occurrence of one precludes the occurrence of the other.

The event A or B means that A occurs or both A and B occur or B occurs.

- (c) if A and B are mutually exclusive events then we have the **additive law**

$$Pr(A \text{ or } B) = Pr(A) + Pr(B).$$

(d)

$$Pr(A \text{ or } B) = Pr(A) + Pr(B) - Pr(A \text{ and } B)$$

Note that if A and B are mutually exclusive, then the event " A and B " is impossible, and thus $Pr(A \text{ and } B) = 0$.

Complement " A does not occur" $= A^*$.

(e)

$$Pr(A^*) = 1 - Pr(A).$$

Probability Summary : independent events

We say that A and B are **independent** events if

(f)

$$Pr (A \text{ and } B) = Pr(A) \cdot Pr(B).$$

*Now we introduce a new element: **Conditionally Probable Events***

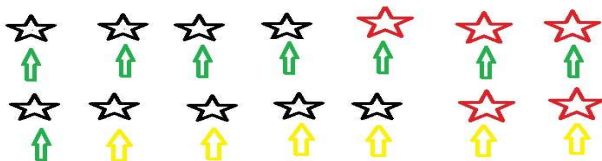


Conditionally Probable Events: a Mendelian introduction



Conditionally Probable Events: a Mendelian introduction

In the figure we have 14 (abstract) peas: nine have white flowers, five have red flowers, six have yellow pods, eight have green pods.



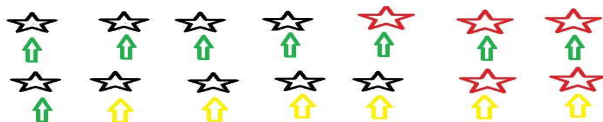
Two peas are selected at **random** and **without replacement**. **Q:** What is the probability for the event that the first selection has a green pod and the second selection has a yellow pod ?

Conditionally Probable Events: a Mendelian introduction



- First selection: $Pr(\text{green pod}) = \frac{8}{14}$

Conditionally Probable Events: a Mendelian introduction



- First selection: $Pr(\text{green pod}) = \frac{8}{14}$
- Second selection: $Pr(\text{yellow pod}) = \frac{6}{13}$

Conditionally Probable Events: a Mendelian introduction



- First selection: $Pr(\text{green pod}) = \frac{8}{14}$
- Second selection: $Pr(\text{yellow pod}) = \frac{6}{13}$
- $Pr(\text{1st pea with green pod and 2nd pea yellow pod}) = \frac{8}{14} \cdot \frac{6}{13} \approx 0.264$.

Conditionally Probable Events

Two peas are selected at **random** and **without replacement**.

$$Pr(\text{1st pea green pod and 2nd pea with yellow pod}) = \frac{8}{14} \cdot \frac{6}{13} \approx 0.264.$$

The key point here is that *we must adjust the probability of the second event to reflect the outcome of the first event*. One pea was removed in the first selection, there are only 13 left to choose for the second selection.

Conditional Probability Events: The Principle & Notation

(Notation for Conditional Probability)

$Pr(B \mid A)$ represents the probability of the event B occurring after it is assumed that the event A has already occurred .

We read $B \mid A$ as " B given A " or as " event B after that event A has already occurred " .

Multiplication Rule

$$Pr (A \text{ and } B) = Pr(A) \cdot Pr(B | A)$$

But also

$$Pr (A \text{ and } B) = Pr(B) \cdot Pr(A | B)$$

Confusing !

Conditionally Probable Events

Now we have in fact already used the multiplication rule

$$\begin{aligned} &Pr(\text{1st pea green pod and 2nd pea yellow pod}) = \\ &Pr(\text{1st pea green pod}) \cdot Pr(\text{2nd pea yellow pod} \mid \text{1st pea green pod}) = \\ &\quad \frac{8}{14} \cdot \frac{6}{13} \end{aligned}$$

Multiplication Rule

If $Pr(A) > 0$

$$\frac{Pr(A \text{ and } B)}{Pr(A)} = Pr(B | A)$$

and if $Pr(B) > 0$

$$\frac{Pr(A \text{ and } B)}{Pr(B)} = Pr(A | B)$$



KTH Matematik

Probability: independent events

$$Pr (A \text{ and } B) = Pr(A) \cdot Pr(B | A)$$

But if A and B are **independent** events

$$Pr (A \text{ and } B) = Pr(A) \cdot Pr(B).$$

Hence, if A and B are independent events

$$P(B | A) = Pr(B),$$

i.e., we need not adjust the probability of B to reflect the outcome of A !

