



# **Risk and Safety**

**in**

# **Civil, Surveying and Environmental**

# **Engineering**

**Prof. Dr. Michael Havbro Faber**  
**Swiss Federal Institute of Technology**  
**ETH Zurich, Switzerland**



## Contents of Today's Lecture

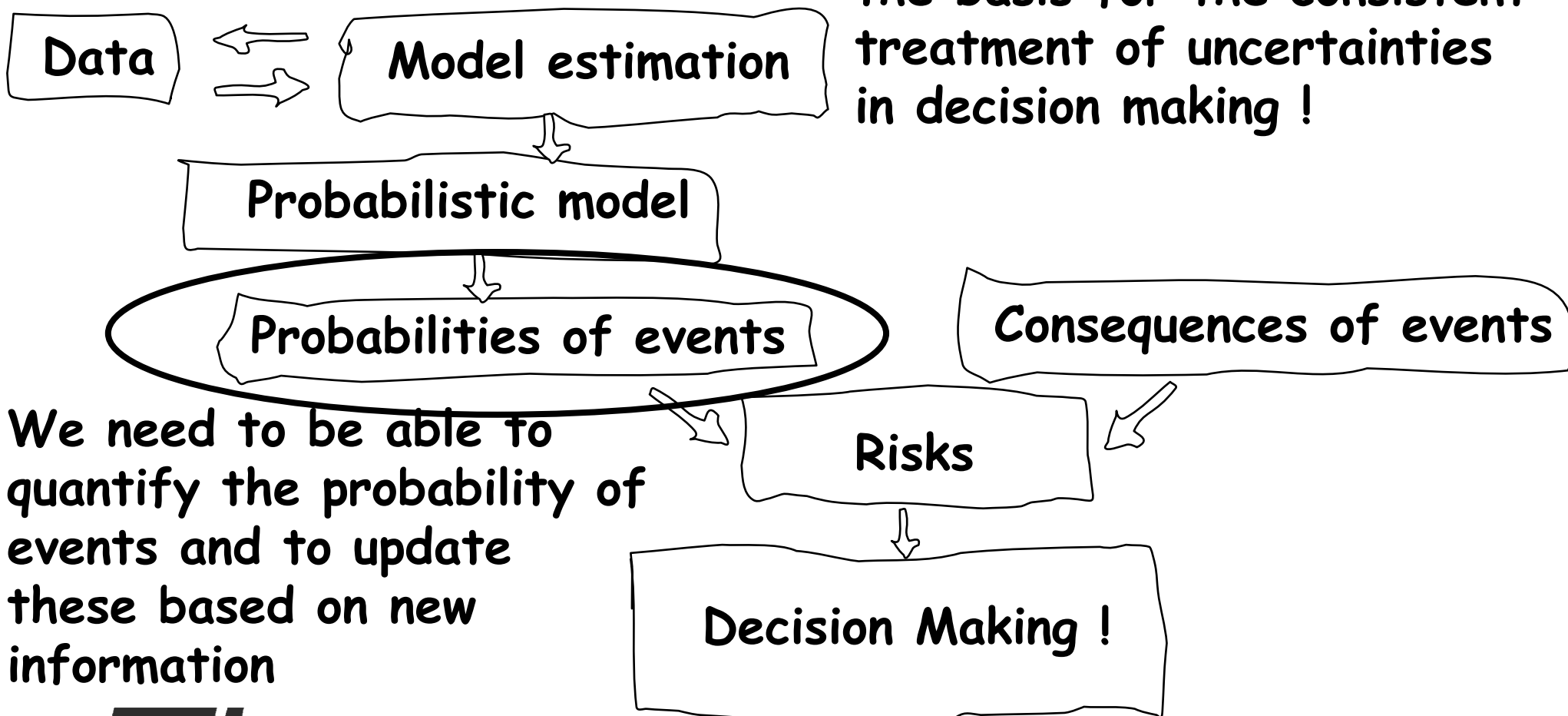
- **Probability theory**
- **Descriptive statistics**
- **Uncertainties in engineering decision making**
- **Probabilistic modelling**
- **Engineering model building**



# Overview of Probability Theory

- What are we aiming for ?

The probability theory provides the basis for the consistent treatment of uncertainties in decision making !



We need to be able to quantify the probability of events and to update these based on new information



# Interpretation of Probability

States of nature of which we have interest such as:

- a bridge failing due to excessive traffic loads
- a water reservoir being over-filled
- an electricity distribution system „falling out“
- a project being delayed

are in the following denoted „events“

we are generally interested in quantifying the probability that such events take place within a given „time frame“



# Interpretation of Probability

- There are in principle three different interpretations of probability

- **Frequentistic**  $P(A) = \lim_{n_{\text{exp}} \rightarrow \infty} \frac{N_A}{n_{\text{exp}}}$  for  $n_{\text{exp}} \rightarrow \infty$
- **Classical**  $P(A) = \frac{n_A}{n_{\text{tot}}}$
- **Bayesian**  $P(A) =$  degree of belief that  $A$  will occur



# Interpretation of Probability

Consider the probability of getting a „head“ when flipping a coin

- Frequentistic

$$P(A) = \frac{510}{1000} = 0.51$$

- Classical

$$P(A) = \frac{1}{2}$$

- Bayesian

$$P(A) = 0.5$$





# Conditional Probability and Bayes's Rule





# Conditional Probability and Bayes's Rule





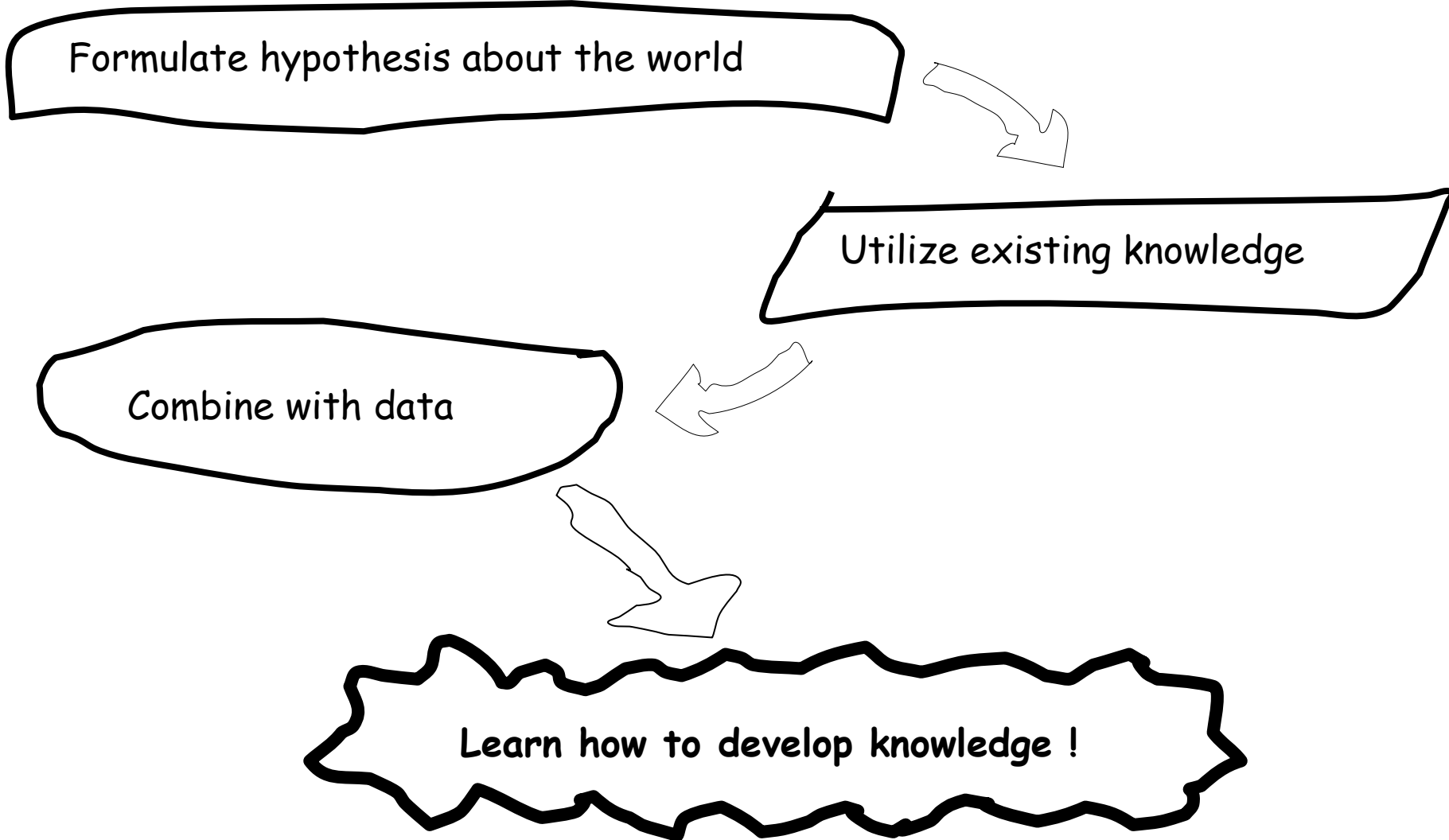


# Conditional Probability and Bayes's Rule





# Conditional Probability and Bayes's Rule





# Conditional Probability and Bayes's Rule

Conditional probabilities are of special interest as they provide the basis for utilizing new information in decision making.

