

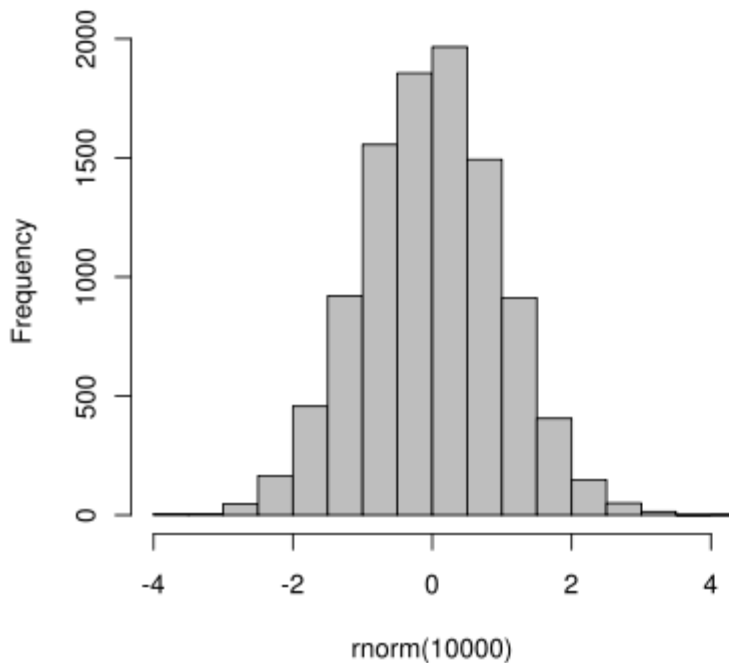
# BASICS ON DISTRIBUTIONS



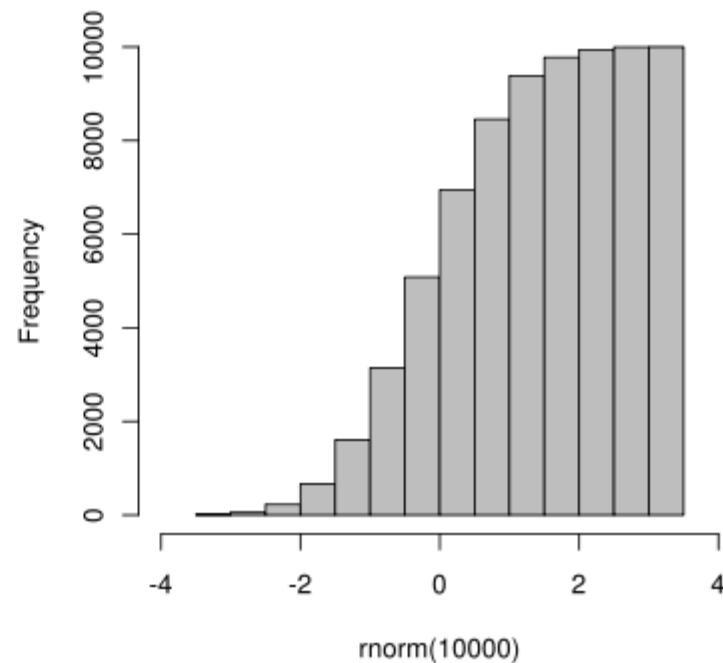
# Histograms

Consider an experiment in which different outcomes are possible (ex. Dice tossing) . The probability of all the outcomes can be represented in an histogram

Ordinary histogram



Cumulative histogram



# Distributions

Discrete distribution

- Probabilities are described with distributions
  - Discrete variables: data are valued on a discrete set:  $x \in I = \{x_i, i \in \mathbb{N}\}$ 
    - Probability distribution:  $f(x_i)$  = probability for the value  $x_i$
    - Cumulative distribution:  $F(x_i)$  = probability for a value  $\leq x_i$
    - Normalization:  $\sum_{i \in I} f(x_i) = 1$

$$F(x_{HIGHEST}) = 1$$

# Examples

4

Consider a coin with a tail probability equal to  $p$

Flip a coin three times.

*Plot the probability distribution for the number of tails*

you have to find the probability of getting any one of the 4 possible results

$$P(t_0) = (1-p)^3$$

$$p(t_1) = 3P(1-p)^2$$

$$p(t_2) = 3P^2(1-p)$$

$$P(t_3) = P^3$$

Flip a coin until the first tail appears.

*Plot the probability distribution for the number of flips*

$$P(1) = p$$

$$P(2) = (1-p)p$$

$$p(3) = (1-p)^2P$$

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

the process stops when we get the first tail

# Empirical Distributions

- Experimental data are represented with normalized histograms
  - ▣ Discrete variables: data are valued on a discrete set:  $x \in I = \{x_i, i \in N\}$ 
    - Probability distribution:  $f(x_i) = N(x=x_i)/N_{\text{tot}}$
    - Cumulative distribution:  $F(x_i) = N(x \leq x_i)/N_{\text{tot}}$
    - Normalization: 
$$\sum_{i \in I} f(x_i) = 1$$

# Characterization: Mode

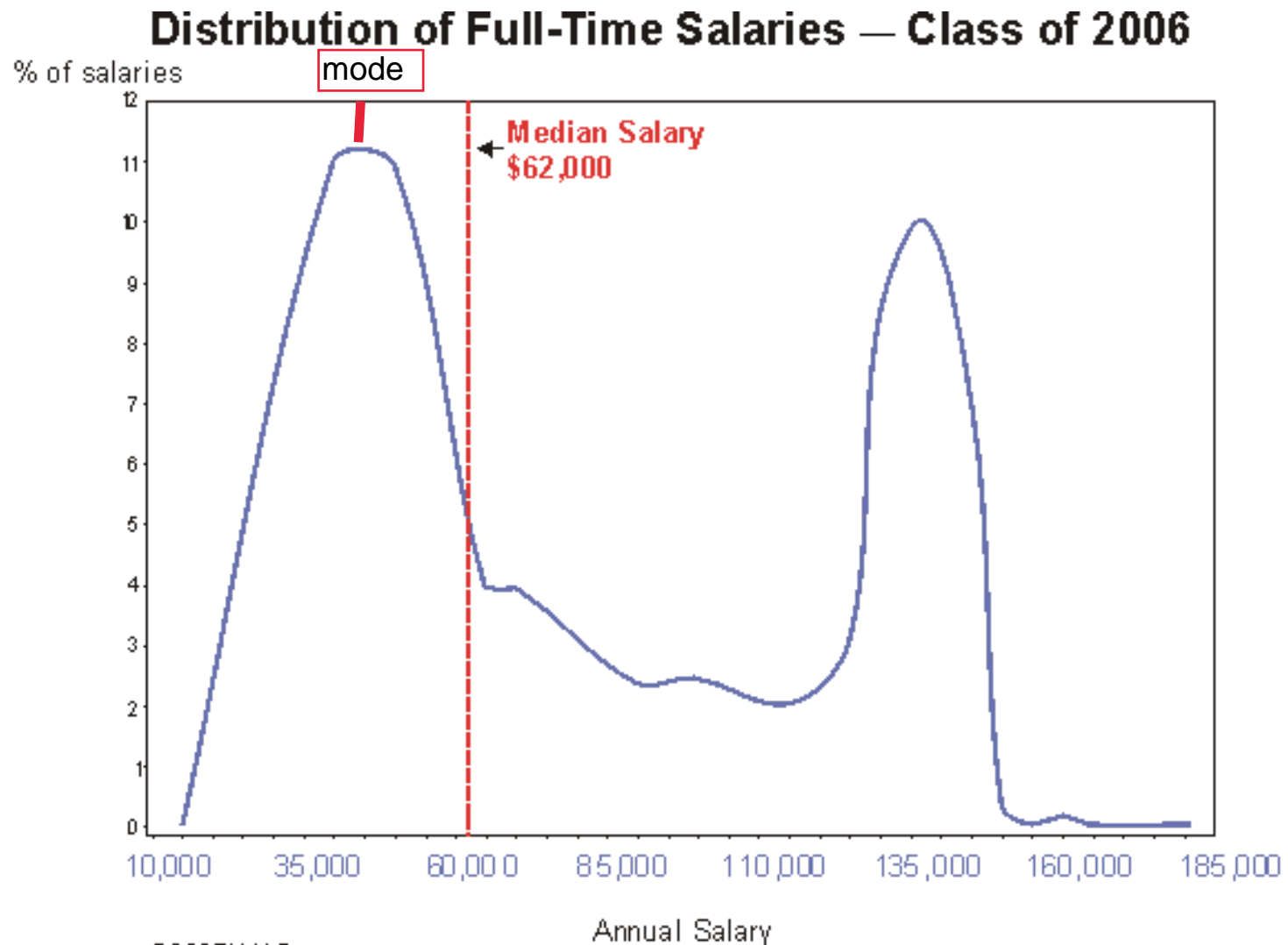
- The mode is the value that occurs the most frequently in a probability distribution.
- For empirical distributions, the mode is the value that most frequently occurs in a data set.
- More than one mode can be present

# Characterization: Median

split the probability in two half so 50% to get a result on either side

- a **median** is described as the numeric value separating the higher half of a sample, a population, or a probability distribution, from the lower half. The *median* of a finite list of numbers can be found by arranging all the observations from lowest value to highest value and picking the middle one. If there is an even number of observations, then there is no single middle value; the median is then usually defined to be the mean of the two middle values.

# Characterization: Median





## Averages (The Median)

The median is the middle value of a set of data once the data has been **ordered**.

**Example 1.** Robert hit 11 balls at Grimsby driving range. The recorded distances of his drives, measured in yards, are given below. Find the median distance for his drives.

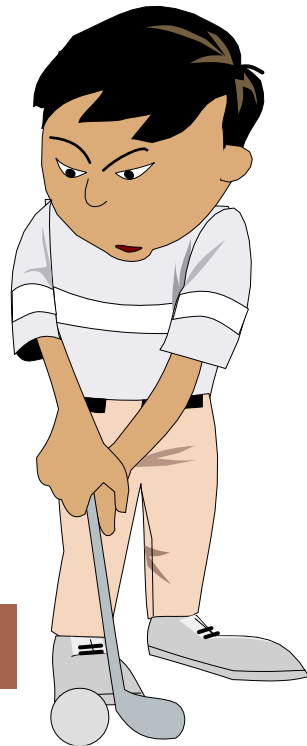
85, 125, 130, 65, 100, 70, 75, 50, 140, 95, 70

50, 65, 70, 70, 75, 85, 95, 100, 125, 130, 140

Single middle value

Ordered data

Median drive = 85 yards



# Averages (The Median)

The median Algorithm.