We consider a syndrom of disease and a diagnostic test. With positive result we mean that the test indicates the presence of the syndrom, in case of negative result the test indicates the absence of the syndrom. The table gives the probabilities of various outcomes. To explain it we have for example

$$Pr(\text{Positive Result and Syndrom}) = 0.81$$

	Positive Result	Negative Result
Syndrom	0.81	0.05
No syndrom	0.03	0.11

	Positive Result	Negative Result
Syndrom	0.81	0.05
No syndrom	0.03	0.11

$$Pr(\text{Positive Result} \mid \text{Syndrom}) = \frac{Pr(\text{Positive Result and Syndrom})}{Pr(\text{Syndrom})}$$

By the additive law $Pr(\text{Syndrom}) = 0.81 + 0.05 = 0.86$. Thus

$$Pr(\text{Positive Result} \mid \text{Syndrom}) = \frac{0.81}{0.86} \approx 0.94$$

	Positive Result	Negative Result
Syndrom	0.81	0.05
No syndrom	0.03	0.11

$$Pr(\text{Syndrom} \mid \text{Positive Result}) = \frac{Pr(\text{Positive Result and Syndrom})}{Pr(\text{Positive Result})}$$

By the additive law $Pr(\text{Positive Result}) = 0.81 + 0.03 = 0.84$. Thus

$$Pr(\text{Syndrom} \mid \text{Positive Result}) = \frac{0.81}{0.84} \approx 0.96$$

Hence

$$Pr(\text{Positive Result} \mid \text{Syndrom}) \neq Pr(\text{Syndrom} \mid \text{Positive Result}).$$

Inversion: Th. Bayes 1702 - 1761

We shall next relate $Pr(A | B)$ and $Pr(B | A)$ using **Bayes' Rule**



Bayes' Rule

Bayes' rule goes from probability of $B \mid A$ to probability of $A \mid B$. (Inverts the roles of A and B as conditionally probable events)

Bayes' Rule

$$Pr(A \mid B) = \frac{Pr(A)Pr(B \mid A)}{Pr(A)Pr(B \mid A) + Pr(A^*)Pr(B \mid A^*)}.$$

This is a formidable expression, but is relatively easy to use.



Bayes' Rule

You can check (exercise)

$$Pr(\text{Syndrom} \mid \text{Positive Result}) \approx 0.96$$

by using Bayes' Rule and the probabilities above

$$\frac{Pr(\text{Syndrom})Pr(\text{Pos Res} \mid \text{Syndrom})}{Pr(\text{Syndrom})Pr(\text{Pos. Res} \mid \text{Syndrom}) + Pr(\text{No Syndrom})Pr(\text{Pos. Res} \mid \text{No Syndrom})}$$



Bayes' Rule & Cancer Test

Let us assume that 0.4 % of the population has a form of cancer. In biostatistics 0.4 % is called the **prevalence** of this disease, we write $Pr(\text{Cancer}) = 0.004$. Biomedical studies have shown that for a diagnostic test $Pr(\text{Positive Result} \mid \text{Cancer}) = 0.995$ and $Pr(\text{Positive Result} \mid \text{No Cancer}) = 0.01$.

We take at random an individual from the population and the diagnostic test gives a positive result. We want to find

$$Pr(\text{Cancer} \mid \text{Positive Result})$$

The rule of the complement gives

$$Pr(\text{No Cancer}) = 1 - Pr(\text{Cancer}) = 0.996.$$



Bayes' Rule: Numbers

By Bayes' rule

$$Pr(\text{Cancer} \mid \text{Positive Result}) =$$

$$\frac{Pr(\text{Canc})Pr(\text{Pos Res} \mid \text{Canc})}{Pr(\text{Canc})Pr(\text{Pos. Res} \mid \text{Canc}) + Pr(\text{No Canc})Pr(\text{Pos. Res} \mid \text{No Canc})}.$$

We insert the numbers from the above

$$Pr(\text{Cancer} \mid \text{Positive Result}) = \frac{0.004 \cdot 0.995}{0.004 \cdot 0.995 + 0.996 \cdot 0.01} \\ \approx 0.29$$

What does this mean !?

Data → Statistics → Information/Knowledge

Bayes' rule : screening tests

John Allen Paulos: The Math behind Screening Tests. What a positive result really means. Scientific American, January 2012.

<http://www.scientificamerican.com/article.cfm?id=weighing-the-positives>



Bayes' rule : vocabulary

$$Pr(A | B) = \frac{Pr(A)Pr(B | A)}{Pr(A)Pr(B | A) + Pr(A^*)Pr(B | A^*)}.$$

The probability $Pr(A)$ is called the **prior probability** for A . The probability $Pr(A | B)$ is called the **posterior probability** for A . Bayes' rule shows the priori probability is transformed to the posterior probability.

Data → Statistics → Information/Knowledge

Bayes' rule & Bioinformatics

Wilkinson, D.J. :Bayesian methods in bioinformatics and computational systems biology. Briefings in bioinformatics, pp. 109–116, 8, 2007.