The conditional probability of an event  $E_1$  given that event  $E_2$  has occurred is written as:

$$P(E_1|E_2) = \frac{P(E_1 \cap E_2)}{P(E_2)} \quad \text{Not defined if } P(E_2) = 0$$

The events  $E_1$  and  $E_2$  are said to be statistically independent if:

$$P(E_1|E_2) = P(E_1)$$



## Conditional Probability and Bayes's Rule

From 
$$P(E_1|E_2) = \frac{P(E_1 \cap E_2)}{P(E_2)}$$

it follows that 
$$P(E_1 \cap E_2) = P(E_2)P(E_1|E_2)$$

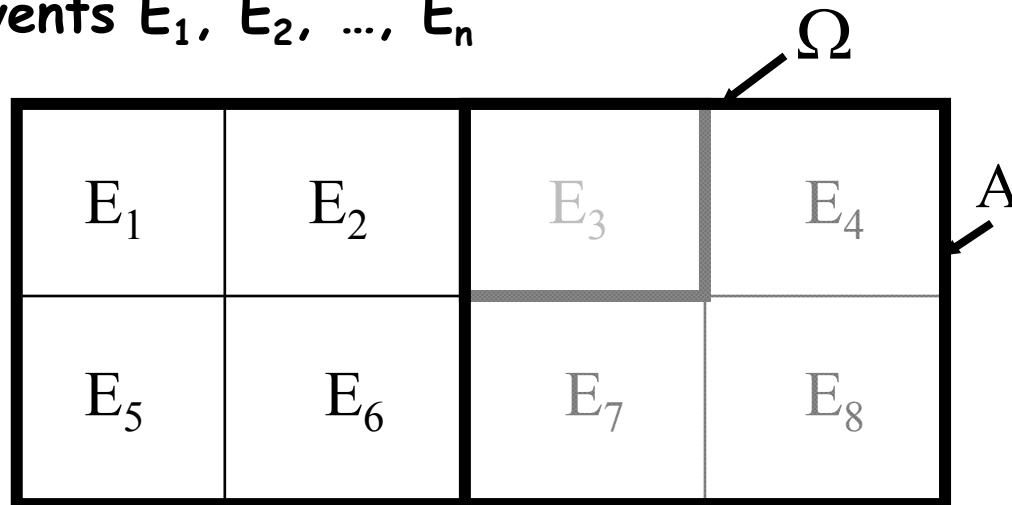
and when  $E_1$  and  $E_2$  are statistically independent it is

$$P(E_1 \cap E_2) = P(E_2)P(E_1)$$



# Conditional Probability and Bayes's Rule

Consider the sample space  $\Omega$  divided up into  $n$  mutually exclusive events  $E_1, E_2, \dots, E_n$



$$P(A) = P(A \cap E_1) + P(A \cap E_2) + \dots + P(A \cap E_n)$$

$$P(A|E_1)P(E_1) + P(A|E_2)P(E_2) + \dots + P(A|E_n)P(E_n) =$$

$$\sum_{i=1}^n P(A|E_i)P(E_i)$$



# Conditional Probability and Bayes's Rule

as there is  $P(A \cap E_i) = P(A|E_i)P(E_i) = P(E_i|A)P(A)$

we have

Likelihood

Prior

$$P(E_i|A) = \frac{P(A|E_i)P(E_i)}{P(A)} = \frac{P(A|E_i)P(E_i)}{\sum_{i=1}^n P(A|E_i)P(E_i)}$$

Posterior

Bayes Rule



Reverend Thomas Bayes  
(1702-1764)



# Conditional Probability and Bayes's Rule

## Example – inspection of degrading concrete structure

A reinforced concrete structure is considered

It is assumed (known) that the probability that corrosion of the reinforcement has initiated is:  $P(CI) = 0.01$

The state of the reinforcement of the considered beam is unknown and NDE tests are invoked

The quality of the test is specified by the probabilities

- that the test will indicate corrosion given that corrosion has initiated  $P(I|CI)$

- that the test will indicate corrosion given that corrosion has not initiated  $P(I|\overline{CI})$





# Conditional Probability and Bayes's Rule

## Example – inspection of degrading concrete structure

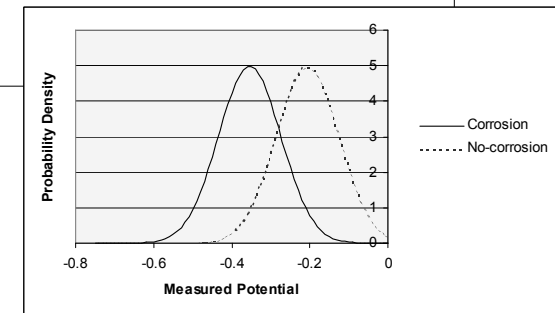
By comparison of a large number of NDE measurements with the real condition of concrete structures it has been found that

$$P(I|CI) = 0.8$$

$$P(I|\overline{CI}) = 0.1$$

We now seek the probability of corrosion given that we get an indication of corrosion by the NDE inspection i.e.

$$P(CI|I) = ?$$



$$P(CI|I) = \frac{\overset{\text{Likelihood}}{P(I|CI)} \overset{\text{Prior}}{P(CI)}}{\overset{\text{Posterior}}{P(I|CI)P(CI) + P(I|\overline{CI})P(\overline{CI})}}$$

$$P(CI|I) = \frac{0.008}{0.107} = 0.075$$

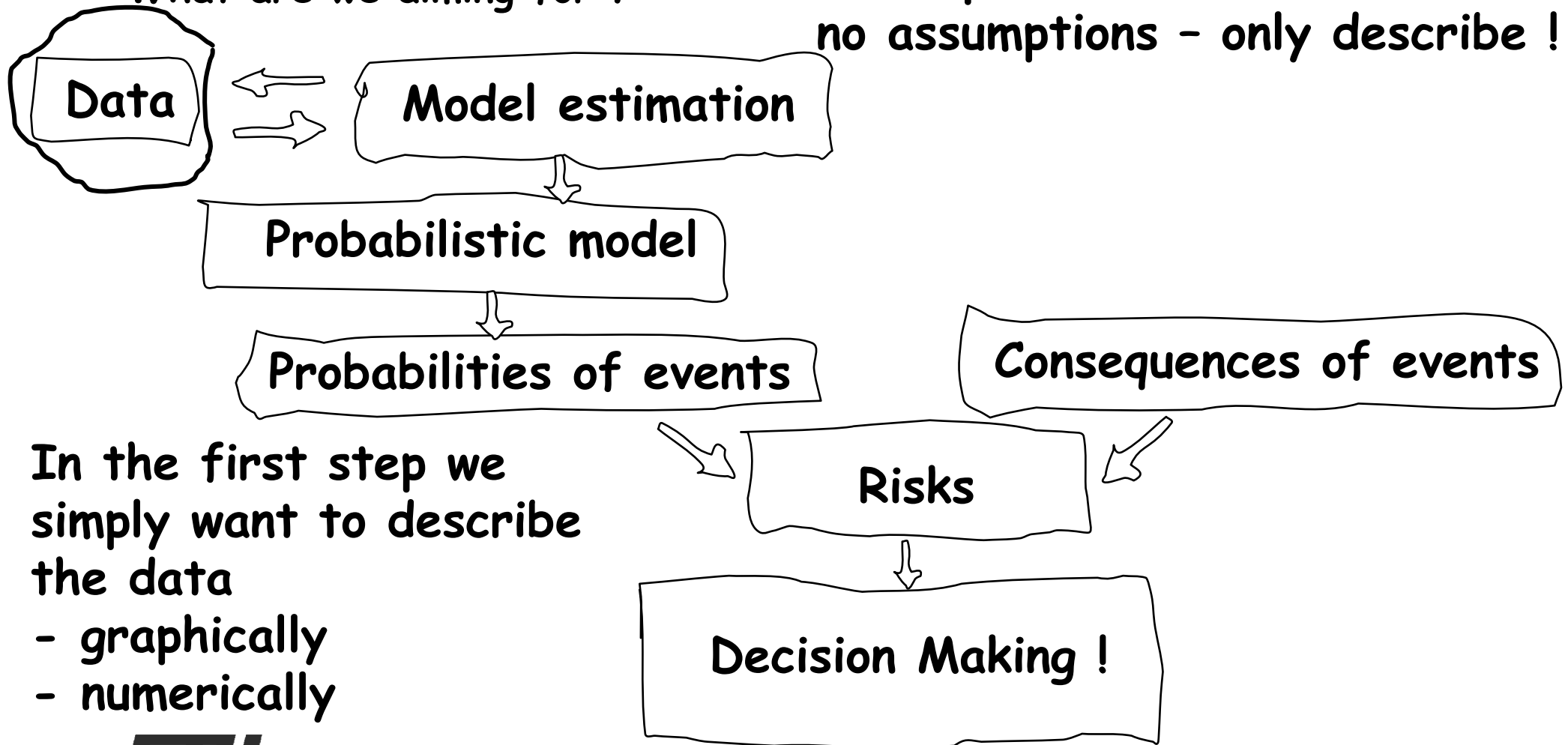




# Overview of Descriptive Statistics

- What are we aiming for ?

Descriptive statistics make no assumptions - only describe !



In the first step we simply want to describe the data

- graphically
- numerically



# Numerical Summaries

- Central measures:

Sample mean :

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

If one number should be given to represent a data set typically the sample mean would be chosen

Median : The 0.5 quantile (obtained from ordered data sets, see quantile plots)

Mode : Most frequent value - obtained from histograms



# Numerical Summaries

- Dispersion measures:

**Sample variance:** 
$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$
  $S$  : standard deviation

Indicator of variability around the sample mean

**Sample coefficient of variation (CoV):** 
$$v = \frac{S}{\bar{x}}$$

Indicator of variability relative to the sample mean



# Numerical Summaries

- Other measures:

Sample skewness: 
$$\eta = \frac{1}{n} \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})^3}{s^3}$$
 Measure of symmetry

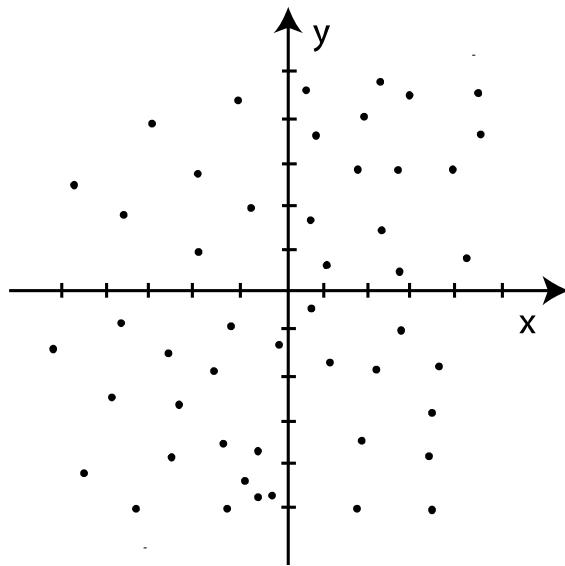
Sample kurtosis 
$$K = \frac{1}{n} \cdot \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{s^4}$$
 Measure of peakedness



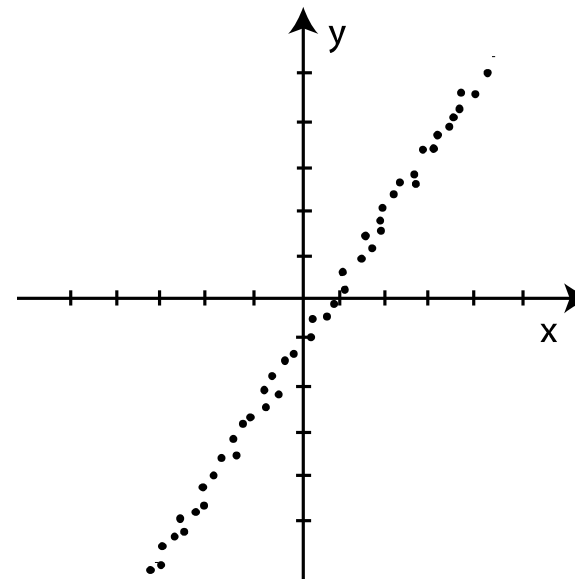
# Numerical Summaries

- Measures of correlation (linear dependency between data pairs):