1 Sort the data

2 If the number of data is odd:

Take the middle value

Else:

Take the average between the two  
central values

# Measuring the Spread with Median

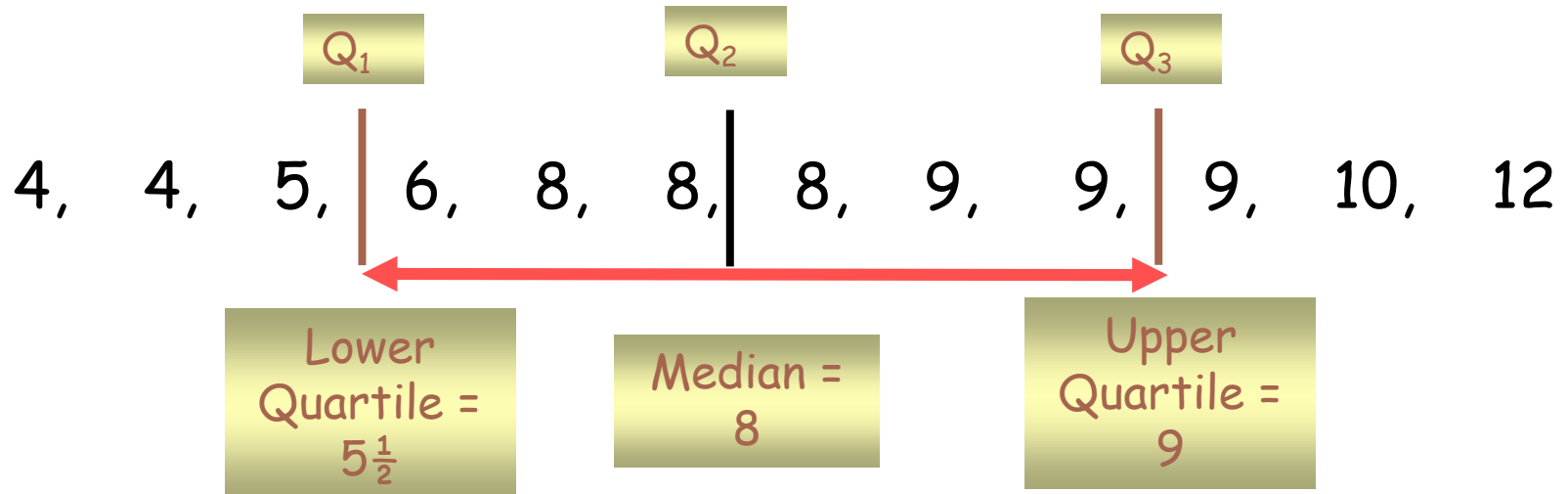
## Finding the median, quartiles and inter-quartile range.

quartiles=median of the medians

**Example 1:** Find the median and quartiles for the data below.

12, 6, 4, 9, 8, 4, 9, 8, 5, 9, 8, 10

Order the data



$$\text{Inter-Quartile Range} = 9 - 5\frac{1}{2} = 3\frac{1}{2}$$

# Expected Value

For a given function  $g(X)$  defined over discrete random variables with probability distribution  $p(X)$  the expected value of a discrete variable is computed as:

$$E[g(X)] = \sum_{i \in I} g(x_i) p(x_i)$$

# Variance and Moments of a Random Variable

- The 1th moment of a random variable  $X$  is  $E[X]$

- **Definition**

- The  $k$  th moment of a random variable  $X$  is  $E[X^k]$

- **Definition**

- The Variance of a random variable  $X$  is defined as

$$\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - (E[X])^2.$$

mean of the squares minus the square of the mean

- The standard deviation of a random variable  $X$  is

$$\sigma[X] = \sqrt{\text{Var}[X]}.$$

# Characterization: Mean

If  $X$  is a random variable with probability distribution  $p(X)$

each value \* the number of times you got it summed together divided by the total number of data

□ Mean

$$\mu = E[X] = \sum_{i \in I} x_i p(x_i)$$

E=expected value(weighted mean)

□ For empirical distributions

$$\mu = E[X] = \sum_{i \in I} x_i \frac{N(x_i)}{N_{TOT}} = \frac{\sum_{j \in Data} x_j}{N_{TOT}}$$

# Mean and Variance

## □ Variance

$$\begin{aligned} \text{Var}(X) &= E[(X - \mu)^2] = E[X^2 - 2X\mu + \mu^2] = \\ &= E[X^2] - 2\mu \cdot E[X] + \mu^2 = E[X^2] - \mu^2 \end{aligned}$$



# Properties of Expected Values

If  $X$  and  $Y$  are random variables,

- 1)  $E[X+Y]=E[X]+E[Y]$
- 2)  $E[X+c]=E[X]+c$
- 3)  $E[aX+bY]=aE[X]+bE[Y]$

For example, if  $E[X]=5$  and  $E[Y]=6$ , then

$$E[X+5]=5+5=10$$

$$E[2X+5]=2*5+5=15$$

$$E[3X+2Y]=3*5+2*6=27$$

Prove 1), 2) and 3)

# Variance and Moments of a Random Variable

□ Proof 1:

$$\begin{aligned} E[X + Y] &= \sum_{i,j} (x_i + y_j) p(x_i, y_j) = \\ &= \sum_i x_i \sum_j p(x_i, y_j) + \sum_j y_j \sum_i p(x_i, y_j) = \\ &= \sum_i x_i p(x_i) + \sum_j y_j p(y_j) = E[X] + E[Y] \end{aligned}$$

# Variance and Moments of a Random Variable

## □ Theorem

- ▣ If  $X$  and  $Y$  are two independent random variables, then

$$E[X \cdot Y] = E[X] \cdot E[Y].$$

# Variance and Moments of a Random Variable

□ Proof:

$$E[X \cdot Y] = \sum_{i,j} x_i y_j p(x_i, y_j) =$$

**Only if the events are independent**

$$= \sum_{i,j} x_i y_j p(x_i) p(y_j) =$$

$$= \sum_i x_i p(x_i) \sum_j y_j p(y_j) = E[X] \cdot E[Y]$$

# Properties of Expected Values and Variances

NB

1.  $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$ , if  $X$  and  $Y$  are independent.
2.  $\text{Var}(aX) = (a^2)\text{Var}(X)$
3.  $\text{Var}(X+c) = \text{Var}(X)$
4.  $\text{Var}(aX+bY) = (a^2)\text{Var}(X) + (b^2)\text{Var}(Y)$ , if  $X$  and  $Y$  are independent.

# Variance and Moments of a Random Variable

□ **Definition** is the variance for two variable

▣ The covariance of two random variable  $X$  and  $Y$  is

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

$$\begin{aligned} \text{Cov}[X, Y] &= E[(X - E[X])(Y - E[Y])] = \\ &= E[XY] - E[X]E[Y] - E[Y]E[X] + E[X]E[Y] = \\ &= E[XY] - E[X]E[Y] \end{aligned}$$

# Variance and Moments of a Random Variable

## □ Definition

- ▣ The covariance of two random variable  $X$  and  $Y$  is

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

## □ Theorem

- ▣ For any two random variables  $X$  and  $Y$ .

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}(X, Y).$$

# Variance and Moments of a Random Variable

□ Proof:

directly proof by  $\text{Var}[X] = E[(X - E[X])^2]$

$$\begin{aligned}\text{Var}[X + Y] &= E[(X + Y - E[X + Y])^2] \\&= E[(X + Y - E[X] - E[Y])^2] \\&= E[(X - E[X])^2 + (Y - E[Y])^2 + 2(X - E[X])(Y - E[Y])] \\&= E[(X - E[X])^2] + E[(Y - E[Y])^2] + 2E[(X - E[X])(Y - E[Y])] \\&= \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}(X, Y)\end{aligned}$$



# Covariance and Variance

## □ Corollary

- If  $X$  and  $Y$  are independent random variables, then

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y].$$

(already proved)

→ THEN

$$\text{Cov}(X, Y) = 0$$

# Binomial distribution

Compute the probability of obtaining exactly 3 tails tossing 5 coins. These are all the possibilities

TTTTT	HTTTT	HHTTT	TTHHH	THHHH	HHHHH
	THTTT	HTHTT	THTHH	HTHHH	
	TTHTT	HTTHT	THHTH	HHTHH	
	TTTHT	HTTTT	THHHT	HHHTH	
	TTTTH	THHTT	HTTTH	HHHHT	
		THTHT	HTHTH		
		THTTH	HTHHT		
		TTHHT	HHTTH		
		TTHTH	HHTHT		
		TTTHH	HHHTT		

$$f(5)=10/32$$

And what if the probabilities of head and tail are different?

# Binomial distribution

How can we count without enumerating?

Total number of configurations: 5 events, with two possibilities each  $\rightarrow 2^5=32$

Number of configurations with two heads:

- 1) choose 3 of the 5 positions:  $5 \times (5-1) \times (5-2)$  possibilities
  - 2) the order of choice is irrelevant: there are 6 possible permutations of 3 objects ( $3!$ )
- $\rightarrow 5 \times 4 \times 3 / 6 = 5! / (3! \times 2!) = 10$

And what if the probabilities of head and tail are different?

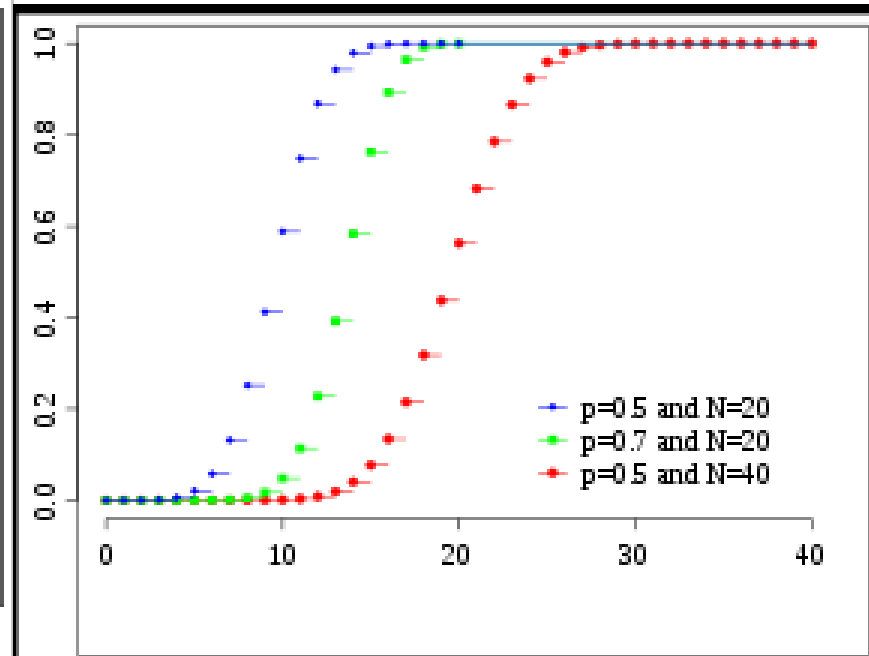
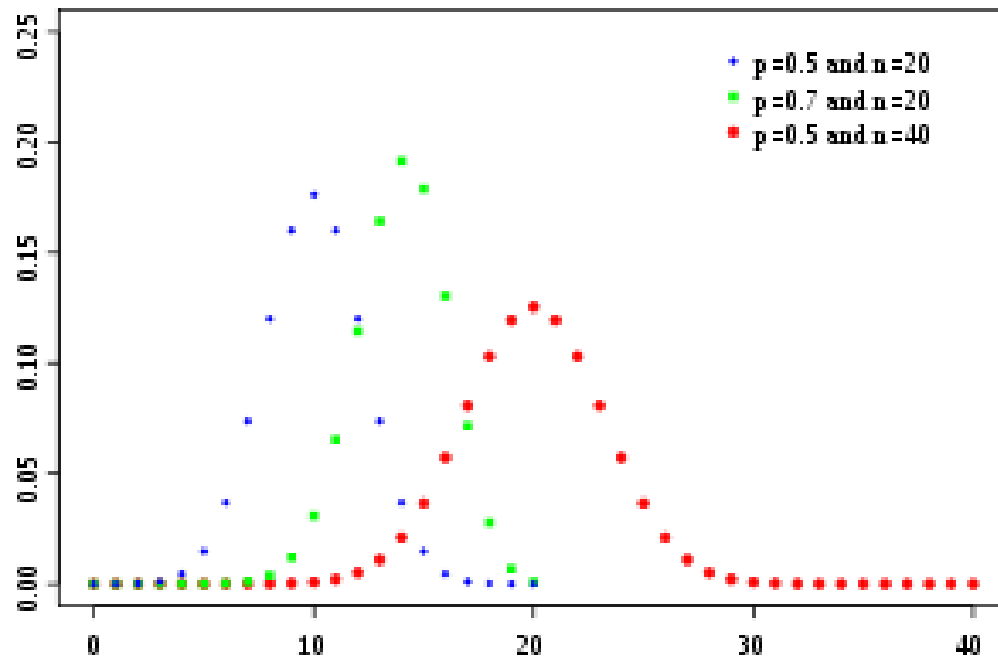
# Binomial distribution