Heredity

Mendel and His Peas

Recall from Lect. 1:

$$Pr(50 \leq \text{the number of matches in a box} \leq 53) = 0.80$$

(here match= tändsticka in Swedish) or

$$Pr(11 \text{ matches in two sequences of 15 nucleotides} = 11)$$

(here match= något som matchar in Swedish)

Random Variables: a technical device

Things become easier if we define

$X =$ the number of matches in a box

and then write

$$Pr(50 \leq X \leq 53) = 0.80$$

and

$$Pr(X = 50)$$

or any other probability connected to the counting of matches in the matchbox experiment in Lect. 1..



Random Variables: Definition

In the same way

$X =$ number of matches in two sequences of 15 nucleotides

and write

$$Pr(X = 11)$$

Random Variables & Probability Distributions

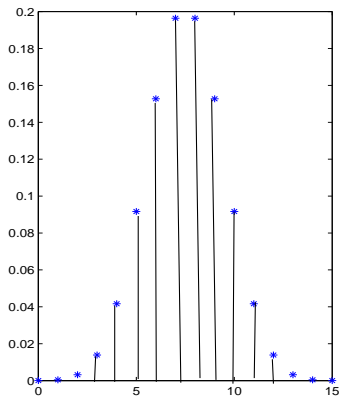
Definition

A **random variable** is a variable represented by X (or Y, Z, \dots) that has a single numerical value, determined by chance, for each outcome of a procedure.

Definition

A **probability distribution** is graph, table, formula or algorithm that gives the probability of each value of the random variable.

Binomial probability distribution with $n = 15$ and $p = 0.2$, $p = 0.5$ as a graph



Definition

A **probability distribution** is graph, table, formula or algorithm that gives the probability of each value of the random variable.

A Formula

$$Pr(X = k), \quad k = 0, 1, 2, \dots,$$

Note that

$$0 \leq Pr(X = k) \leq 1, \quad \sum_{k=0}^{\infty} Pr(X = k) = 1.$$

Binomial Distribution

Definition

A random variable X has the binomial distribution with the parameters n och p , if

$$Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}, \text{ för } k = 0, 1, \dots, n.$$

We write this as $X \sim \text{Bin}(n, p)$ and say that X is $\text{Bin}(n, p)$ - distributed.

Binomial Distribution: The Conditions

A **binomial distribution** results from a procedure that meets all the following conditions:

The procedure has a fixed number of random events.

1

The events have outcomes in two categories.

2

The events are independent.

3

The probabilities are constant for each event.

4

Binomial Distribution: Example

X = number of matches in two sequences of 15 nucleotides

If we are willing to assume that the nucleotides are random (DNA dice !)
and independent, and that the probabilities are the same at each site, then
the probability distribution of the number of matches (=successes) in an
alignment of two sequences is a binomial distribution with parameters
 $n = 15$, $p = \frac{1}{4}$.

$$X \sim \text{Bin}(15, 1/4)$$

If a procedure meets all the conditions of a binomial distribution except that the number of events is not fixed in advance, then the **geometric distribution** can be used.

$$Pr(X = k) = p \cdot (1 - p)^{k-1}, k = 1, 2, \dots,$$

$$X \sim \text{Ge}(p).$$

Geometric Distribution: an Example

Think of tossing the DNA dice with $p = Pr(G)$. If X = the number of tosses of before you get G for the first time **including** the successful toss. Then $X \sim \text{Ge}(1/4)$ and

$$Pr(X = k) = \frac{1}{4} \cdot \left(\frac{3}{4}\right)^{k-1}, k = 1, 2, \dots,$$

Geometric Distribution: an Example

Usually in statistics p is the probability of success in $\text{Ge}(p)$.
In bioinformatics (e.g., in the theory underlying BLAST, c.f., Lect. 10)
 $(1 - p)$ is the probability of success.

Then

$$\Pr(X = k) = p \cdot (1 - p)^{k-1}, k = 1, 2, \dots,$$

*is the probability that there were $k - 1$ successes (a **success run** of length $k - 1$) before the first failure at event k . X = number of successes plus the first failure.*

Binomial Distribution for small p and large n

In bioinformatics it happens often that $X \sim \text{Bin}(n, p)$ with p being small and n being large, or $p = \lambda/n$. Then

$$\binom{n}{k} p^k (1-p)^{n-k} = \frac{n(n-1)\dots(n-k+1)}{k!} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ \approx \frac{\lambda^k}{k!} e^{-\lambda}.$$

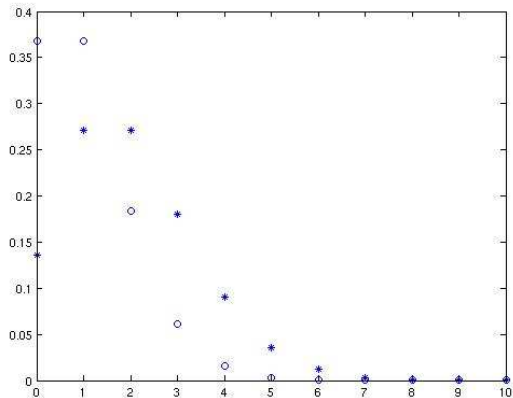
Definition

A random variable X has a Poisson distribution with parameter $\lambda > 0$, if

$$\Pr(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}, \text{ för } k = 0, 1, 2, \dots$$

We write this with $X \sim \text{Po}(\lambda)$ and say that X is $\text{Po}(\lambda)$ -distributed.