## 2-dimensional scatter plots



Almost no dependency



Almost full dependency



# Numerical Summaries

- Measures of correlation (linear dependency between data pairs):

Sample covariance: 
$$s^2_{XY} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})$$

The sum will get positive contributions in case of low-low or high-high data pairs

Sample coefficient of correlation: 
$$r_{XY} = \frac{1}{n} \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{s_X \cdot s_Y}$$

$r_{XY}$  is limited in the interval -1 to +1



# Numerical Summaries

- **Summary:**

**Central measures:**

- **sample mean value:** The center of gravity of a data set
- **sample median:** The mid value of a data set
- **sample mode:** The most frequent value/range of a data set

**Dispersion measures:**

- **sample variance:** The distribution around the sample mean
- **sample CoV:** The variability relative to the sample mean

**Other measures:**

- **sample skewness:** The skewness relative to the sample mean
- **sample kurtosis:** The peakedness around the sample mean

**Measures of correlation:**

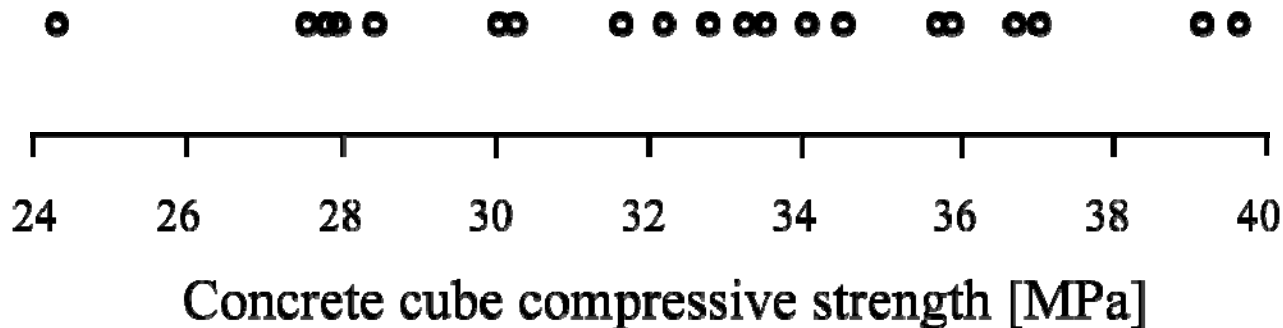
- **sample covariance:** Tendency for high-high, low-low and high-low pairs in two data sets
- **sample coefficient of correlation :** Normalized coefficient between -1 and +1



# Graphical Representations

- Assume that we have a set of data (observations of concrete compressive strength)

The simplest representation of the data is the one-dimensional scatter plot



$i$	Unordered $x_i$	Ordered $x_i^o$
1	35.8	24.4
2	39.2	27.6
3	34.6	27.8
4	27.6	27.9
5	37.1	28.5
6	33.3	30.1
7	32.8	30.3
8	34.1	31.7
9	27.9	32.2
10	24.4	32.8
11	27.8	33.3
12	33.5	33.5
13	35.9	34.1
14	39.7	34.6
15	28.5	35.8
16	30.3	35.9
17	31.7	36.8
18	32.2	37.1
19	36.8	39.2
20	30.1	39.7



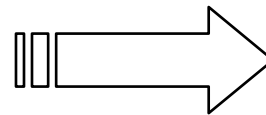


# Graphical Representations

- Histograms

The data are grouped into intervals

$i$	Unordered $x_i$	Ordered $x_i^o$
1	35.8	24.4
2	39.2	27.6
3	34.6	27.8
4	27.6	27.9
5	37.1	28.5
6	33.3	30.1
7	32.8	30.3
8	34.1	31.7
9	27.9	32.2
10	24.4	32.8
11	27.8	33.3
12	33.5	33.5
13	35.9	34.1
14	39.7	34.6
15	28.5	35.8
16	30.3	35.9
17	31.7	36.8
18	32.2	37.1
19	36.8	39.2
20	30.1	39.7



Interval	Midpoint	Number of observations	Frequency [%]	Cumulative frequency
23-26	24.5	1	5	0.05
26-29	27.5	4	20	0.25
29-32	30.5	3	15	0.40
32-35	33.5	6	30	0.70
35-38	36.5	4	20	0.90
38-41	39.5	2	10	1.00

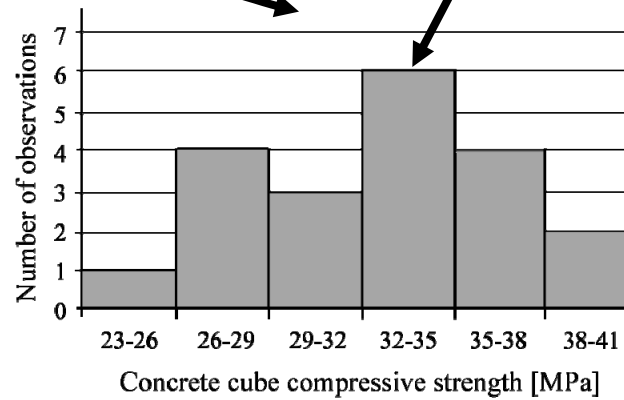


# Graphical Representations

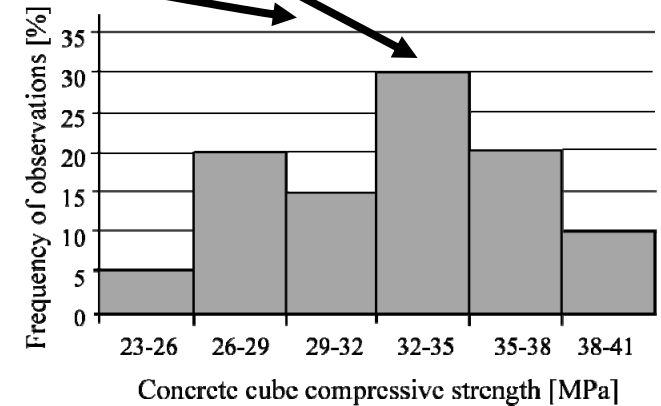
- Histograms

The grouped data are plotted

Interval	Midpoint	Number of observations	Frequency [%]	Cumulative frequency
23-26	24.5	1	5	0.05
26-29	27.5	4	20	0.25
29-32	30.5	3	15	0.40
32-35	33.5	6	30	0.70
35-38	36.5	4	20	0.90
38-41	39.5	2	10	1.00



Simple histogram



Frequency distribution



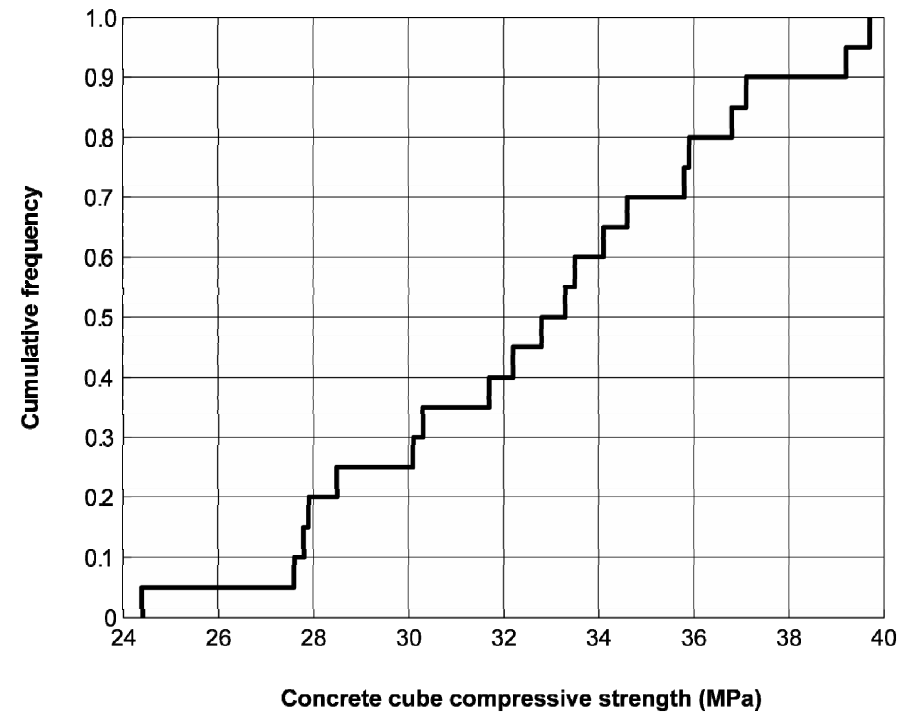
# Graphical Representations

- Histograms

The grouped data are plotted



Interval	Midpoint	Number of observations	Frequency [%]	Cumulative frequency
23-26	24.5	1	5	0.05
26-29	27.5	4	20	0.25
29-32	30.5	3	15	0.40
32-35	33.5	6	30	0.70
35-38	36.5	4	20	0.90
38-41	39.5	2	10	1.00





# Graphical Representations

- Quantile plots

**Definition : the Q-quantile corresponds to the value in a data set which is exceeded by  $100\% - Q \times 100\%$  of the data**

**e.g. the 0.75 quantile is exceeded by  $100\% - 0.75 \times 100\%$   
= 25% of the data**

**Quantile plots are generated by plotting the data against their quantile values**



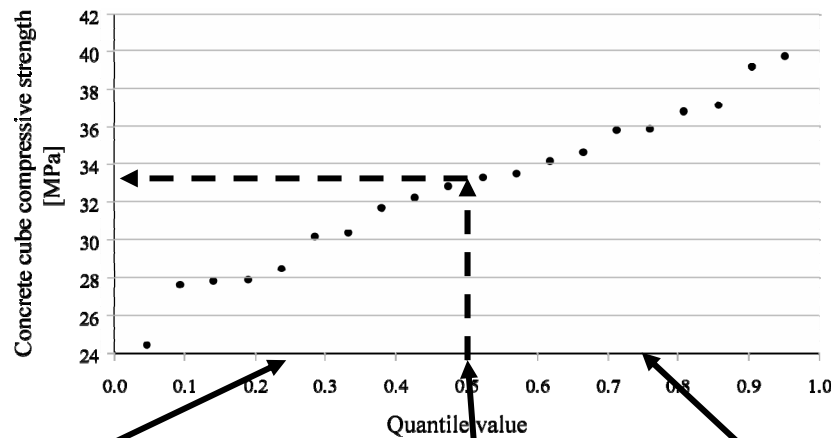
# Graphical Representations

- Quantile plots

The quantiles are calculated from the ordered data set as:

$$Q_i = \frac{i}{1+n}$$

<i>i</i>	Ordered $x_i^o$	$Q_i$
1	24.4	0.048
2	27.6	0.095
3	27.8	0.143
4	27.9	0.190
5	28.5	0.238
6	30.1	0.286
7	30.3	0.333
8	31.7	0.381
9	32.2	0.429
10	32.8	0.476
11	33.3	0.524
12	33.5	0.571
13	34.1	0.619
14	34.6	0.667
15	35.8	0.714
16	35.9	0.762
17	36.8	0.810
18	37.1	0.857
19	39.2	0.905
20	39.7	0.952