Probability of having  $k$  positive events out of  $n$  independent experiments.  
The probability of a positive event is  $p$

$P$ =probability of single positive event

$$f(k; n, p) = \Pr(K = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

Binomial coefficient= $\frac{n!}{k!(n-k)!}$



# Binomial distribution (mean)

$$\mu = E[k] = \sum_{k \in N} kf(k) = \text{NP}$$

$$= \sum_{k=0}^n k \cdot \binom{n}{k} p^k (1-p)^{n-k} = \sum_{k=1}^n k \cdot \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k} =$$

$$= \sum_{k=1}^n \frac{(n-1)!n}{(k-1)!(n+1-k-1)!} p p^{(k-1)} (1-p)^{(n-k)} =$$

Define  $j=k-1$

$$= np \sum_{j=0}^{n-1} \frac{(n-1)!}{(j)!(n-1-j)!} p^j (1-p)^{(n-1-j)} = np$$

# Binomial distribution (variance)

---

$$\text{Var}(k) = E[(k - \mu)^2] = np(1 - p)$$

# Variance and Moments of a Random Variable

## □ Example

### ▣ Variance of a Binomial Random Variable

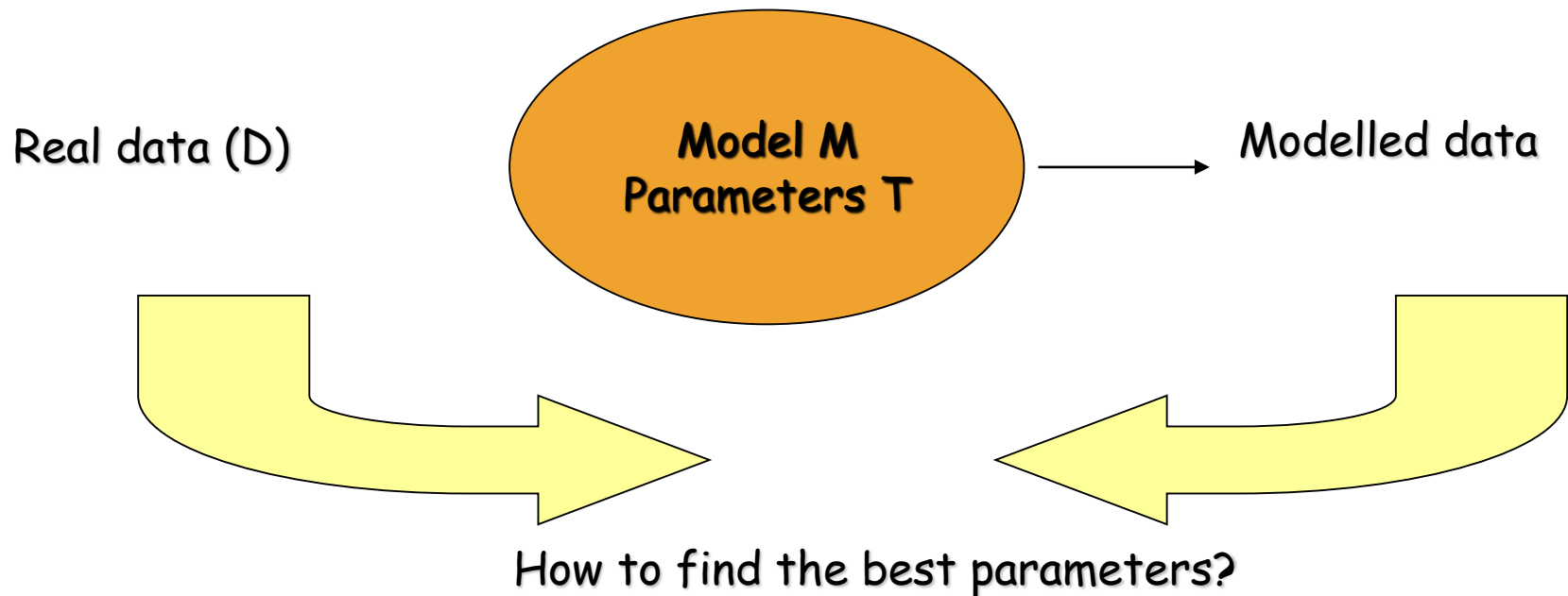
- The variance of a binomial random variable  $X$  with parameters  $n$  and  $p$  can be determined directly by computing  $E[X^2]$ .

$$E[X^2] = \sum_j^n C_j^n p^j (1-p)^{n-j} j^2 = n(n-1)p^2 + np$$

$$\begin{aligned} \text{Var}[X] &= E[X^2] - (E[X])^2 = n(n-1)p^2 + np - (np)^2 \\ &= np(1-p) \end{aligned}$$

# Evaluation of parameters from empirical data

Generally speaking, a parametric model  $M$  aims to reproduce a set of known data



We search for a  $P$  that maximizes the similarity with our data



# Maximum likelihood

$$T^* = \underset{T}{\operatorname{argmax}} P(D|T, M) =$$

$$T^* = \underset{T}{\operatorname{argmax}} \log(P(D|T, M))$$

*D=data, M= model, T=model parameters*

# Evaluation of $p$

- Given  $n$  independent experiments, with  $k$  positive events, how to evaluate the best  $p$ ?
- MAXIMUM LIKELIHOOD definition:
  - ▣ Best  $p \rightarrow p$  that maximises the total probability of the experiment results
- Example
  - ▣ Which is the  $p$  that maximises the outcome of 10010001100111100010 ?

# Evaluation of p

$$f(k | n, p) = \binom{n}{k} p^k (1-p)^{(n-k)}$$

evaluating the best P to get our data

$$\left. \frac{df}{dp} \right|_{p=p_{ML}} = 0$$

maximum of  
likelihood  
because the  
derivate=0

$$\Rightarrow k p_{ML}^{k-1} (1-p_{ML})^{n-k} - p_{ML}^k (n-k) (1-p_{ML})^{(n-k-1)} = 0$$

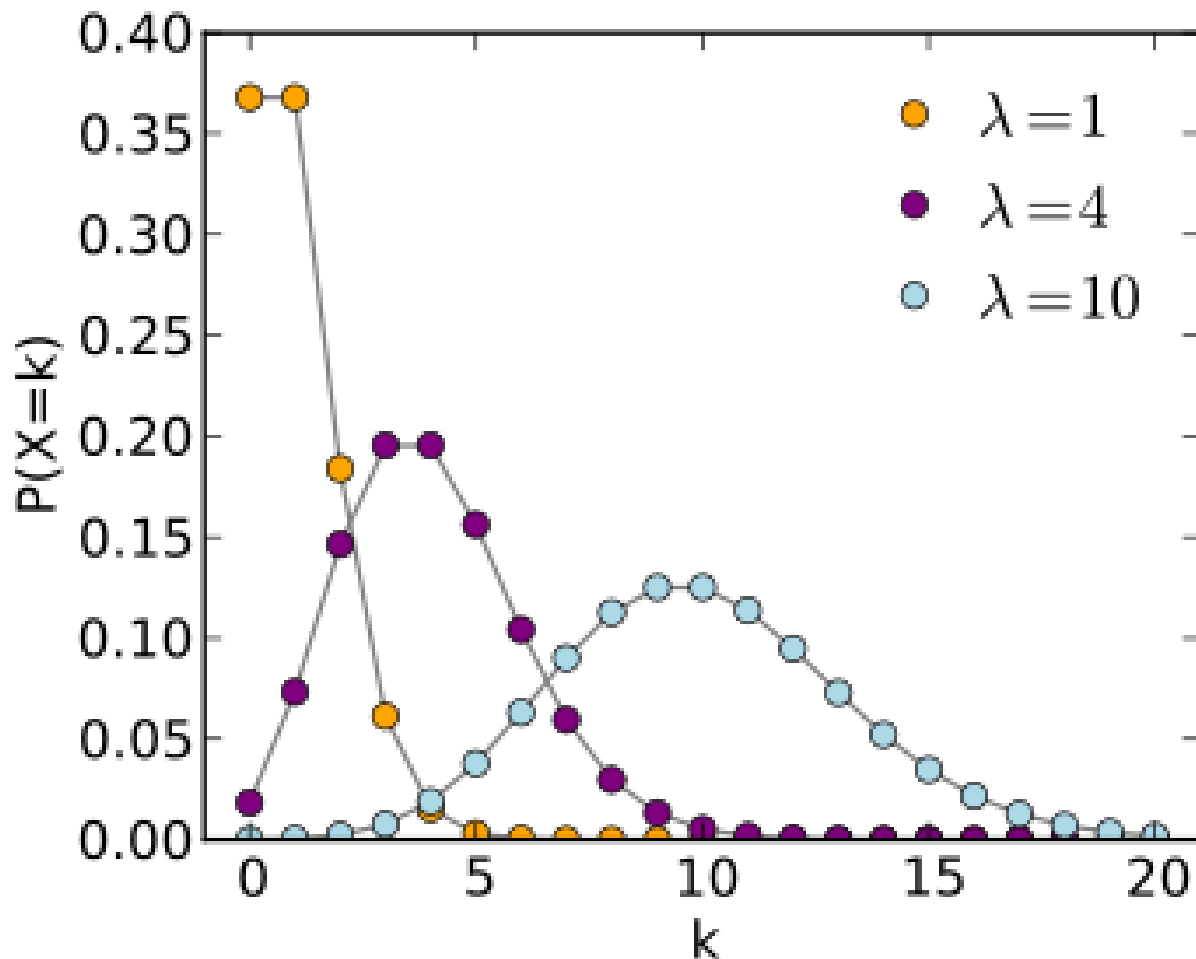
$$k(1-p_{ML}) - p_{ML}(n-k) = 0$$

$$p_{ML} = \frac{k}{n}$$

# Poisson distribution

the **Poisson distribution** is a discrete probability distribution that expresses the probability of a number of events occurring in a fixed period of time if these events occur with a known average rate ( $\lambda$ ) and independently of the time since the last event. (The Poisson distribution can also be used for the number of events in other specified intervals such as distance, area or volume.)