$Pr(X = k)$ for $Po(2)$ och $Po(1)$



* $\leftrightarrow Po(2)$, ○ $\leftrightarrow Po(1)$

Poisson distribution

Poisson distribution applies to occurrences of some event over a specified "unit". The random variable X is the number of random occurrences of the event in that "unit". The "unit" can be time, distance, area, volume or similar.

Example: The number of patients arriving the emergency room on Fridays between 10.00 p.m. and 11.00 p.m.

Note that

$$\begin{aligned}\sum_{k=0}^{\infty} \Pr(X = k) &= \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} \\ &= e^{-\lambda} \underbrace{\sum_{k=0}^{\infty} \frac{\lambda^k}{k!}}_{=e^{\lambda}} = e^{-\lambda} \cdot e^{\lambda} = 1.\end{aligned}$$

Poisson distribution & Sequencing

Genomic projects are generally guided by probabilistic models. In shotgun sequencing a large DNA molecule is broken into a collection of *fragments* which are sampled as randomly as possible.

The fragments are cloned and sequenced individually. The sampled fragments are then used for assembly by determining their relative orientations and overlaps and aligned to form a column-by-column matrix.

Poisson distribution & Sequencing

Suppose that a collection of fragment sequences $\mathbf{x}^1, \dots, \mathbf{x}^m$ has been aligned by some procedure. The symbol \emptyset is used as a place holder for nonaligned positions beyond the ends of the fragment. The depth of coverage of position i is defined as the number of fragments contributing sequence information at position i

$$d = \sum_{l=1}^t I(x_i^l \neq \emptyset),$$

where $I(x_i^l \neq \emptyset) = 1$ if $x_i^l \neq \emptyset$ and zero otherwise. In addition

$$\bar{p} = \sum_{x_j \neq x_l, x_j \in \mathcal{X}} Pr(x_i^l = x_j \mid x_l),$$

where $Pr(x_i^l = x_j \mid x_l)$ is the conditional probability that the shotgun sequenced fragment \mathbf{x}^l is equal to x_j in position i , given that the true sequence symbol is x_l . It is assumed that \bar{p} is the same for all x_l and i .



Poisson distribution & Sequencing

For a position covered at depth d , we compute the probability that at most one half of the bases are correct. This is the binomial probability given by

$$Pr(\text{error}|d) = \sum_{k=0}^{\lfloor d/2 \rfloor} \binom{d}{k} \bar{p}^{d-k} (1 - \bar{p})^k,$$

where the probability of success is taken as $1 - \bar{p}$ and $\lfloor d/2 \rfloor$ is largest integer smaller than or equal to $d/2$.

Poisson distribution & Sequencing

In shotgun sequencing the distribution of d can be taken as $\text{Po}(\lambda)$. Positions not covered by any fragments are ignored and thus d is at least 1. Then we obtain

$$Pr(\text{error}) = \frac{1}{1 - e^{-\lambda}} \sum_{d=1}^{\infty} \frac{\lambda^d}{d!} e^{-\lambda} Pr(\text{error}|d) \quad (1)$$

as a practical measure of the accuracy of shotgun sequencing.

The strategy for the shotgun approach to sequencing follows

*Lander E.S., Waterman M.S. Genomic mapping by fingerprinting random clones: a mathematical analysis *Genomics* 2 (3): 231-239 (1988)*

The "Lander-Waterman model" provides statistical estimates for the read depth in a whole genome shotgun (WGS) sequencing experiment via the Poisson approximation to the Binomial distribution. Although originally intended for estimating the redundancy when mapping by fingerprinting random clones, the Lander-Waterman model has served as an essential tool for estimating sequencing requirements for modern WGS experiments.

Discrete Data, Discrete Random Variables

Binomial, Geometric and Poisson distributions are examples of probability distributions for discrete data (represented by respective discrete random variables)

A **discrete random variable** *has either a finite number of values or countable number of values, where "countable" refers to the fact that there might be infinitely many values, but they can be associated with a counting process.*

Mean and Variance of Discrete Random Variables

The **mean** of a discrete random variable X is denoted by $E(X)$ and is defined by

$$E(X) \stackrel{\text{def}}{=} \sum_{k=0}^{\infty} kPr(X = k).$$

The **variance** of a discrete random variable X is denoted by $V(X)$ and is defined by

$$V(X) \stackrel{\text{def}}{=} \sum_{k=0}^{\infty} (k - E(X))^2 Pr(X = k).$$

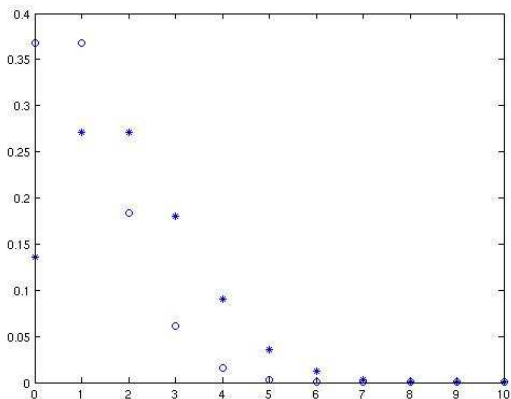
The mean is a measure of the central tendency of the random variable and the variance is a measure of dispersion