Lower quartile = 0.25 quartile value

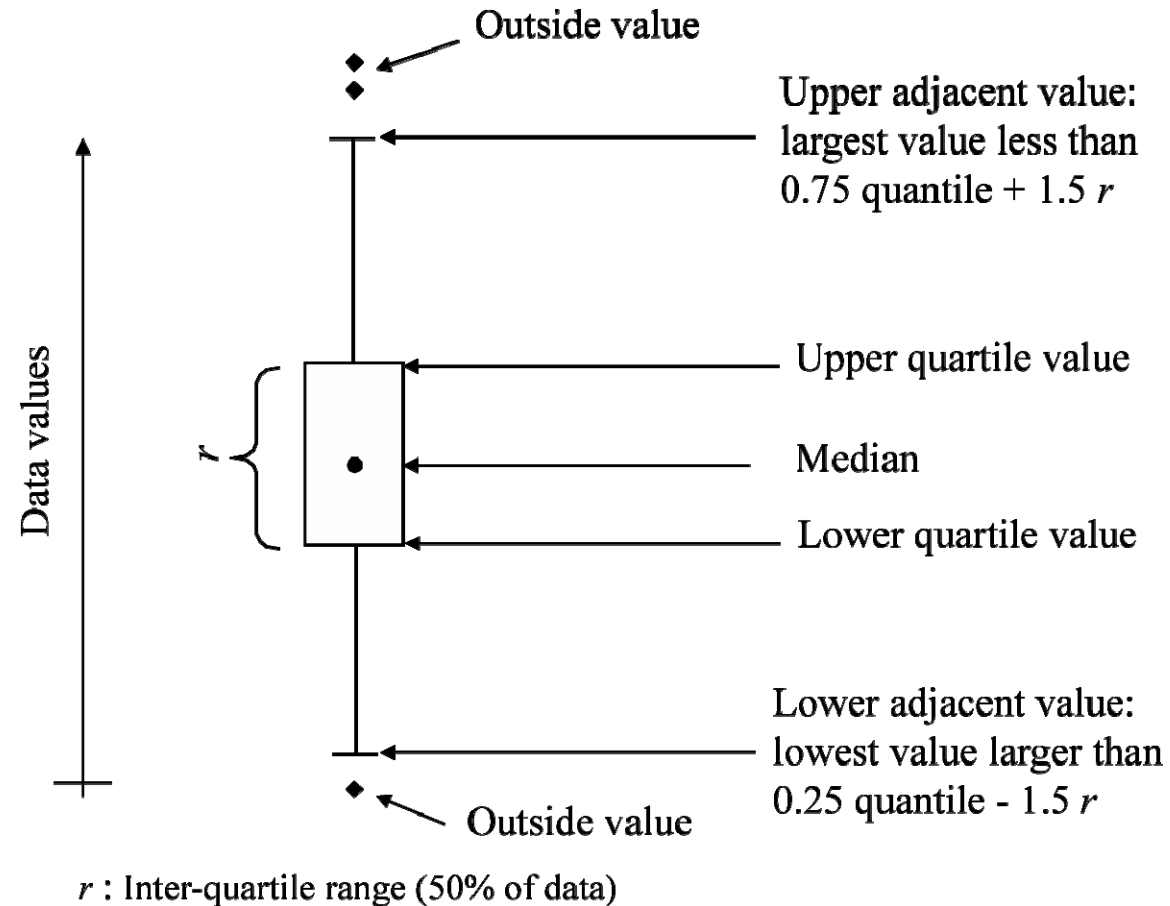
Median = 0.5 quartile value

Upper quartile = 0.75 quartile value



# Graphical Representations

- Tukey Box plots

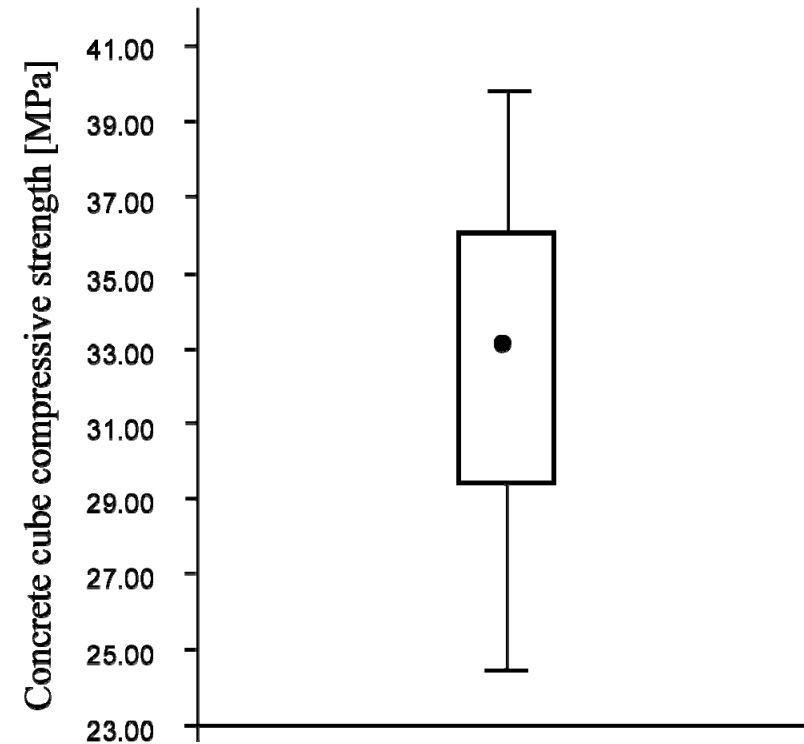




# Graphical Representations

- Tukey Box plots

Statistic	Value
Lower quartile	29.30
Lower adjacent value	24.40
Median	33.05
Upper adjacent value	39.70
Upper quartile	35.85





# Graphical Representations

- **Summary**

**One-dimensional scatter plots** : illustrate the range and distribution of a data sets along one axis, indicate symmetry.

**Histograms:** illustrate how the data are distributed over the range of data, indicate mode and symmetry.

**Quantile plots:** Illustrate median, distribution and symmetry

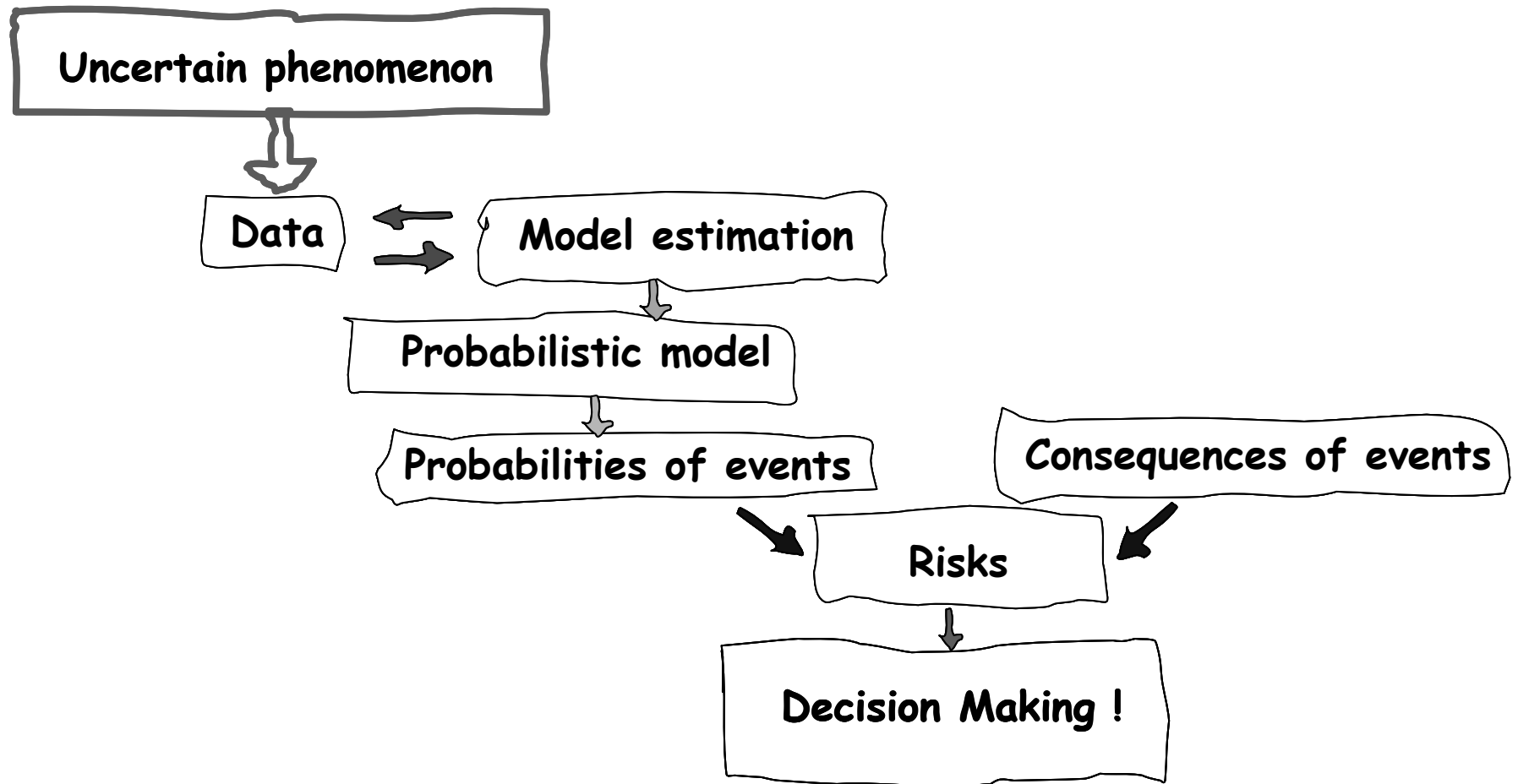
**Tukey - Box plots:** Illustrate median, upper/lower quartiles, symmetry and distribution