# Poisson distribution



$$f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

$$E[k] = \lambda$$

$$Var[k] = \lambda$$

probability to have (Xaxis) event in a unit of time given the Lambda frequency of it

# Continuous Distributions

- Data are described with distributions
  - ▣ Continuous variables: data are valued on a continuous interval:  $x \in I \subseteq \mathbb{R}$ 
    - Cumulative probability:  $F(x) = N(x \leq X)/N_{\text{tot}}$
    - Prob( $A < x \leq B$ ) = <sup>cumulative</sup>  $F(B) - F(A) = \int_A^B f(x) dx$
    - Probability density function (limit)  $f(x) = \frac{dF(x)}{dx}$
    - Normalization:  $\int_I f(x) dx = 1$

# Mean and Variance

## □ Mean

$$\mu = E[x] = \int_{x \in I} x \cdot f(x) dx$$

## □ Variance

$$\text{Var}[x] = E[(x - \mu)^2] = \int_{x \in I} (x - \mu)^2 \cdot f(x) dx$$

# The Uniform Probability Distribution

## □ Uniform Probability Density Function

$$f(x) = 1/(b - a) \text{ for } a \leq x \leq b \\ = 0 \text{ elsewhere}$$

where

$a$  = smallest value the variable can assume

$b$  = largest value the variable can assume

The probability of the continuous random variable assuming a specific value is 0.

$$P(x=x_1) = 0$$



# Example: Buffet

Customers are charged for the amount of salad they take. Sampling suggests that the amount of salad taken is uniformly distributed between 5 ounces and 15 ounces.

## □ Probability Density Function

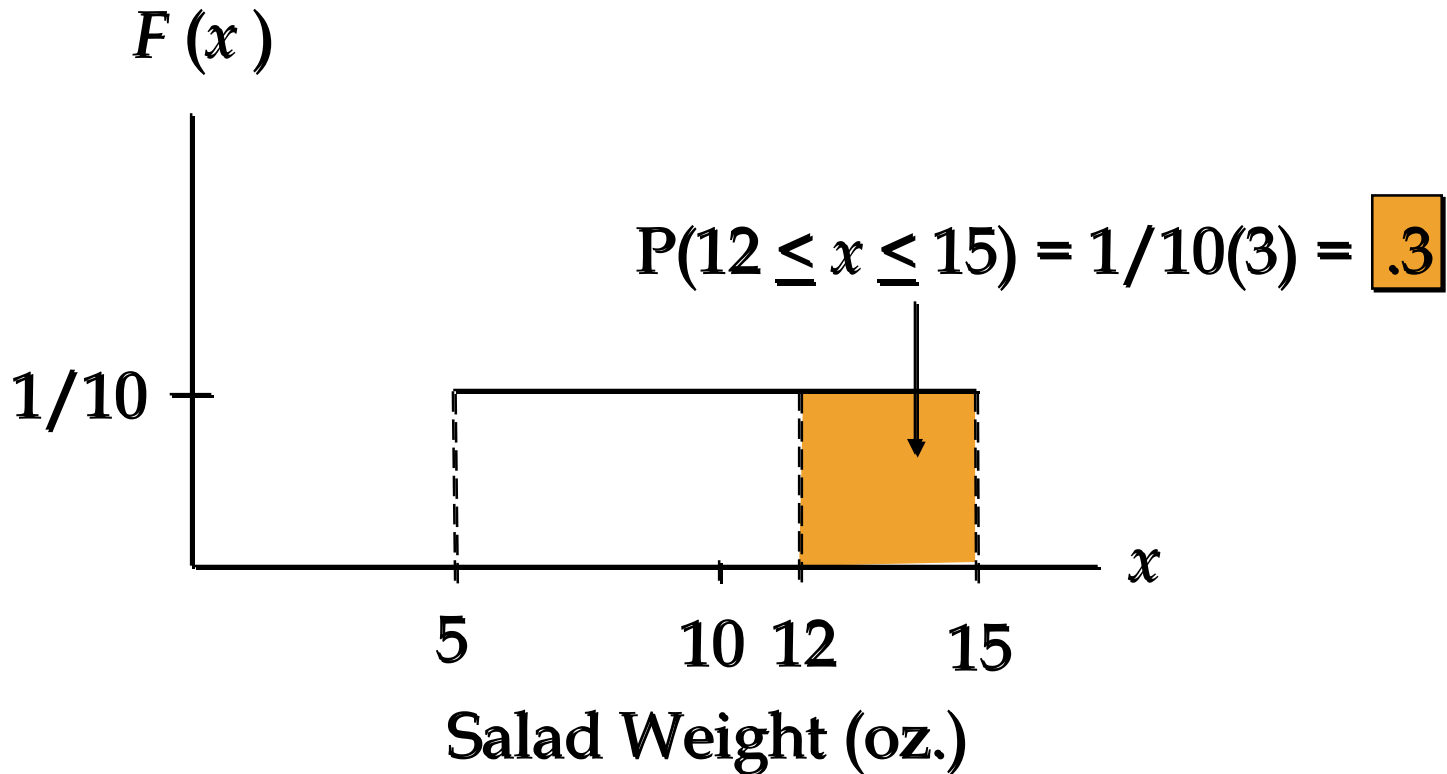
$$f(x) = 1/10 \text{ for } 5 \leq x \leq 15$$
$$= 0 \text{ elsewhere}$$

where

$x$  = salad plate filling weight

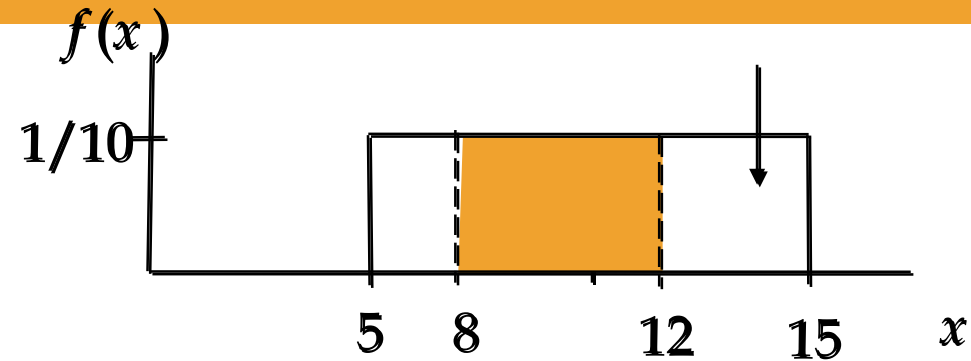
# Example: Buffet

What is the probability that a customer will take between 12 and 15 ounces of salad?



# The Uniform Probability Distribution

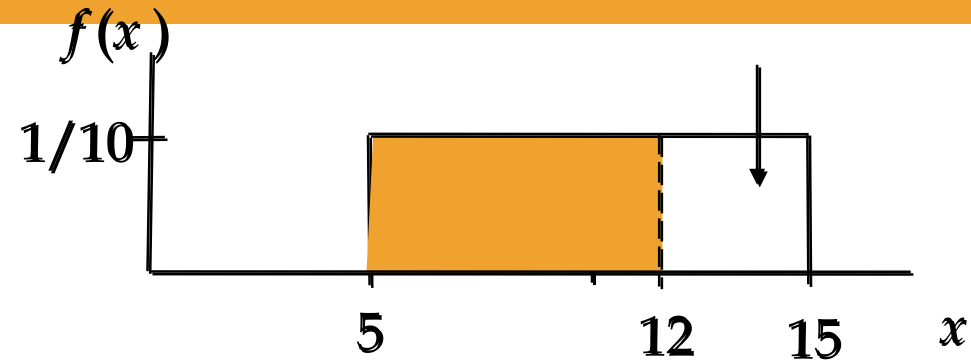
$$P(8 \leq x \leq 12) = ?$$



$$P(8 \leq x \leq 12) = (1/10)(12-8) = .4$$

# The Uniform Probability Distribution

$$P(0 \leq x \leq 12) = ?$$



$$\begin{aligned} P(0 \leq x \leq 12) &= P(5 \leq x \leq 12) = \\ &= (1/10)(12-5) = .7 \end{aligned}$$

# The Uniform Probability Distribution

## □ Uniform Probability Density Function

$$f(x) = 1/(b - a) \text{ for } a \leq x \leq b \\ = 0 \text{ elsewhere}$$

## □ Expected Value of $x$

$$E(x) = (a + b)/2$$

## □ Variance of $x$

$$\text{Var}(x) = (b - a)^2/12$$

where

$a$  = smallest value the variable can assume

$b$  = largest value the variable can assume

Prove it!!

# Normal distribution

The normal distribution is considered the most “basic” continuous probability distribution. Specifically, by the central limit theorem, under certain conditions, the sum of a number of random variables with finite means and variances approaches a normal distribution as the number of variables increases.

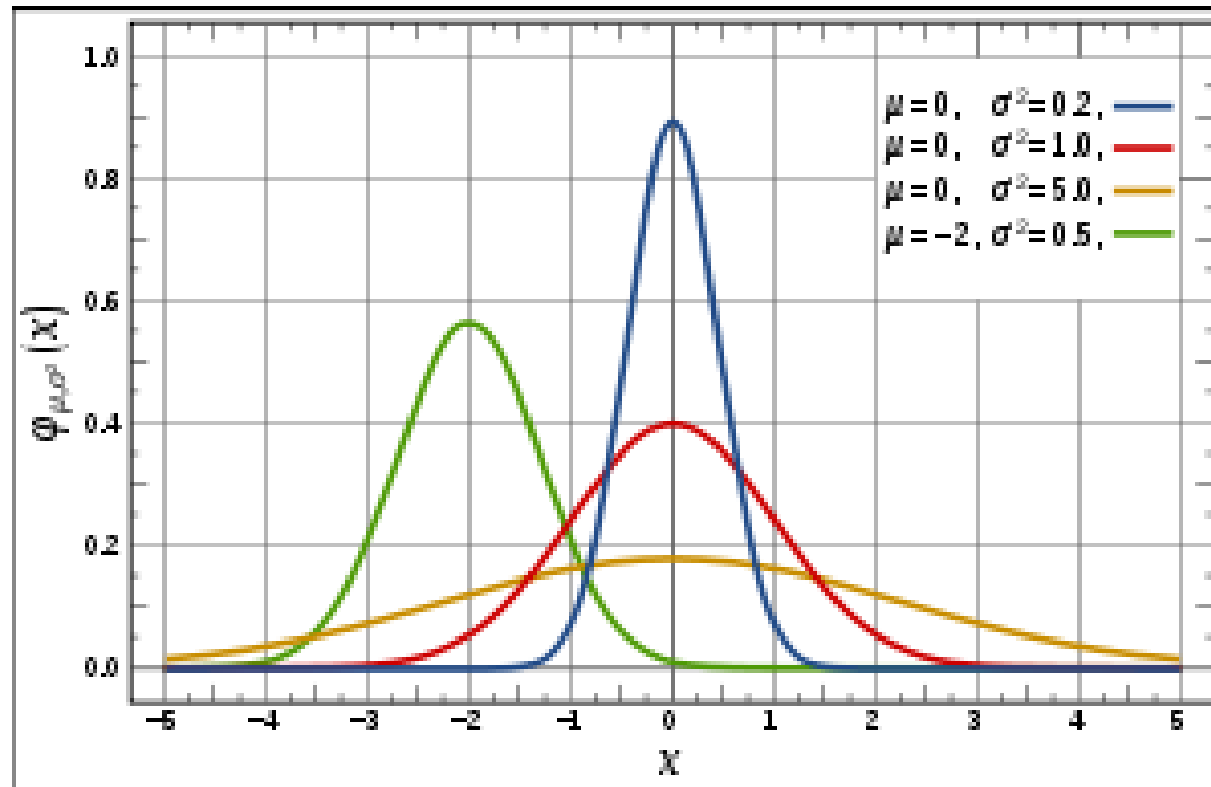
For this reason, the normal distribution is commonly encountered in practice, and is used throughout statistics, natural sciences, and social sciences as a simple model for complex phenomena. For example, the observational error in an experiment is usually assumed to follow a normal distribution, and the propagation of uncertainty is computed using this assumption.