Mean and Variance of Discrete Random Variables

- $X \sim \text{Bin}(n, p)$: $E(X) = np$, $V(X) = np(1 - p)$
- $X \sim \text{Ge}(p)$: $E(X) = 1/p$, $V(X) = (1 - p)/p^2$
- $X \sim \text{Po}(\lambda)$: $E(X) = \lambda$, $V(X) = \lambda$

$Pr(X = k)$ for $Po(2)$ och $Po(1)$

We see that λ influences the location and the shape of the mass of



probability.

* $\leftrightarrow Po(2)$, ○ $\leftrightarrow Po(1)$

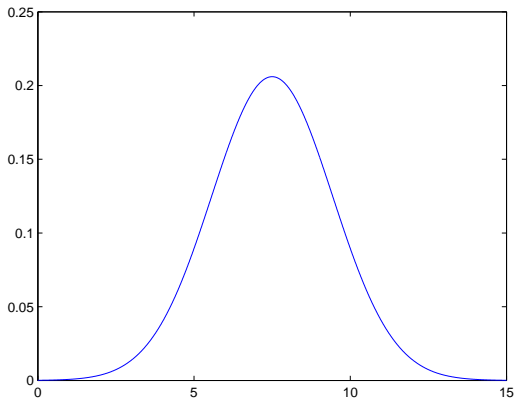
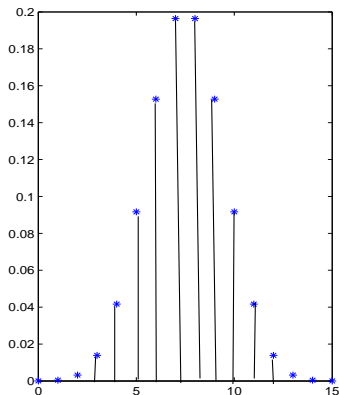
A **continuous random variable** *has infinitely many values, and those values can be associated with measurements on a continuous scale without gaps or interruptions.*

Continuous random variables can often be obtained as approximative forms of discrete distributions in a natural way.

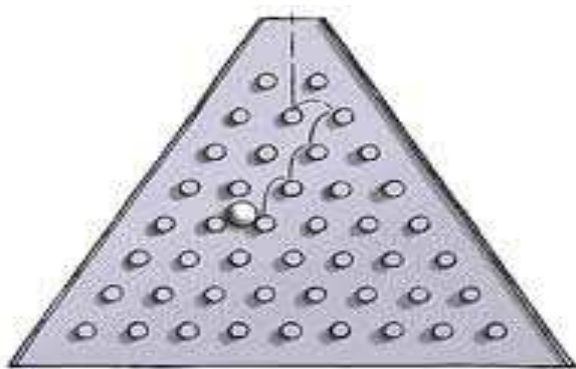
Binomial distribution: approximation

For $p = 0.5$ and $n = 15$ the Binomial distribution with mean $15 \cdot 0.5 = 7.5$, and variance $15 \cdot 0.5 \cdot 0.5 = 3.75$. The curve in the right hand figure $\mu = 7.5$ and $\sigma^2 = 3.75$.

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

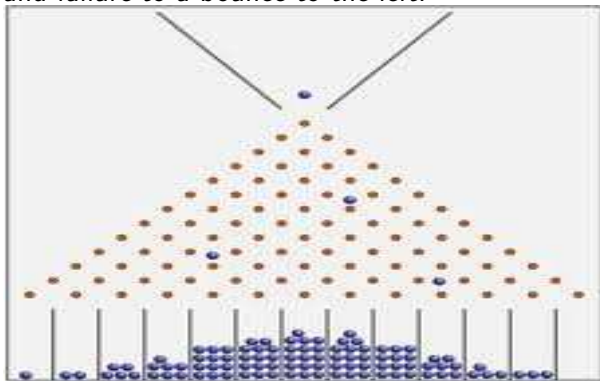


The Galton Board

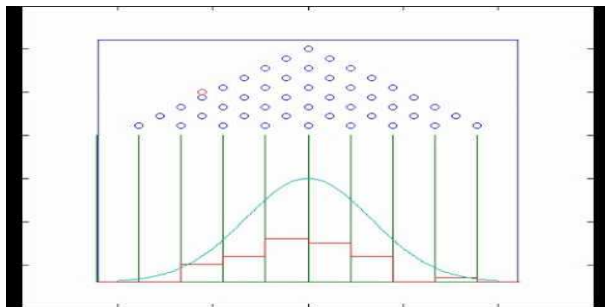


The Galton Board Experiment

The Galton board experiment consists of performing n trials with probability of success p . The trial outcome are represented graphically as a path in the Galton board: success corresponds to a bounce to the right and failure to a bounce to the left.