# Overview of Uncertainty Modelling

- Why uncertainty modelling





# Uncertainties in Engineering Problems

Different types of uncertainties influence decision making

- **Inherent natural variability - aleatory uncertainty**
  - result of throwing dices
  - variations in material properties
  - variations of wind loads
  - variations in rain fall
- **Model uncertainty - epistemic uncertainty**
  - lack of knowledge (future developments)
  - inadequate/imprecise models (simplistic physical modelling)
- **Statistical uncertainties - epistemic uncertainty**
  - sparse information/small number of data



# Uncertainties in Engineering Problems

- Consider as an example a dike structure
  - the design (height) of the dike will be determining the frequency of floods
  - if exact models are available for the prediction of future water levels and our knowledge about the input parameters is perfect then we can calculate the frequency of floods (per year) - a deterministic world !
  - even if the world would be deterministic - we would not have perfect information about it - so we might as well consider the world as random



# Uncertainties in Engineering Problems

In principle the so-called

inherent physical uncertainty (aleatory - Type I)

is the uncertainty caused by the fact that the world is random, however, another pragmatic viewpoint is to define this type of uncertainty as

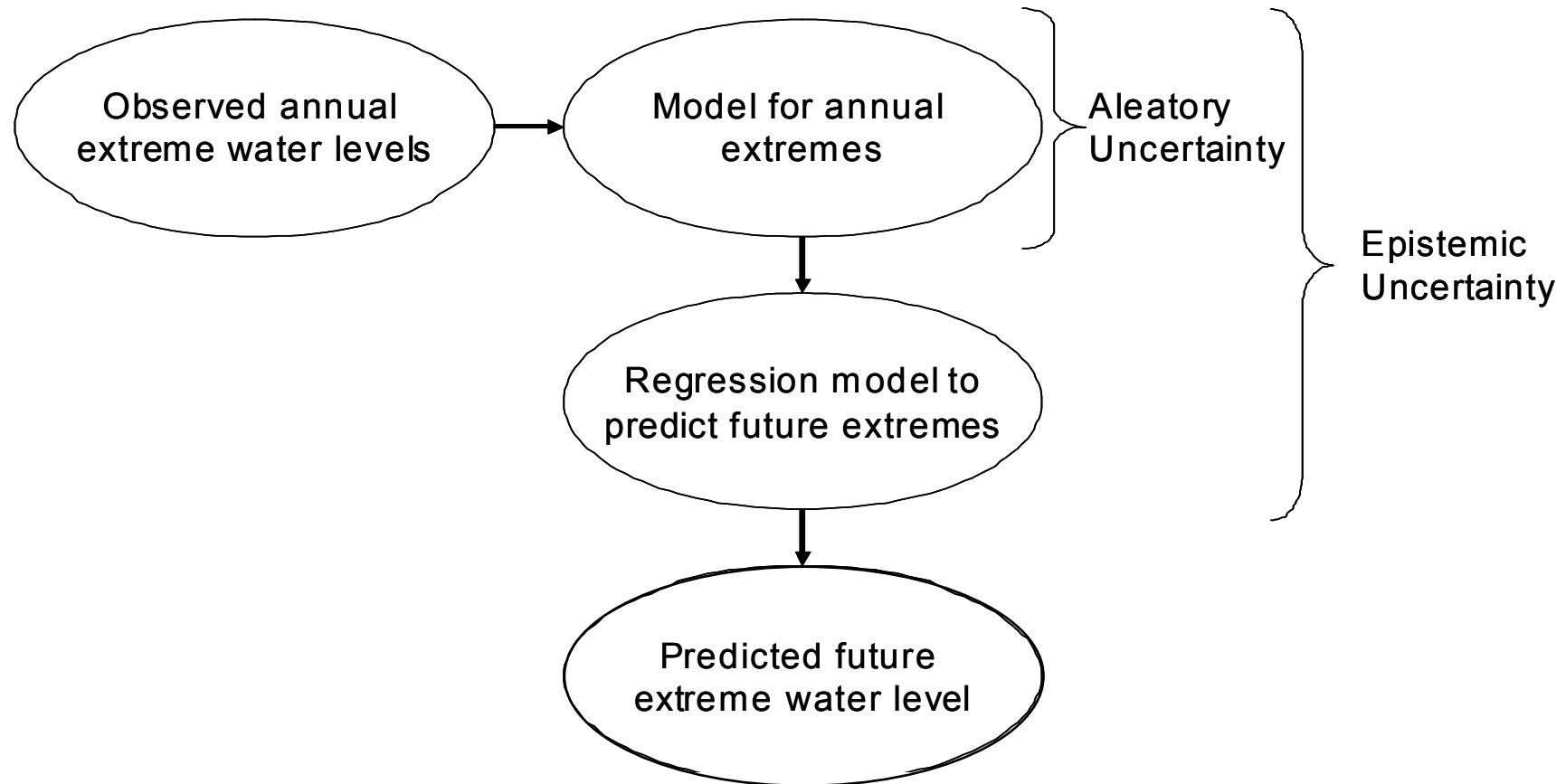
any uncertainty which cannot be reduced by means of collection of additional information

the uncertainty which can be reduced is then the

model and statistical uncertainties (epistemic - Type II)



# Uncertainties in Engineering Problems





## Uncertainties in Engineering Problems

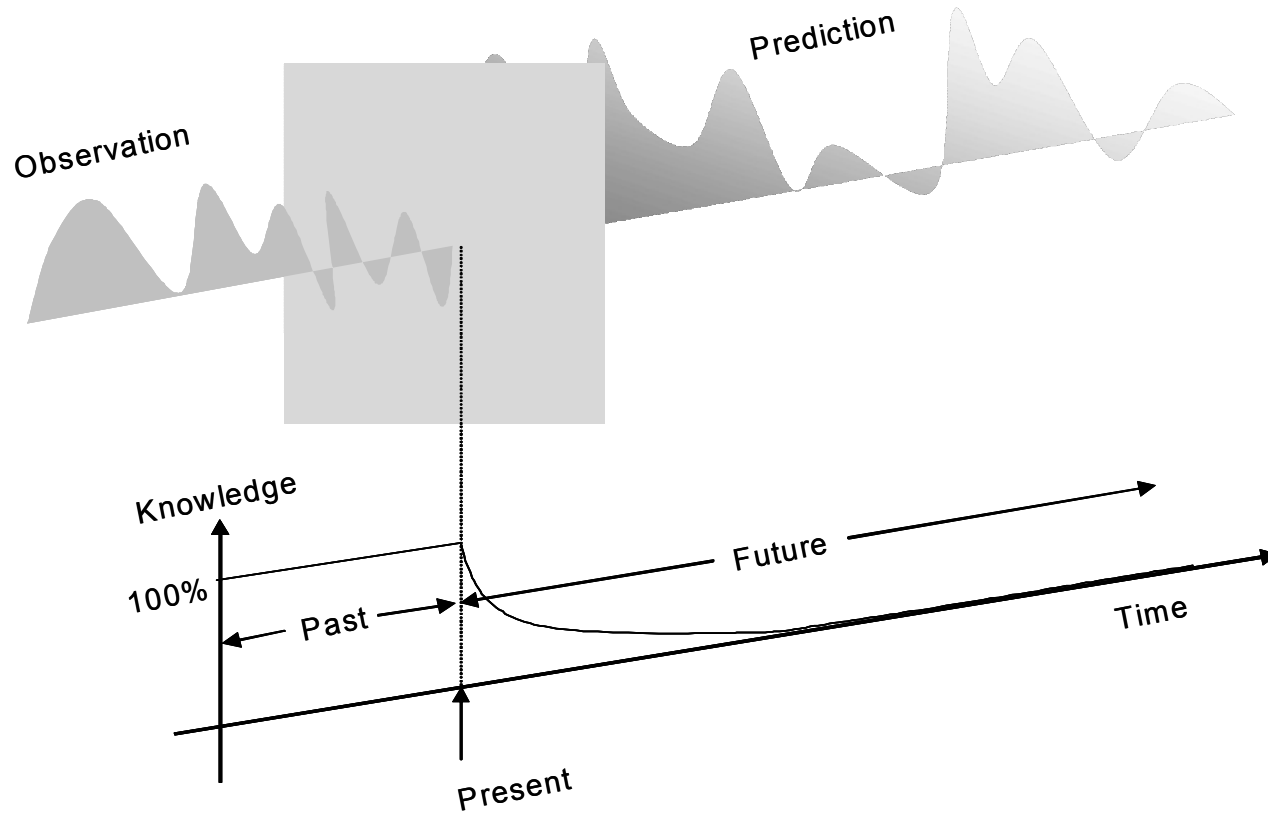
The relative contribution of aleatory and epistemic uncertainty to the prediction of future water levels is thus influenced directly by the applied models

refining a model might reduce the epistemic uncertainty - but in general also changes the contribution of aleatory uncertainty

the uncertainty structure of a problem can thus be said to be scale dependent !



# Uncertainties in Engineering Problems



The uncertainty structure changes also as function of time - is thus time dependent !