# Normal distribution

$$f(x) = N(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \cdot e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

$\mu$  : mean

$\sigma^2$  : variance



# Normal distribution: continuous

$$f(x) = N(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \cdot e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

normalize

$\mu$  : mean

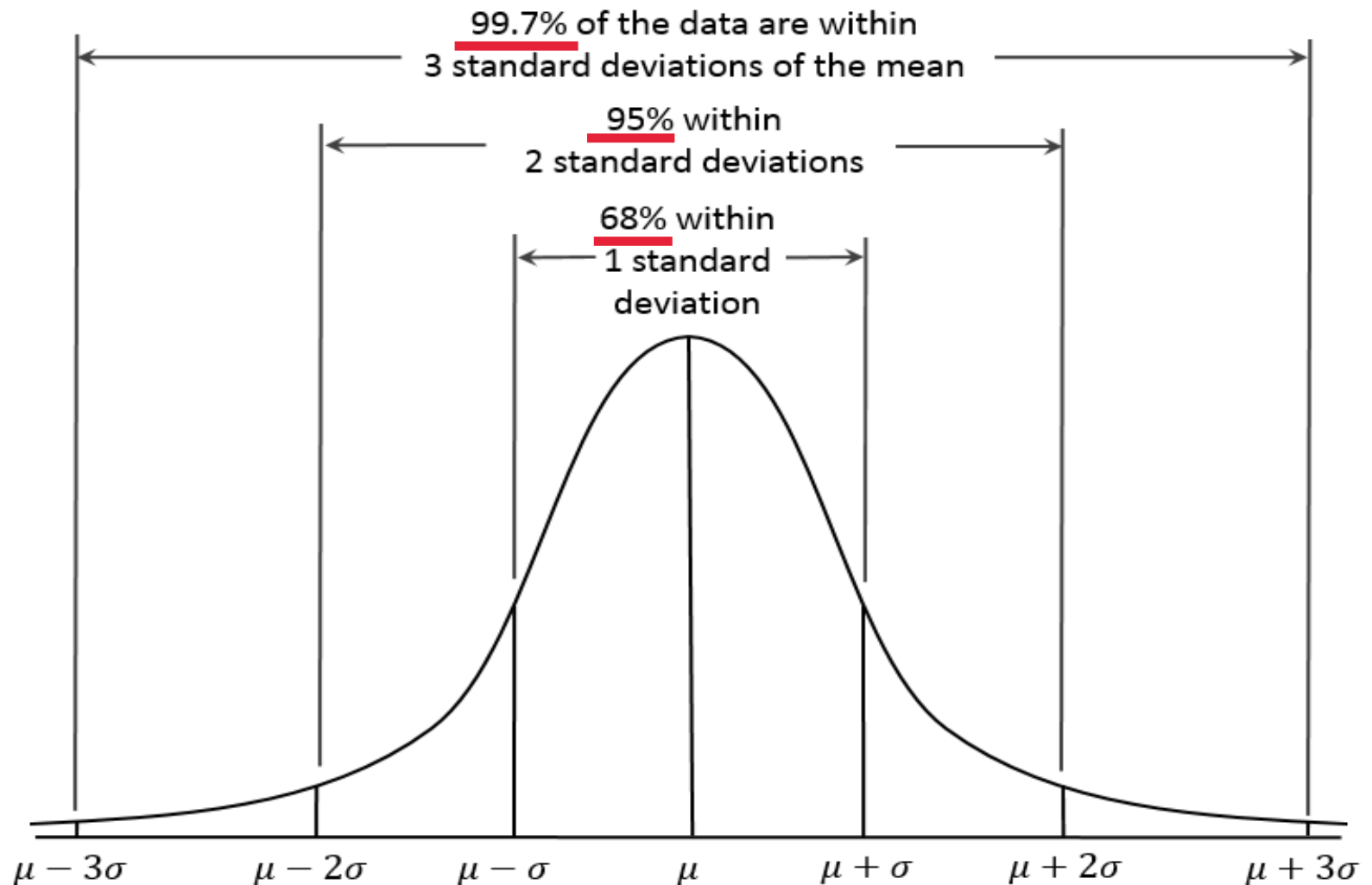
$\sigma^2$  : variance

$$E[x] = \int x N(x | \mu, \sigma^2) dx = \mu$$

$$E[x^2] = \int x^2 N(x | \mu, \sigma^2) dx = \mu^2 + \sigma^2$$



# Normal distribution



# Normal distribution in standard form

Let  $x$  be a Normally distributed random variable with mean  $\mu$  and variance  $\sigma^2$

Introduce the variable

$$z = (x - \mu) / \sigma$$

$$Z = (x - \text{mean}) / \sqrt{\text{variance}}$$

$$E[z] = E\left[\frac{x - \mu}{\sigma}\right] = \frac{E[x] - \mu}{\sigma} = 0$$

$$E[z^2] = E\left[\frac{(x - \mu)^2}{\sigma^2}\right] = \frac{E[x^2] - \mu^2}{\sigma^2} = 1$$

$$\text{Var}[z] = E[z^2] - E[z]^2$$

**$z$  is a Normally distributed random variable with mean 0 and variance 1**

Entries in the table give the area under the curve between the mean and  $z$  standard deviations above the mean. For example, for  $z = 1.25$  the area under the curve between the mean (0) and  $z$  is 0.3944.

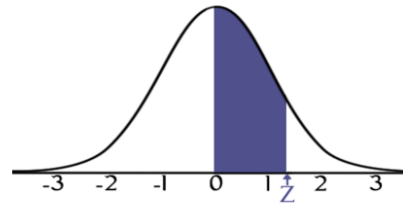
Heights of females are normally distributed with a mean of 63.6 inches and a standard deviation of 2.5 inches. What is the probability that a female selected at random will be taller than 70 inches ?

$$Z = \frac{(x - 63.6)}{2.5} \quad X = 70$$

$$Z = 2.56$$

look on the table to see the area under a 2.56 Z

$A = 0.4948$  it is the chance to be between 0 and 63.6 and so  $0.5 - 0.4948$  is the chance to be higher than 70

$$z = (x - 63.6) / 2.5$$


Entries in the table give the area under the curve between the mean and  $z$  standard deviations above the mean. For example, for  $z = 1.25$  the area under the curve between the mean (0) and  $z$  is 0.3944.

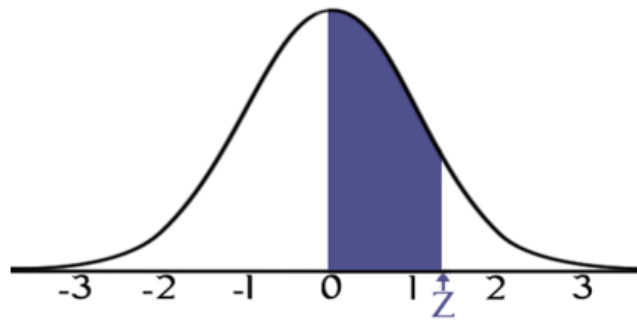
For  $x = 70 \rightarrow z = 2.56$

Prob [0-2.56] = 0.4948

$$\text{Prob } [z > 2.56] = 0.5 - 0.4948 = 0.0052$$
[illegible]

- The SAT scores are normally distributed with a mean of 500 and a standard deviation of 150 points. If an admissions office wanted to accept students who were in the top 40%, what would the cut-off scores be?

- Top 40%: given the standardized variable
  - ▣ The cutoff is the value of the standard table giving an area equal to 0.1 (between 0 and the cutoff)



### STANDARD NORMAL TABLE (Z)

Entries in the table give the area under the curve between the mean and z standard deviations above the mean. For example, for  $z = 1.25$  the area under the curve between the mean (0) and z is 0.3944.

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.0000	0.0040	0.0080	0.0120	0.0160	0.0190	0.0239	0.0279	0.0319	0.0359
0.1	0.0398	0.0438	0.0478	0.0517	0.0557	0.0596	0.0636	0.0675	0.0714	0.0753
0.2	0.0793	0.0832	0.0871	0.0910	0.0948	0.0987	0.1026	0.1064	0.1103	0.1141
0.3	0.1179	0.1217	0.1255	0.1293	0.1331	0.1368	0.1406	0.1443	0.1480	0.1517
0.4	0.1554	0.1591	0.1628	0.1664	0.1700	0.1736	0.1772	0.1808	0.1844	0.1879
0.5	0.1915	0.1950	0.1985	0.2019	0.2054	0.2088	0.2123	0.2157	0.2190	0.2224

$$0.26 = \text{cutoff}(z) = [\text{cutoff}(x) - \mu] / \sigma$$

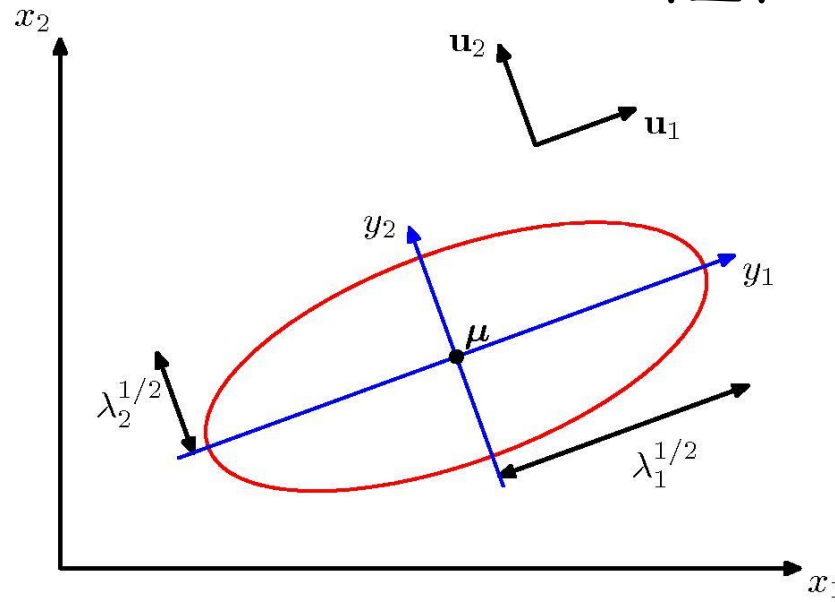
$$\text{cutoff}(x) = 150 * 0.26 + 500 = 539$$

# Gaussian distribution in a D-dimensional space

56

$$N(x | \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

where  $\mu$  is a D-valued vector (means) and  $\Sigma$  is a DxD symmetric matrix (covariance matrix), with determinant  $|\Sigma|$





# Estimating Mean and Variance from sampled data

- Consider a set of  $n$  independent and identically distributed data, following a normal distribution with mean  $\mu$  and variance  $\sigma^2$

$$x_1, x_2 \dots, x_n = X$$

How can we estimate mean and variance?