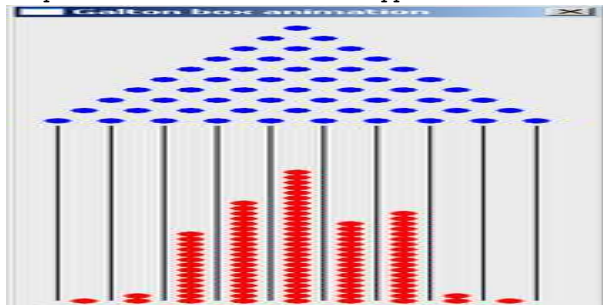
The Galton Board Experiment: Bell Shape



The Galton Board Experiment

An applet for the Galton board experiment

<http://www.math.uah.edu/stat/applets/GaltonBoardExperiment.html>



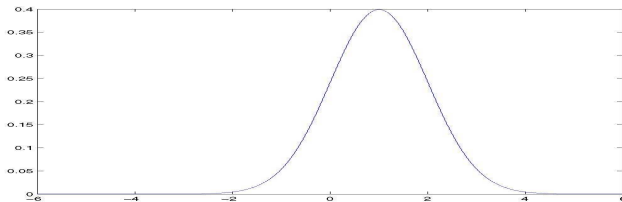
Normal (Probability) Density Curve, the Bell Shape

We call the function $f(x)$, $\sigma > 0$,

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

the normal (probability) density or the normal (probability) density curve.

In the figure $\mu = 1, \sigma = 1$



Normal Density Curve

We have

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

i.e.,

$$\int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2} dx = 1.$$

A Density Curve

A density curve $f(x)$ is graph of a continuous probability distribution. It must satisfy the following properties

- Total area under the curve must equal 1:

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

A Density Curve

A density curve $f(x)$ is graph of a continuous probability distribution. It must satisfy the following properties

- Total area under the curve must equal 1:

$$\int_{-\infty}^{\infty} f(x)dx = 1$$

- Every point on the curve must have a vertical height that is 0 or greater (that is, the curve cannot fall under the x -axis).

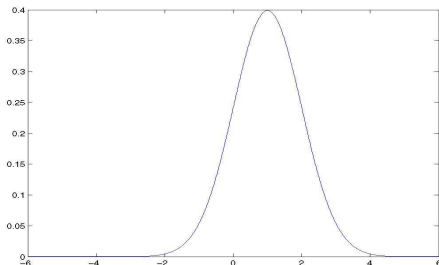
A Density Curve

A density curve $f(x)$ is graph of a continuous probability distribution. It must satisfy the following properties

- Total area under the curve must equal 1:

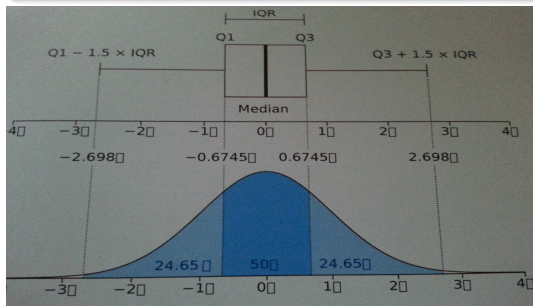
$$\int_{-\infty}^{\infty} f(x) dx = 1$$

- Every point on the curve must have a vertical height that is 0 or greater (that is, the curve cannot fall under the x-axis).



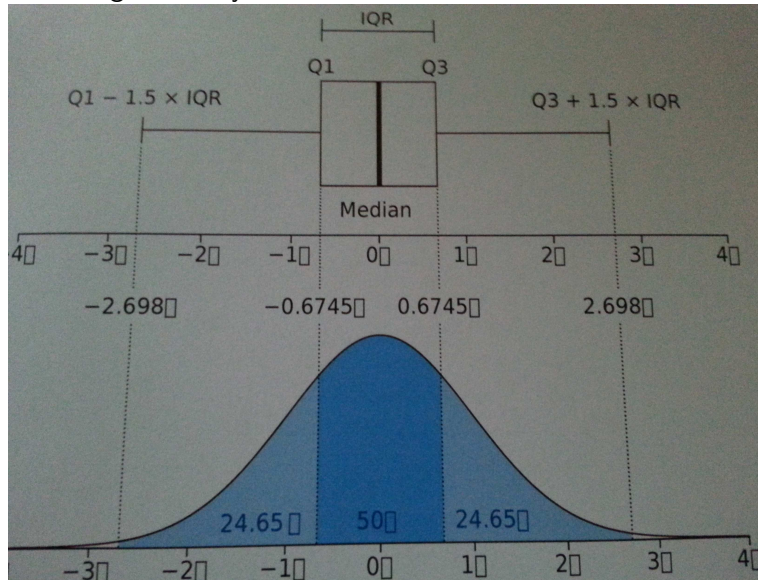
A Density Curve

Because the total area under a density curve is equal to 1, there is correspondence between area and probability.