## Random Variables

- Probability distribution and density functions

A random variable is denoted with capital letters :  $X$

A realization of a random variable is denoted with small letters :  $x$

We distinguish between

- *continuous random variables* : can take any value in a given range
- *discrete random variables* : can take only discrete values





# Random Variables

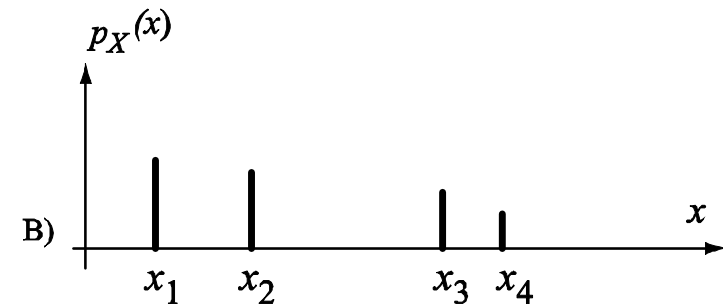
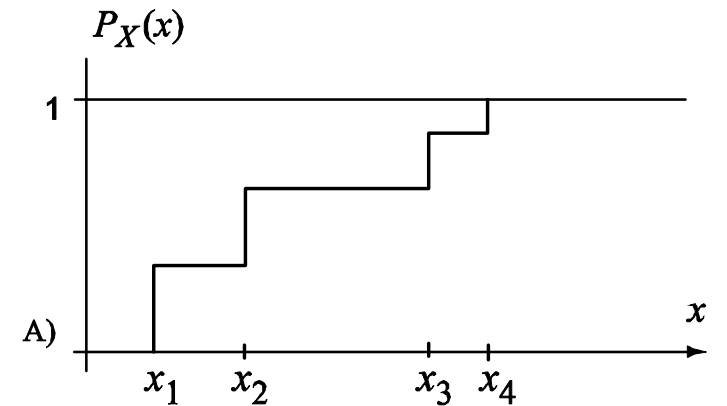
- Probability distribution and density functions

The probability that the outcome of a discrete random variable  $X$  is smaller than  $x$  is denoted the *probability distribution function*

$$P_X(x) = \sum_{x_i < x} p_X(x_i)$$

The *probability density function* for a discrete random variable is defined by

$$p_X(x_i) = P(X = x_i)$$





# Random Variables

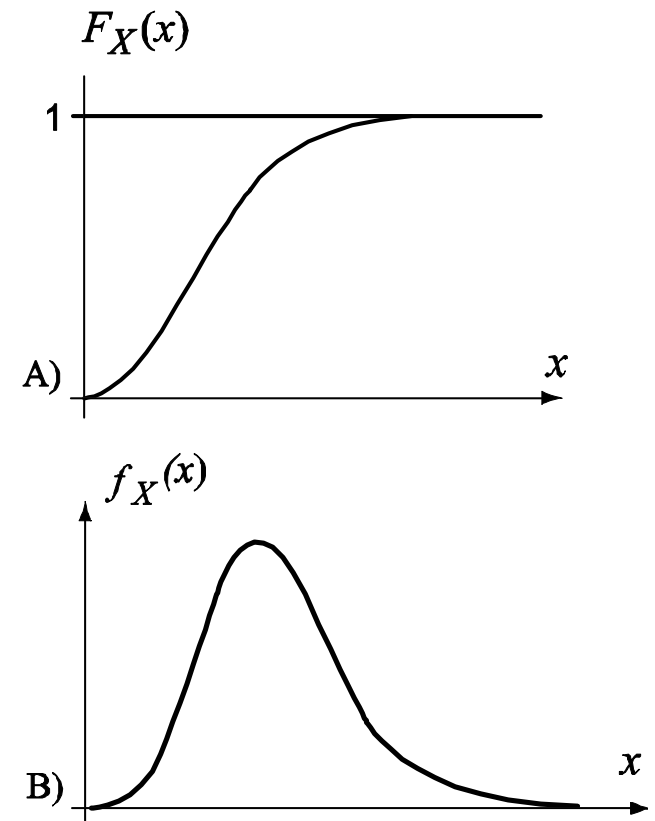
- Probability distribution and density functions

The probability that the outcome of a continuous random variable  $X$  is smaller than  $x$  is denoted the *probability distribution function*

$$F_X(x) = P(X < x)$$

The probability density function for a continuous random variable is defined by

$$f_X(x) = \frac{\partial F_X(x)}{\partial x}$$





## Random Variables

- Moments of random variables and the expectation operator

Probability distribution and density function can be described in terms of their parameters  $\mathbf{p}$  or their moments

Often we write

$$F_X(x, \mathbf{p}) \quad f_X(x, \mathbf{p})$$

Parameters

The parameters can be related to the moments and visa versa



## Random Variables

- Moments of random variables and the expectation operator

The  $i$ 'th moment  $m_i$  for a continuous random variable  $X$  is defined through

$$m_i = \int_{-\infty}^{\infty} x^i \cdot f_X(x) dx$$

The *expected value*  $E[X]$  of a continuous random variable  $X$  is defined accordingly as the first moment

$$\mu_X = E[X] = \int_{-\infty}^{\infty} x \cdot f_X(x) dx$$



## Random Variables

- Moments of random variables and the expectation operator

The  $i$ 'th moment  $m_i$  for a discrete random variable  $X$  is defined through

$$m_i = \sum_{j=1}^n x_j^i \cdot p_X(x_j)$$

The *expected value*  $E[X]$  of a discrete random variable  $X$  is defined accordingly as the first moment

$$\mu_X = E[X] = \sum_{j=1}^n x_j \cdot p_X(x_j)$$





## Random Variables

- Moments of random variables and the expectation operator

The ratio between the standard deviation and the expected value of a random variable is called the *Coefficient of Variation CoV* and is defined as

$$CoV[X] = \frac{\sigma_X}{\mu_X}$$

 **Dimensionless**

a useful characteristic to indicate the variability of the random variable around its expected value



# Random Variables

- Typical probability distribution functions in engineering

Normal : sum of random effects

Log-Normal: product of random effects

Exponential: waiting times

Gamma: Sum of waiting times

Beta: Flexible modeling function

Distribution type	Parameters	Moments
<b>Rectangular</b> $a \leq x \leq b$ $f_X(x) = \frac{1}{b-a}$ $F_X(x) = \frac{x-a}{b-a}$	$a$ $b$	$\mu = \frac{a+b}{2}$ $\sigma = \frac{b-a}{\sqrt{12}}$
<b>Normal</b> $f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$ $F_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right) dx$	$\mu$ $\sigma > 0$	$\mu$ $\sigma$
<b>Shifted Lognormal</b> $x > \varepsilon$ $f_X(x) = \frac{1}{(x-\varepsilon)\zeta\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{\ln(x-\varepsilon)-\lambda}{\zeta}\right)^2\right)$ $F_X(x) = \Phi\left(\frac{\ln(x-\varepsilon)-\lambda}{\zeta}\right)$	$\lambda$ $\zeta > 0$ $\varepsilon$	$\mu = \varepsilon + \exp\left(\lambda + \frac{\zeta^2}{2}\right)$ $\sigma = \exp\left(\lambda + \frac{\zeta^2}{2}\right) \sqrt{\exp(\zeta^2) - 1}$
<b>Shifted Exponential</b> $x \geq \varepsilon$ $f_X(x) = \lambda \exp(-\lambda(x-\varepsilon))$ $F_X(x) = 1 - e^{-\lambda(x-\varepsilon)}$	$\varepsilon$ $\lambda > 0$	$\mu = \varepsilon + \frac{1}{\lambda}$ $\sigma = \frac{1}{\lambda}$
<b>Gamma</b> $x \geq 0$ $f_X(x) = \frac{b^p}{\Gamma(p)} \exp(-bx)x^{p-1}$ $F_X(x) = \frac{\Gamma(bx, p)}{\Gamma(p)}$	$p > 0$ $b > 0$	$\mu = \frac{p}{b}$ $\sigma = \frac{\sqrt{p}}{b}$
<b>Beta</b> $a \leq x \leq b, r, t \geq 1$ $f_X(x) = \frac{\Gamma(r+t)}{\Gamma(r)\Gamma(t)} \frac{(x-a)^{r-1}(b-x)^{t-1}}{(b-a)^{r+t-1}}$ $F_X(x) = \frac{\Gamma(r+t)}{\Gamma(r)\Gamma(t)} \int_a^x \frac{(u-a)^{r-1}(b-u)^{t-1}}{(b-a)^{r+t-1}} du$	$a$ $b$ $r > 1$ $t > 1$	$\mu = a + (b-a) \frac{r}{r+t}$ $\sigma = \frac{b-a}{r+t} \sqrt{\frac{rt}{r+t+1}}$





# Random Variables

- The Normal distribution

The analytical form of the Normal distribution may be derived by repeated use of the result regarding the probability density function for the sum of two random variables

The normal distribution is very frequently applied in engineering modelling when a random quantity can be assumed to be composed as a sum of a number of individual contributions.

A linear combination  $S$  of  $n$  Normal distributed random variables  $X_i, i=1,2,\dots,n$  is thus also a Normal distributed random variable

$$S = a_0 + \sum_{i=1}^n a_i X_i$$



# Random Variables

- The Normal distribution:

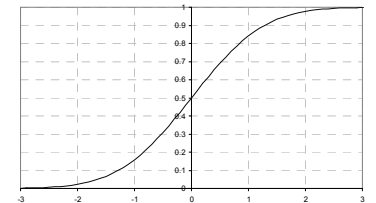
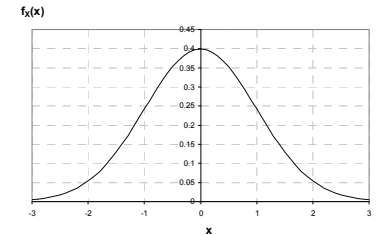
In the case where the mean value is equal to zero and the standard deviation is equal to 1 the random variable is said to be *standardized*.

$$Z = \frac{X - \mu_X}{\sigma_X} \quad \text{Standardized random variable}$$

$$f_Z(z) = \varphi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z^2\right)$$

$$F_Z(z) = \Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z \exp\left(-\frac{1}{2}x^2\right) dx$$

Standard normal





# Stochastic Processes and Extremes

- Random quantities may be “time variant” in the sense that they take new values at different times or at new trials.
  - If the new realizations occur at discrete times and have discrete values the random quantity is called a random sequence

failure events, traffic congestions,...

- If the new realizations occur continuously in time and take continuous values the random quantity is called a random process or stochastic process

wind velocity, wave heights,...



# Stochastic Processes and Extremes

- Random sequences

The Poisson counting process is one of the most commonly applied families of probability distributions applied in reliability theory

The process  $N(t)$  denoting the number of events in a (time) interval  $(t, t+Dt[$  is called a Poisson process if the following conditions are fulfilled:

- 1) the probability of one event in the interval  $(t, t+Dt[$  is asymptotically proportional to  $Dt$ .
- 2) the probability of more than one event in the interval  $(t, t+Dt[$  is a function of higher order of  $Dt$  for  $Dt \rightarrow 0$ .
- 3) events in disjoint intervals are mutually independent.