Example:

$X=0.2;0.18;0.24;0.14;0.05;0.26;0.15;0.29,0.49$

# Estimating mean and Variance from sampled data

## MAXIMUM LIKELIHOOD:

Estimating the parameters as to maximize the probability for the sampled values

$$(\mu_{ML}, \sigma_{ML}^2) = \arg \max_{\mu, \sigma^2} f(x_1, x_2 \dots x_n \mid \mu, \sigma^2)$$

# ML Mean and Variance

- The probability for the sampled data is

$$f(X \mid \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \cdot e^{\frac{-(x_n - \mu)^2}{2\sigma^2}}$$

- It has to be maximised over the variables  $\mu$  and  $\sigma^2$

# ML Mean and Variance

- Since logarithm is a monotonic function

$$\begin{aligned}\arg \max_{\mu, \sigma} [f(X|\mu, \sigma^2)] &= \arg \max_{\mu, \sigma} [\ln f(X|\mu, \sigma^2)] = \\ &= \arg \max_{\mu, \sigma} \left[ -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right]\end{aligned}$$

# ML Mean and Variance

$$\begin{cases} \frac{\partial}{\partial \mu} \left[ -\frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right] \Big|_{\mu_{ML}, \sigma_{ML}^2} = 0 \\ \frac{\partial}{\partial \sigma^2} \left[ -\frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right] \Big|_{\mu_{ML}, \sigma_{ML}^2} = 0 \end{cases}$$

$$\begin{cases} \sum_{i=1}^n (x_i - \mu_{ML}) = 0 \\ -\frac{n}{2\sigma_{ML}^2} + \frac{1}{2(\sigma_{ML}^2)^2} \sum_{i=1}^n (x_i - \mu_{ML})^2 = 0 \end{cases}$$

# ML Mean and Variance

$$\left\{ \begin{array}{l} \mu_{ML} = \frac{1}{n} \sum_{i=1}^n x_i \\ \sigma^2_{ML} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{ML})^2 \end{array} \right.$$

Are these the expressions that you usually apply?

# The $\mu_{ML}$ is an unbiased estimation

- Consider a set of data generated starting from a normal distribution  $N(\mu, \sigma^2)$  the estimation

$$\mu_{ML} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{Is an unbiased estimate of } \mu \text{ since}$$

$$E[\mu_{ML}] =$$

$$\iint \frac{1}{n} \left( \sum_{i=1}^n x_i \right) \cdot \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \cdot e^{-\frac{(x_1-\mu)^2}{2\sigma^2}} \cdot e^{-\frac{(x_2-\mu)^2}{2\sigma^2}} \cdots e^{-\frac{(x_n-\mu)^2}{2\sigma^2}} dx_1 dx_2 \cdots dx_n =$$

$$= \frac{1}{n} \sum_{i=1}^n E[x_i] = \mu$$

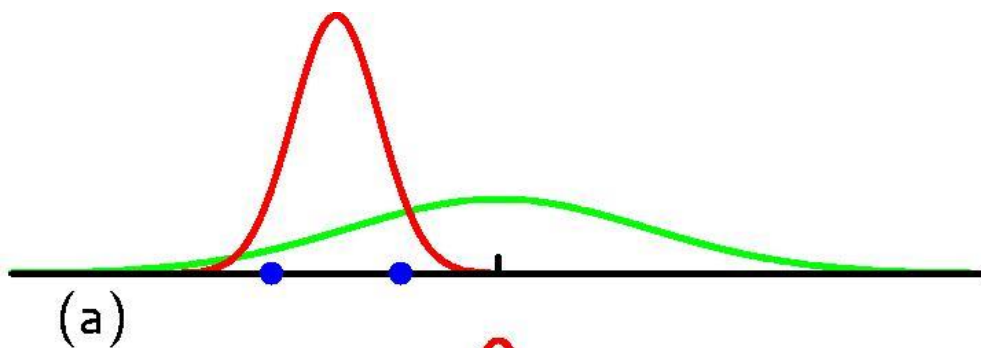
# The $\sigma_{ML}^2$ is a biased estimation

$$\sigma_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - M)^2 \quad \text{Is a biased estimate of } \sigma^2 \text{ since}$$

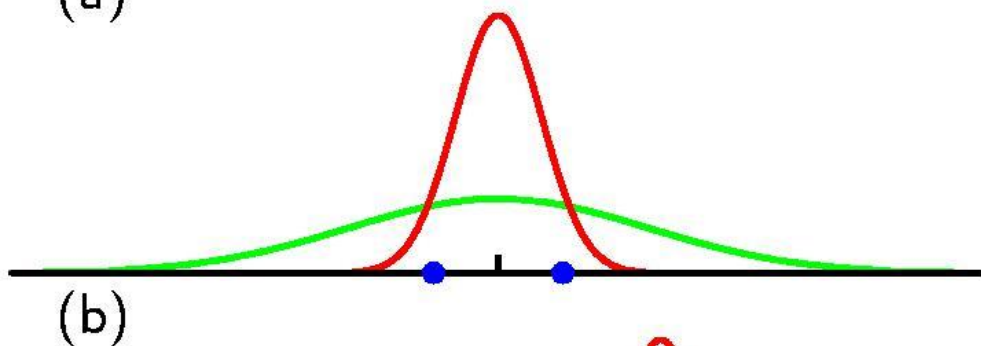
$$\begin{aligned} E[\sigma_{ML}^2] &= E\left[\frac{1}{n} \sum_{i=1}^n \left(x_i - \frac{1}{n} \sum_{j=1}^n x_j\right)^2\right] = \\ &= E\left[\frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{2}{n^2} \sum_{i=1}^n \sum_{j=1}^n x_i x_j + \frac{1}{n^2} \sum_{j=1}^n \sum_{k=1}^n x_j x_k\right] = \\ &= E[x^2] - \frac{1}{n^2} E\left[\sum_{i=1}^n \sum_{j=1}^n x_i x_j\right] = E[x^2] - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\mu^2 + \delta_{ij} \sigma^2) = \\ &= \mu^2 + \sigma^2 - \mu^2 - \frac{1}{n} \sigma^2 = \frac{n-1}{n} \sigma^2 \neq \sigma^2 \end{aligned}$$



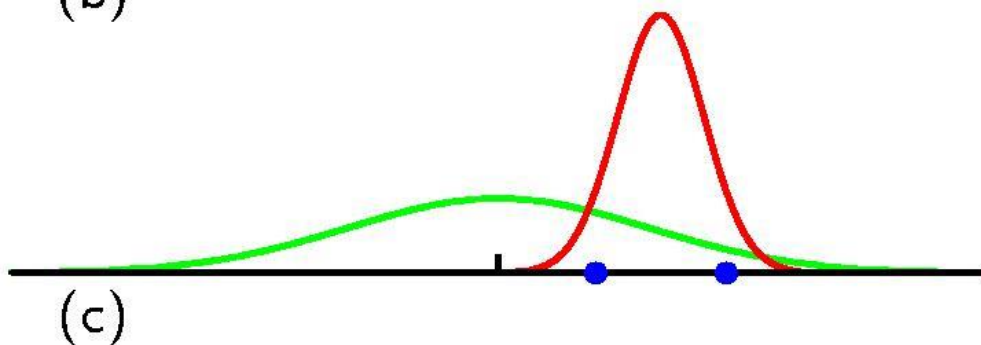
# $\sigma_{ML}^2$ is underestimating $\sigma^2$



Green: "Real" distribution  
Red: Sampled distributions



$\sigma_{ML}^2$  is evaluated with respect to the sample mean



# Unbiased estimation of $\mu$ and $\sigma^2$

$$M = \mu_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$S^2 = \frac{n}{n-1} \sigma^2_{ML} = \frac{1}{n-1} \sum_{i=1}^n (x_i - M)^2$$

# Distribution of the sample mean

- Consider different sets of sampled data  $X^k$  (e.g. different levels of expression of genes in  $n$  different populations of individuals)

A sample mean for each set can be defined

$$M^{(k)} = \frac{1}{n^{(k)}} \sum_{i=1}^{n^{(k)}} x_i^{(k)}$$

How are they distributed?

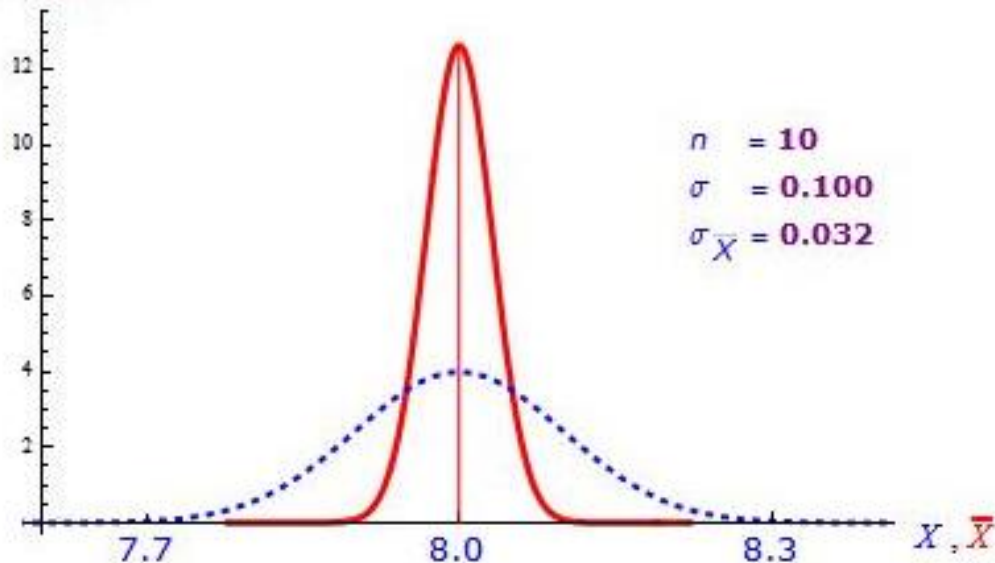
# Distribution of the sample means

$$M^{(k)} = \frac{1}{n} \sum_{i=1}^n x_i^{(k)}$$

M is a linear combination of independent identically distributed normal variables: it is normally distributed