A Density Curve

In this figure the symbol σ should be read as σ .



Probability as Area under Density Curve

X is continuous and represented by a density curve $f(x)$. Then

$$Pr(a < X \leq b) = \int_a^b f(x) dx$$



Normal Distribution (a.k.a. Gaussian Distribution)

Definition

A continuous random variable X is said to have the $N(\mu, \sigma)$ - distribution, where $\sigma > 0$, if

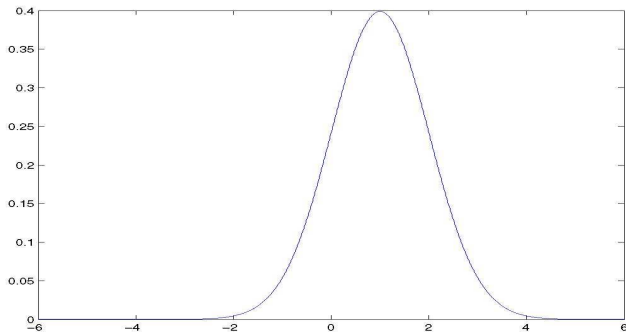
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

We write $X \sim N(\mu, \sigma)$.

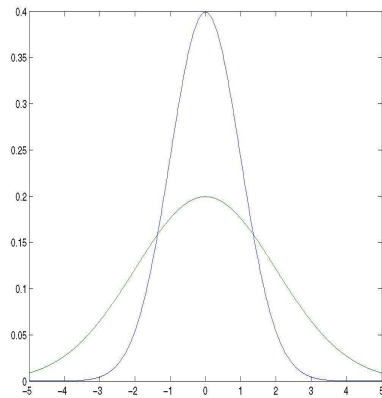
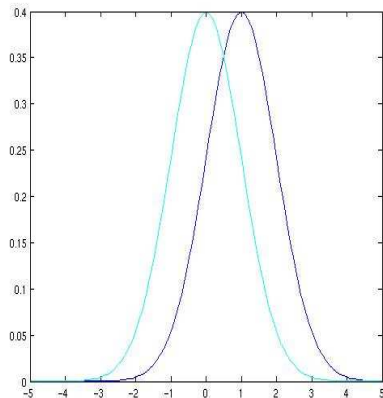
Normal Distribution (a.k.a. Gaussian Distribution)

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

In the figure $\mu = 1, \sigma = 1$

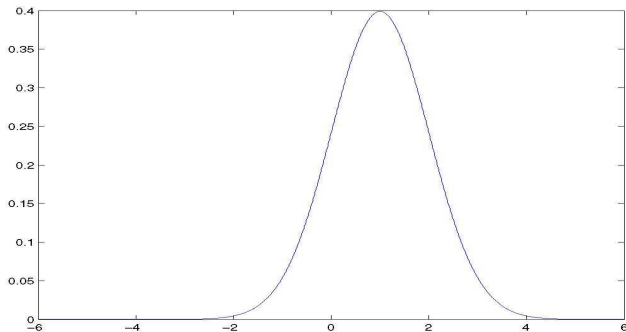


$N(0, 1)$ och $N(1, 1)$ och $N(0, 1)$ och $N(0, 2)$ (from left to the right)



Normal Distribution (a.k.a. Gaussian Distribution)

In the figure $\mu = 1, \sigma = 1$



Normal Distribution: Probability as Area under Density Curve

$X \sim N(\mu, \sigma)$. Then

$$\begin{aligned}Pr(a < X \leq b) &= \int_a^b f(x) dx \\&= \frac{1}{\sigma\sqrt{2\pi}} \int_a^b e^{-(x-\mu)^2/2\sigma^2} dx\end{aligned}$$

This integral must be evaluated by calculator, computer software or old-fashioned tables.

Normal Distribution: Probability as Area under Density Curve

In fact the method of tables corresponds to computing values of

$$\frac{1}{\sigma\sqrt{2\pi}} \int_{(a-\mu)/\sigma}^{(b-\mu)/\sigma} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz$$

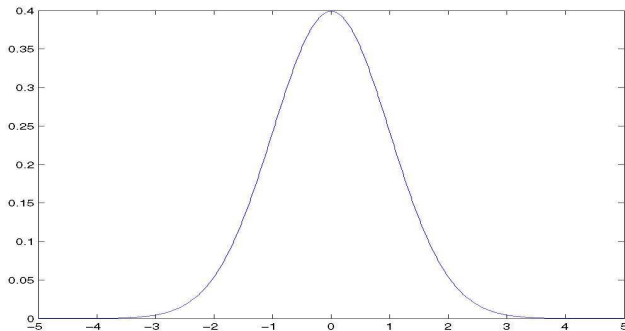
How is this ?