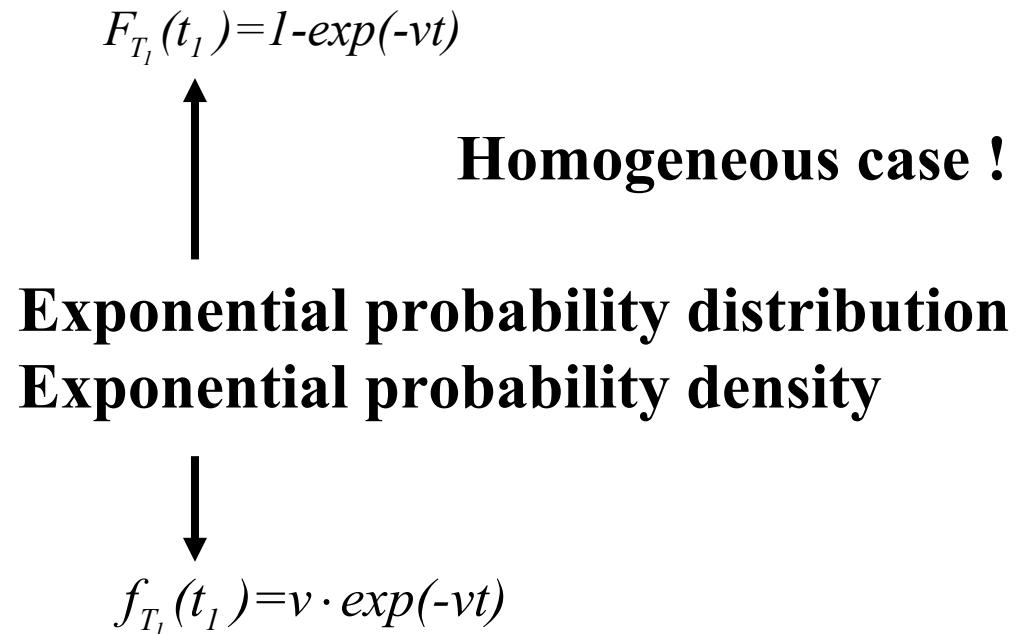


# Stochastic Processes and Extremes

- Random sequences

The probability distribution function of the (waiting) time till the first event  $T_1$  is now easily derived recognizing that the probability of  $T_1 > t$  is equal to  $P_0(t)$  we get:

$$\begin{aligned} F_{T_1}(t_1) &= 1 - P_0(t_1) \\ &= 1 - \exp\left(-\int_0^{t_1} \nu(\tau) d\tau\right) \end{aligned}$$

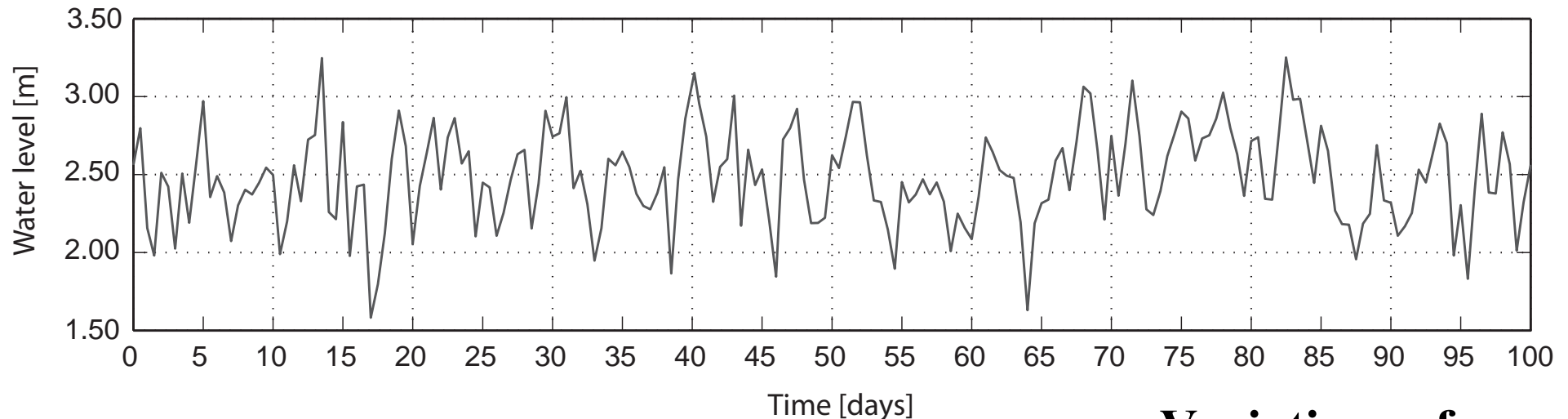




# Stochastic Processes and Extremes

- **Continuous random processes**

**A continuous random process is a random process which has realizations continuously over time and for which the realizations belong to a continuous sample space.**



**Variations of:  
water levels  
wind speed  
rain fall**



# Stochastic Processes and Extremes

- Continuous random processes

The mean value of the possible realizations of a random process is given as:

$$\mu_X(t) = E[X(t)] = \int_{-\infty}^{\infty} x \cdot f_X(x,t) dx$$

↑  
**Function of time !**

The correlation between realizations at any two points in time is given as:

$$R_{XX}(t_1, t_2) = E[X(t_1)X(t_2)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 \cdot x_2 \cdot f_{XX}(x_1, x_2; t_1, t_2) dx_1 dx_2$$

**Auto-correlation function** – refers to a scalar valued random process



# Stochastic Processes and Extremes

## Extreme Value Distributions

In engineering we are often interested in extreme values i.e. the smallest or the largest value of a certain quantity within a certain time interval e.g.:

The largest earthquake in 1 year

The highest wave in a winter season

The largest rainfall in 100 years





# Stochastic Processes and Extremes

## Extreme Value Distributions

We could also be interested in the smallest or the largest value of a certain quantity within a certain volume or area unit e.g.:

The largest concentration of pesticides in a volume of soil

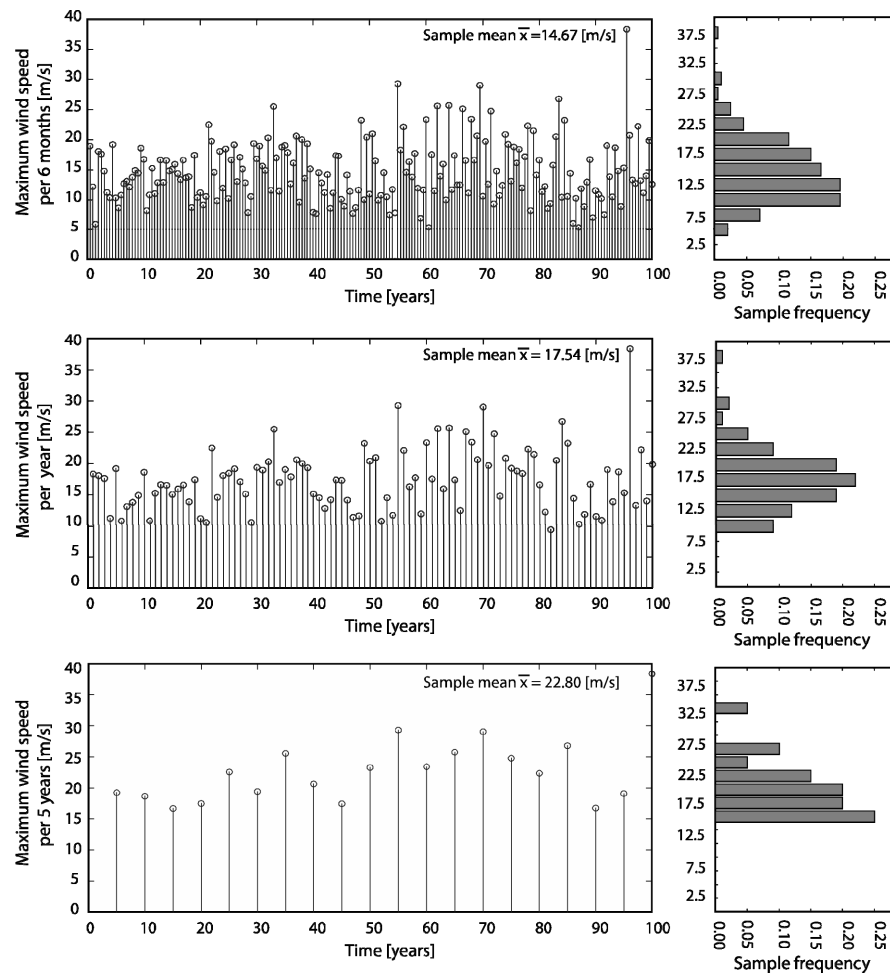
The weakest link in a chain

The smallest thickness of concrete cover



# Stochastic Processes and Extremes

Extremes of a random process:





# Stochastic Processes and Extremes

Return period for extreme events:

The return period for extreme events  $T_R$  may be defined as

$$T_R = n \cdot T = \frac{1}{(1 - F_{X,T}^{\max}(x))}$$

If the probability of exceeding  $x$  during a reference period of 1 year is 0.01 then the return period for exceedances is

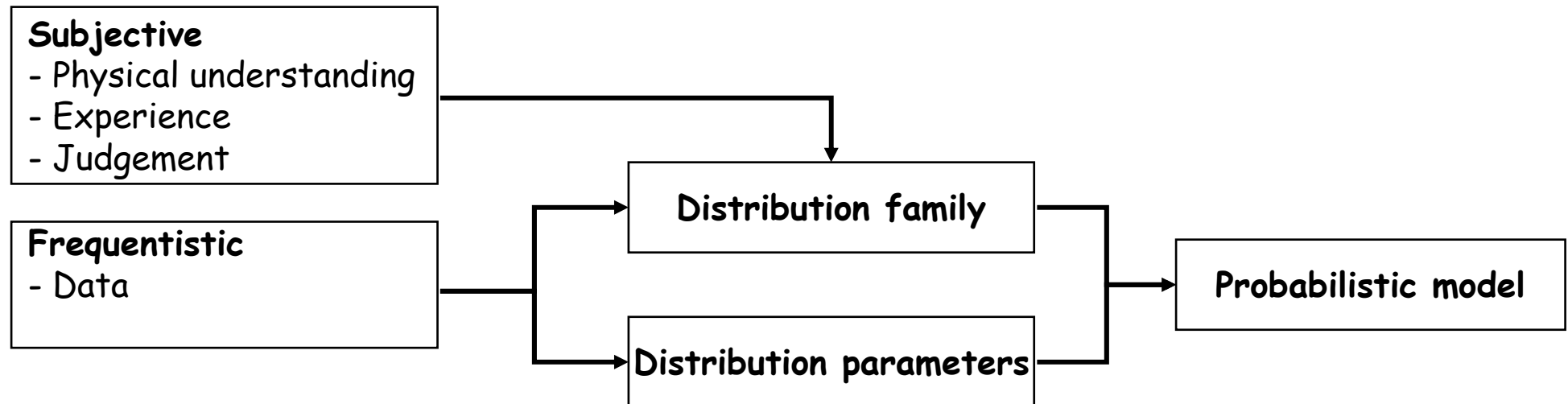
$$T_R = n \cdot T = \frac{1}{0.01} = 100 \cdot 1 = 100$$