PDFs of  $X$  and  $\bar{X}$

$f(X), g(\bar{X})$



# Distribution of the sample means

$$E[M] = \mu$$

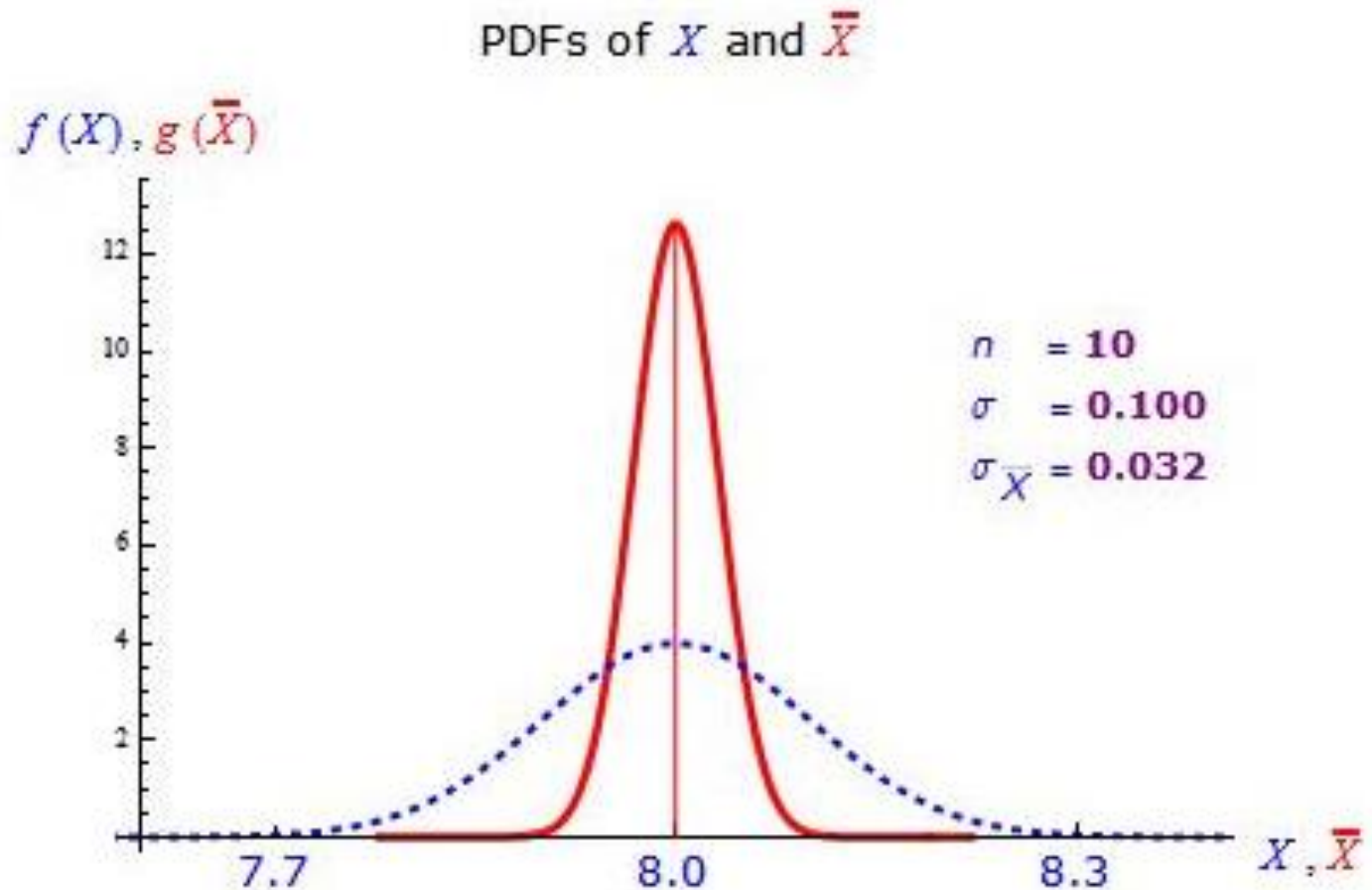
$$\begin{aligned}\text{var}[M] &= E[(M - \mu)^2] = E[M^2 + \mu^2 - 2\mu M] = \\ &= E[M^2] + \mu^2 - 2\mu^2 = E[M^2] - \mu^2\end{aligned}$$

$$E[M^2] = E\left[\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n x_i x_j\right] = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \mu^2 + \delta_{ij} \sigma^2 =$$

$$= \mu^2 + \frac{1}{n} \sigma^2$$

$$\text{var}[M] = \frac{1}{n} \sigma^2$$

# Distribution of the sample means



# Distribution of the sample deviations

- If  $X_i$  are  $n$  independent, normally distributed random variables with mean  $\mu$  and variance  $\sigma^2$ , then the random variable

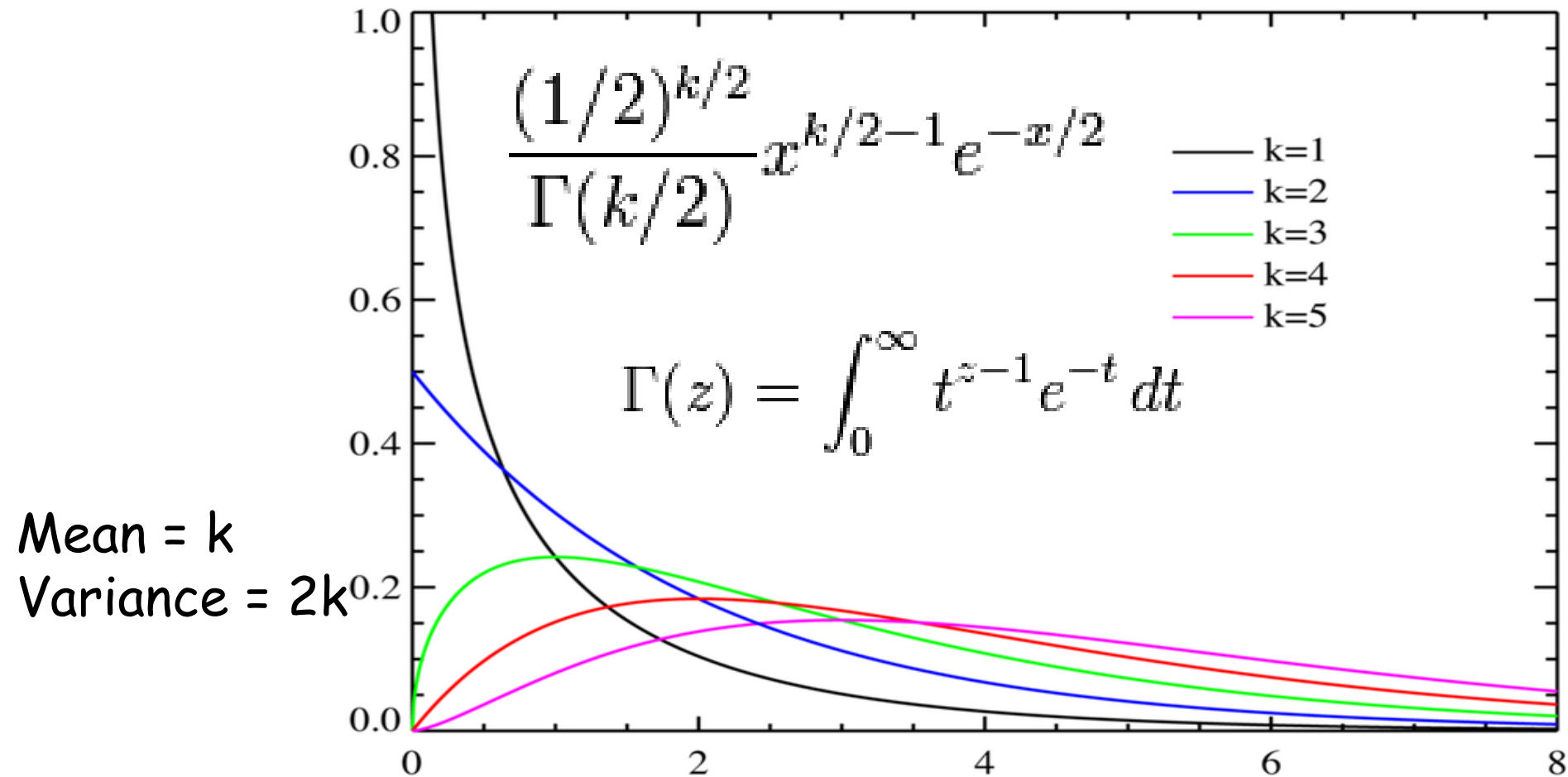
$$Q = \sum_{i=1}^N (X_i - \mu)^2$$

is distributed according to the chi-square distribution with  $n-1$  degrees of freedom. This is usually written as:

$$Q \approx \sigma^2 \chi_{n-1}^2$$

The chi-square distribution has one parameter:  $(n-1)$  that is a positive integer that specifies the number of degrees of freedom (i.e. the number of independent  $X_i - \mu$ )

# Chi-square distribution





# Distribution of the sample deviations

- So, being the unbiased sample variance

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - M)^2$$

is distributed according to the chi-square distribution with  $k$  degrees of freedom.

$$S^2 \approx \frac{\sigma^2}{n-1} \chi_{n-1}^2$$

It follows that

$$E[S^2] \approx \frac{\sigma^2}{n-1} E[\chi_{n-1}^2] = \frac{\sigma^2}{n-1} (n-1) = \sigma^2$$

Unbiased

# Distribution of the sample deviations: Proof for 1 degree of freedom

Let random variable  $Y$  be defined as  $Y = X^2$  where  $X$  has normal distribution with mean 0 and variance 1

We can compute the cumulative function

$$y < 0, P(Y < y) = 0$$

$$\begin{aligned} y \geq 0, P(Y < y) &= P(X^2 < y) = \\ &= P(|X| < \sqrt{y}) = F_x(\sqrt{y}) - F_x(-\sqrt{y}) \end{aligned}$$

# Distribution of the sample deviations: Proof for 1 degree of freedom

$y > 0$

$$f_y(y) = f_x(\sqrt{y}) \frac{\partial(\sqrt{y})}{\partial y} - f_x(-\sqrt{y}) \frac{\partial(-\sqrt{y})}{\partial y}$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{y}{2}} \frac{1}{2y^{1/2}} + \frac{1}{\sqrt{2\pi}} e^{-\frac{y}{2}} \frac{1}{2y^{1/2}}$$

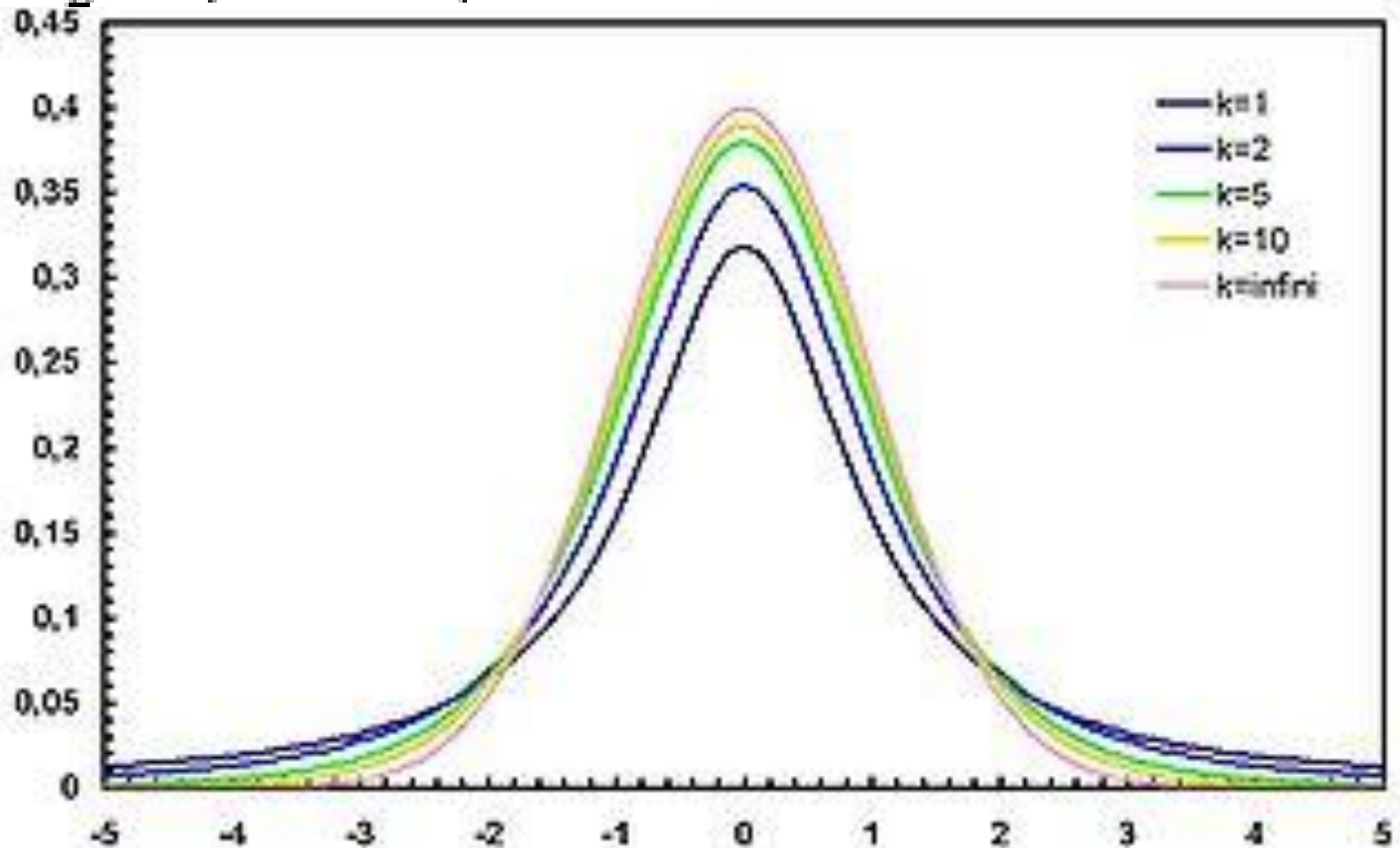
$$= \frac{1}{2^{\frac{1}{2}} \Gamma(\frac{1}{2})} y^{\frac{1}{2}-1} e^{-\frac{y}{2}}$$