Standard Normal distribution, $N(0, 1)$

The distribution $N(0, 1)$ density

$$\varphi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}.$$

is called standard normal distribution



Standard Normal distribution, $N(0, 1)$

Definition

En random variable Z is said to have standard normal distribution if $Z \sim N(0, 1)$, i.e., if it has the density/density curve

$$\varphi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}.$$

We have

$$Pr(Z \leq z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

$$Pr(Z \leq 0) = \frac{1}{2}.$$

A Table for $Pr(Z \leq z)$

Standard normal, cumulative density, $P(Z < z)$

	0	0.01	0.02	0.03	0.04	0.05	0.06	0.07
0	0.500	0.504	0.508	0.512	0.516	0.520	0.524	0.528
0.1	0.540	0.544	0.548	0.552	0.556	0.560	0.564	0.568
0.2	0.579	0.583	0.587	0.591	0.595	0.599	0.603	0.607
0.3	0.618	0.622	0.626	0.629	0.633	0.637	0.641	0.645
0.4	0.655	0.659	0.663	0.666	0.670	0.674	0.677	0.681
0.5	0.691	0.695	0.698	0.702	0.705	0.709	0.712	0.716
0.6	0.726	0.729	0.732	0.736	0.739	0.742	0.745	0.749
0.7	0.758	0.761	0.764	0.767	0.770	0.773	0.776	0.779
0.8	0.788	0.791	0.794	0.797	0.800	0.802	0.805	0.808
0.9	0.816	0.819	0.821	0.824	0.826	0.829	0.831	0.834
1	0.841	0.844	0.846	0.848	0.851	0.853	0.855	0.858
1.1	0.864	0.867	0.869	0.871	0.873	0.875	0.877	0.879
1.2	0.885	0.887	0.889	0.891	0.893	0.894	0.896	0.898
1.3	0.903	0.905	0.907	0.909	0.910	0.911	0.913	0.915

Normal Distribution: Probability as Area under Density Curve

$X \sim N(\mu, \sigma)$. Set

$$Z = \frac{X - \mu}{\sigma}.$$

Then $Z \sim N(0, 1)$ and thus

$$\begin{aligned} Pr(a < X \leq b) &= Pr\left(\frac{a - \mu}{\sigma} < Z \leq \frac{b - \mu}{\sigma}\right) = \\ &= \frac{1}{\sqrt{2\pi}} \int_{\frac{a - \mu}{\sigma}}^{\frac{b - \mu}{\sigma}} e^{-x^2/2} dx \end{aligned}$$

Normal Distribution: Probability as Area under Density Curve

$X \sim N(1, 1)$. Set

$$Z = \frac{X - 1}{1}.$$

Then

$$Pr(X \leq 1.1) = Pr(Z \leq 1.1 - 1) = Pr(Z \leq 0.1) = 0.54$$

Standard normal, cumulative density, $P(Z \leq z)$								
	0	0.01	0.02	0.03	0.04	0.05	0.06	0.07
0	0.500	0.504	0.508	0.512	0.516	0.520	0.524	0.528
0.1	0.540	0.544	0.548	0.552	0.556	0.560	0.564	0.568
0.2	0.579	0.583	0.587	0.591	0.595	0.599	0.603	0.607
0.3	0.618	0.622	0.626	0.629	0.633	0.637	0.641	0.645
0.4	0.655	0.659	0.663	0.666	0.670	0.674	0.677	0.681
0.5	0.691	0.695	0.698	0.702	0.705	0.709	0.712	0.716
0.6	0.726	0.729	0.732	0.736	0.739	0.742	0.745	0.748
0.7	0.758	0.761	0.764	0.767	0.770	0.773	0.776	0.779
0.8	0.788	0.791	0.794	0.797	0.800	0.802	0.805	0.808
0.9	0.816	0.819	0.821	0.824	0.826	0.829	0.831	0.834

Exponential distribution $\text{Exp}(\lambda)$

$X \in \text{Exp}(\lambda)$

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{för } x \geq 0, \\ 0 & \text{för } x < 0. \end{cases}$$

Exponential distribution $\text{Exp}(\lambda)$

$X \in \text{Exp}(\lambda)$. Then for $x > 0$

$$\Pr(X \leq x) = \int_0^x f(u) du = \int_0^x \lambda e^{-\lambda u} du = 1 - e^{-\lambda x}.$$

By the law of the complement

$$\Pr(X \geq x) = 1 - \Pr(X \leq x) = e^{-\lambda x}.$$

In bioinformatics (sequence alignment) one often approximates the distribution of the alignment score S by the exponential distribution so that

$$\Pr(S \geq x) = e^{-\lambda x}.$$

End of Lecture 2