# Overview of Estimation and Model Building

Different types of information is used when developing engineering models

- subjective information
- frequentististic information





# Overview of Estimation and Model Building

Model building may be seen to consist of five steps

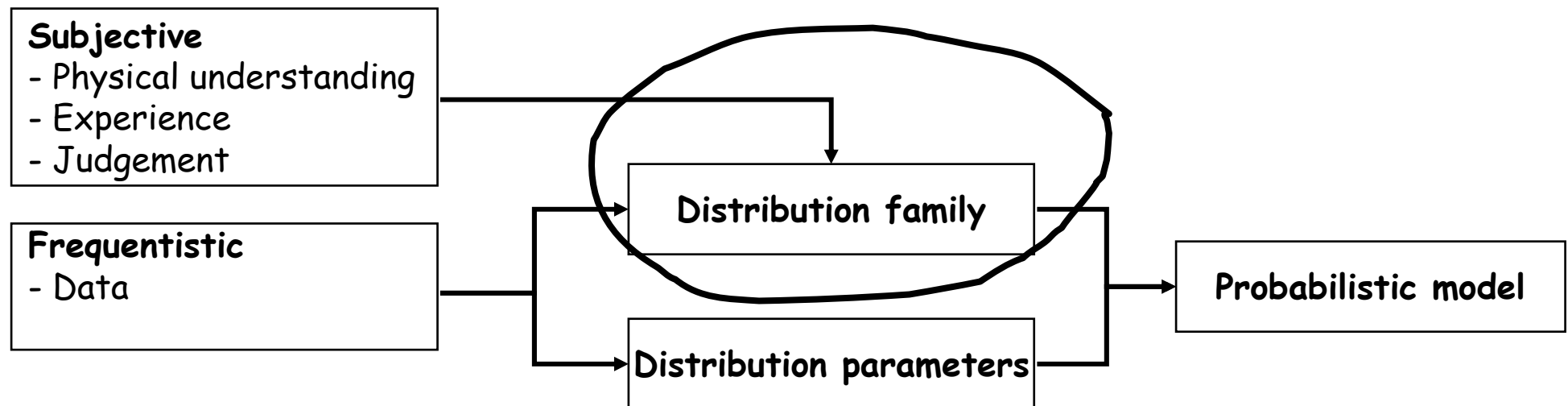
- 1) Assessment and statistical quantification of the available data
- 2) Selection of distribution function
- 3) Estimation of distribution parameters
- 4) Model verification
- 5) Model updating



# Overview of Estimation and Model Building

Different types of information is used when developing engineering models

- subjective information
- frequentististic information





# Estimation and Model Building

## Selection of probability distribution function

In engineering application it is often the case that

the available data is too sparse

to be able to support/reject the hypothesis of a given probability distribution - with a reasonable significance

Therefore it is necessary to use common sense i.e. :

First to consider physical reasons for selecting a given distribution

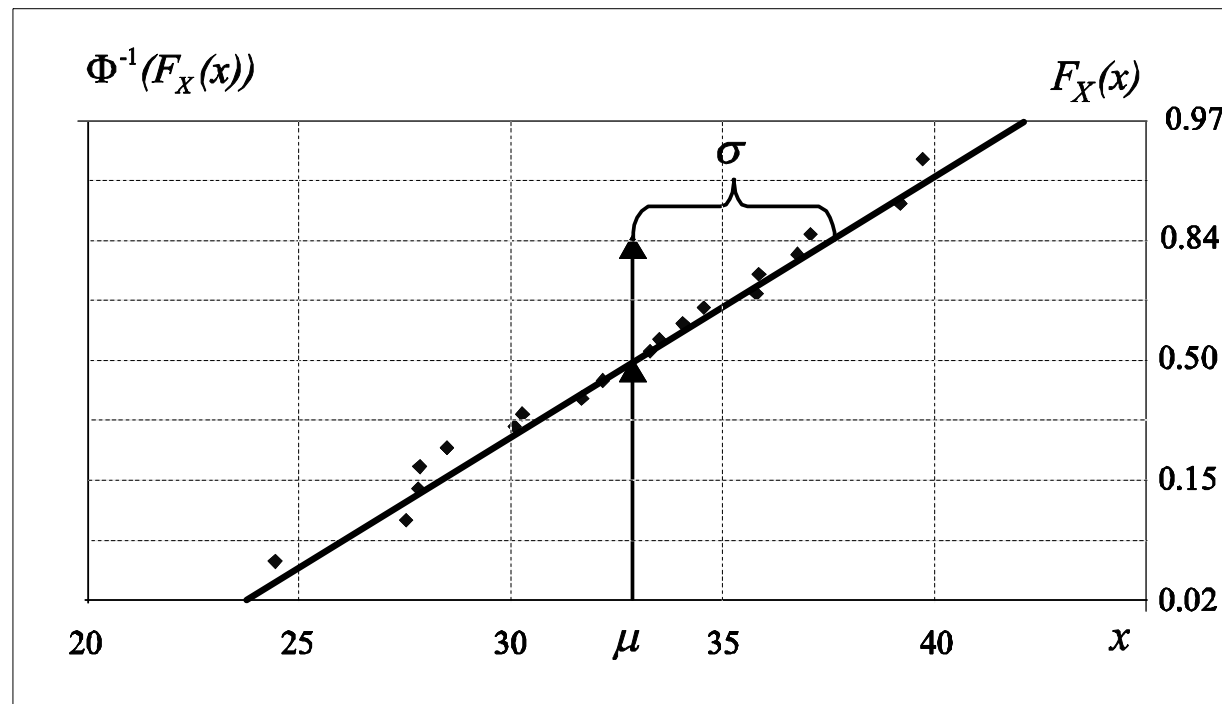
Thereafter to check if the available data are in gross contradiction with the selected distribution



# Estimation and Model Building

Model selection by use of probability paper

Plotting the sample probability distribution function in the probability paper yields



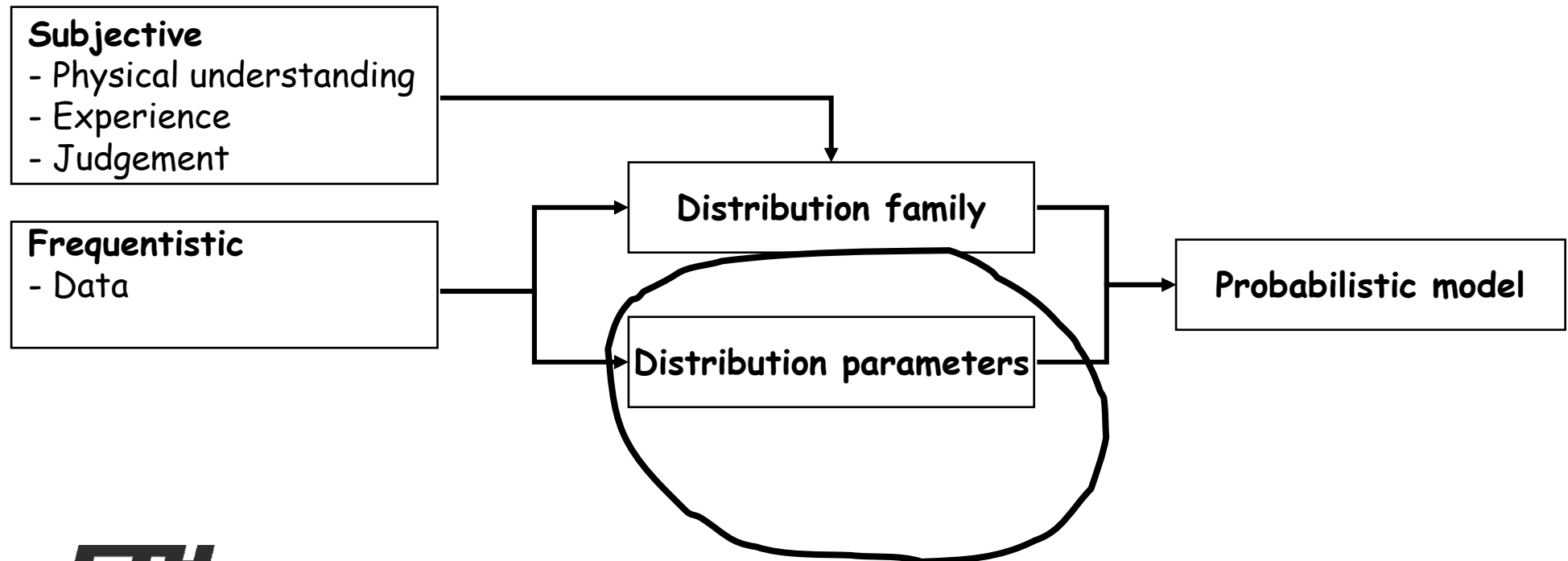




# Overview of Estimation and Model Building

Different types of information is used when developing engineering models

- subjective information
- frequentististic information





# Estimation of Distribution Parameters

We assume that we have identified a plausible family of probability distribution functions - as an example :

## Normal Distribution

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$$

## Weibull distribution

$$f_X(x) = \frac{k}{u-\varepsilon} \left(\frac{x-\varepsilon}{u-\varepsilon}\right)^{k-1} \exp\left(-\left(\frac{x-\varepsilon}{u-\varepsilon}\right)^k\right)$$

and thus now need to determine - estimate - its parameters

$$\theta = (\theta_1, \theta_2, \dots, \theta_k)^T$$



# Estimation of Distribution Parameters

The method of moments (MoM)

To start with we assume that we have data on the basis of which we can estimate the distribution parameters  $\hat{\mathbf{x}} = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n)^T$

The idea behind the method of moments is to determine the distribution parameters such that the sample moments (from the data) and the analytical moments (from the assumed distribution) are identical.

$$m_j = \frac{1}{n} \sum_{i=1}^n x_i^j$$

Sample moments

$$\begin{aligned} \lambda_j &= \int_{-\infty}^{\infty} x^j \cdot f_X(x|\boldsymbol{\theta}) dx \\ &= \lambda_j(\theta_1, \theta_2, \dots, \theta_k) \end{aligned}$$

Analytical moments



# Estimation of Distribution Parameters

## The Maximum Likelihood Method (MLM)

The idea behind the method of maximum likelihood is that the parameters are determined such that the likelihood of the observations is maximized

The likelihood can be understood as the probability of occurrence of the observed data conditional on the model

The Maximum Likelihood Method may seem to be more complicated than the MoM but has a number of attractive properties which we shall see later



# Estimation of Distribution Parameters

## Summary

Method of Moments provides point estimates of the parameters

- No information about the uncertainty with which the parameter estimates are associated.

Maximum Likelihood Method provides point estimates of the estimated parameters

- Full distribution information - normal distributed parameters, mean values and covariance matrix.