$$Y = X^2 \sim \chi_1^2$$

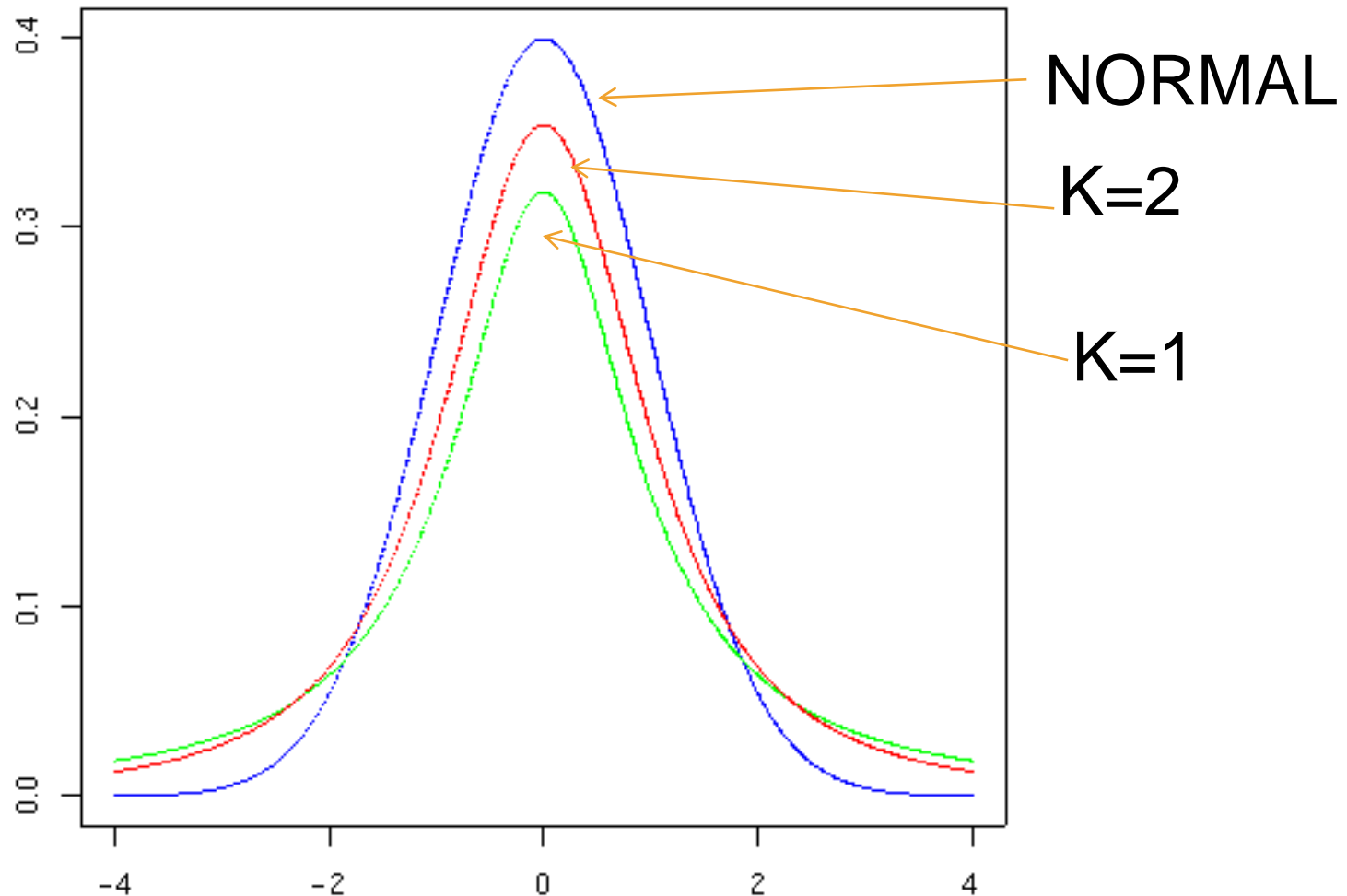
# Student's t-distribution

$$\frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi} \Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{x^2}{\nu}\right)^{-\left(\frac{\nu+1}{2}\right)}$$

$\nu$  : degrees of freedom



# Normal Vs t- distributions



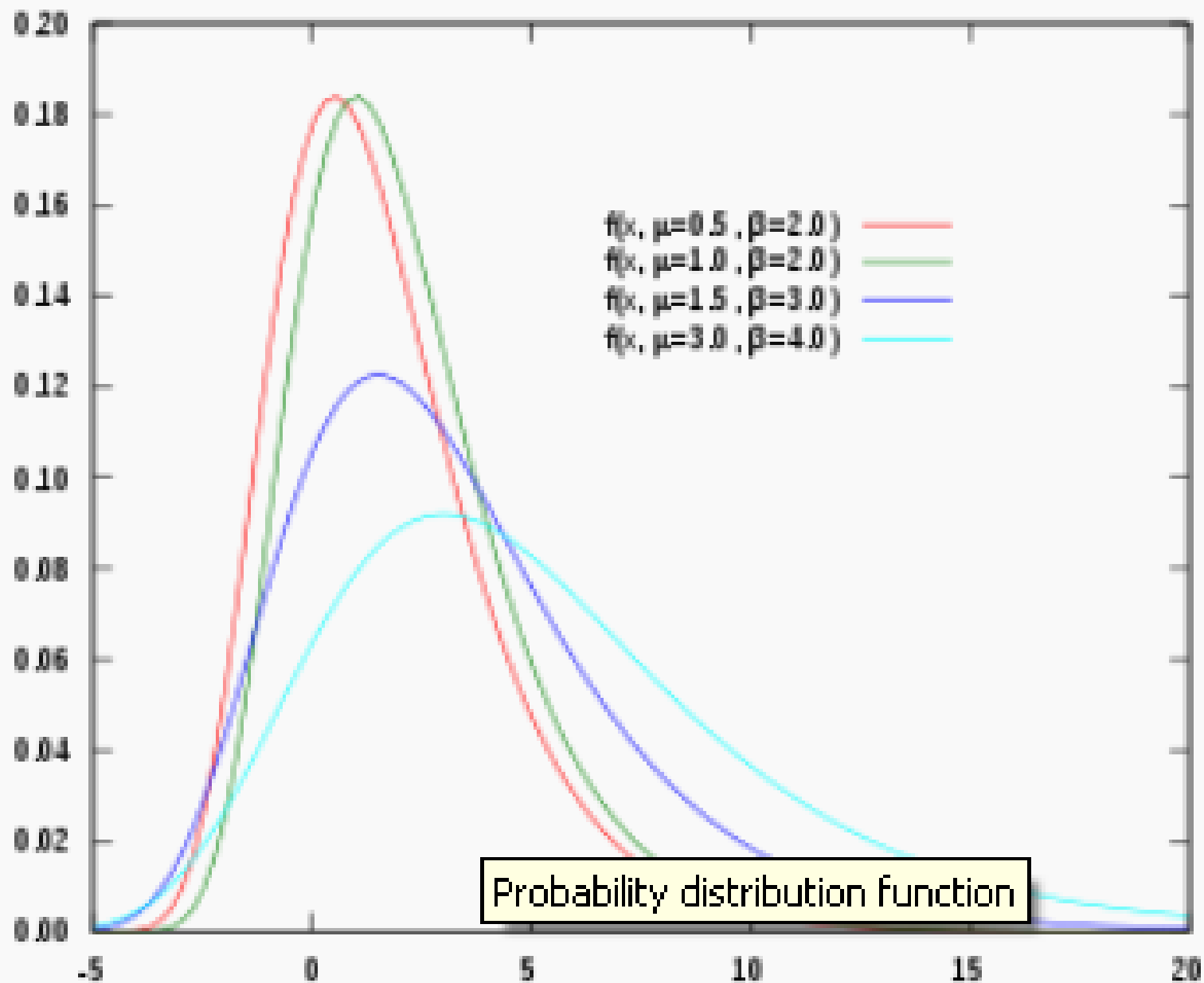
# Extreme value distribution

- the **Gumbel distribution** is used to model the distribution of the maximum (or the minimum) of a number of samples of various distributions.
- For example we would use it to represent the distribution of the maximum level of a river in a particular year if we had the list of maximum values for the past ten years. It is useful in predicting the chance that an extreme earthquake, flood or other natural disaster will occur.

# Extreme value distribution

<b>support:</b>	$x \in (-\infty; +\infty)$
<b>pdf:</b>	$\frac{z e^{-z}}{\beta}$ <p>where <math>z = e^{-\frac{x-\mu}{\beta}}</math></p>
<b>cdf:</b>	$\exp(-e^{-(x-\mu)/\beta})$
<b>mean:</b>	$\mu + \beta \gamma$
<b>median:</b>	$\mu - \beta \ln(\ln(2))$
<b>mode:</b>	$\mu$
<b>variance:</b>	$\frac{\pi^2}{6} \beta^2$

# Probability density function





# BLAST

- In accordance with the Gumbel EVD, the probability  $p$  of observing a score  $S$  equal to or greater than  $x$  is given by the equation

$$p(S \geq x) = 1 - \exp\left(-e^{-\lambda(x-\mu)}\right)$$

$$\mu = [\log(Km'n')]/\lambda$$

# BLAST: E-value

- The expect score  $E$  of a database match is the number of times that an unrelated database sequence would obtain a score  $S$  higher than  $x$  by chance. The expectation  $E$  obtained in a search for a database of  $D$  sequences is given by

$$E \approx pD$$