# Basics of Statistics

## Lecture 1b

# Probability

How likely it is that some event will happen?

**Idea**:

- Experiment
- Outcomes (sample points) $O_1, O_2, \ldots O_n$
- Sample space $\Omega$
- Event A
- Probability function P: Events $\rightarrow [0,1]$

# Probability

Example: Tossing a coin two times

Example:

- $p(A)$ frequency of observing A
- $p(A, B)$ frequency of observing A and B
- $p(B|A)$ frequency of observing B given A

# Properties and definitions

- One can think of events as sets
  - Set operations are defined: $A \cup B, A \cap B, \bar{A} \backslash B$
- $P(A \cup B) = P(A) + P(B)$ if $A \cap B = \emptyset$

- Independence $P(A, B) \equiv P(A \cap B) = P(A)P(B)$

- Conditional probability $P(A|B) = \dfrac{P(A,B)}{P(B)}$

# Bayes theorem

Example:

- We have constructed spam filter that
  - identifies spam mail as spam with probability 0.95
  - Identifies usual mail as spam with probability 0.005
- This kind of spam occurs once in 100,000 mails
- If we found that a letter is a spam, what is the probability that it is actually a spam?

# Bayes theorem

- We have some knowledge about event B
  - Prior probability P(B) of B
- We get new information A
  - P(A)
  - P(A|B) probability of A can occur given B has occured
- New (updated) knowledge about B
  - Posterior probability P(B|A)

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

# Random variables

- Instead of having outcomes, we can have a variable X:

  – Outcomes➔$\mathbb{R}$  Continuous random variables
  – Outcomes➔$\mathbb{N}$ Discrete random variables
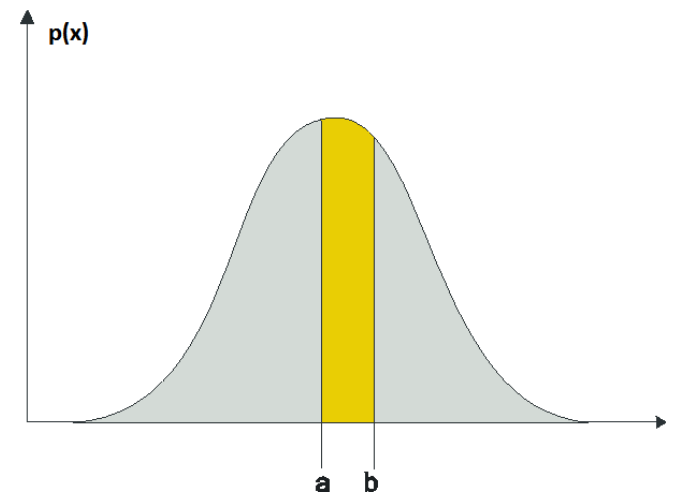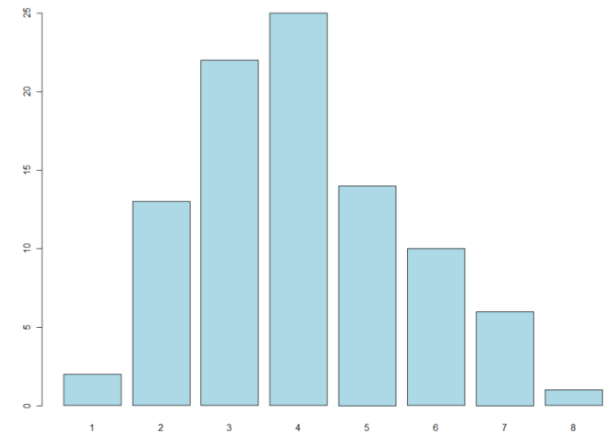
## Examples:

- X={amount of times the word "crisis" can be found in  financial documents}
  –   P(X=3)
- X={Time to download a specific file to a specific computer}
  –   P(X=0.36 min)

# Distributions

- Discrete
  - Probability mass function P(x) for all feasible x

- Coninuous
  - Probability density function p(x)
    - $p(x \in [a,b]) = \int_a^b p(x)dx$
    - $p(x) \geq 0, \int_{-\infty}^{+\infty} p(x)dx = 1$
  - Cumulative distribution function $F(x) = \int_0^x p(t)dt$

# Expected value and variance

- Expected value = mean value
  - $E(X) = \sum_{i=1}^{n} X_i P(X_i)$
  - $E(X) = \int X p(X) dX$

- Variance how much values of random variable can deviate from mean value

  - $Var(X) = E\left(X - E(X)\right)^2 = E(X^2) - E(X)^2$

# Probabilities

- **Laws of probabilities**
  - Sum rule (compute **marginal** probability)

$$p(X) = \sum_Y p(X, Y)$$

$$p(X) = \int p(X, Y) \, dY$$

  - Product rule

$$p(X, Y) = p(X|Y)p(Y)$$

Combination 1:

$$p(X) = \sum_Y p(X|Y)p(Y)$$

$$p(X) = \int p(X|Y)p(Y) \, dY$$
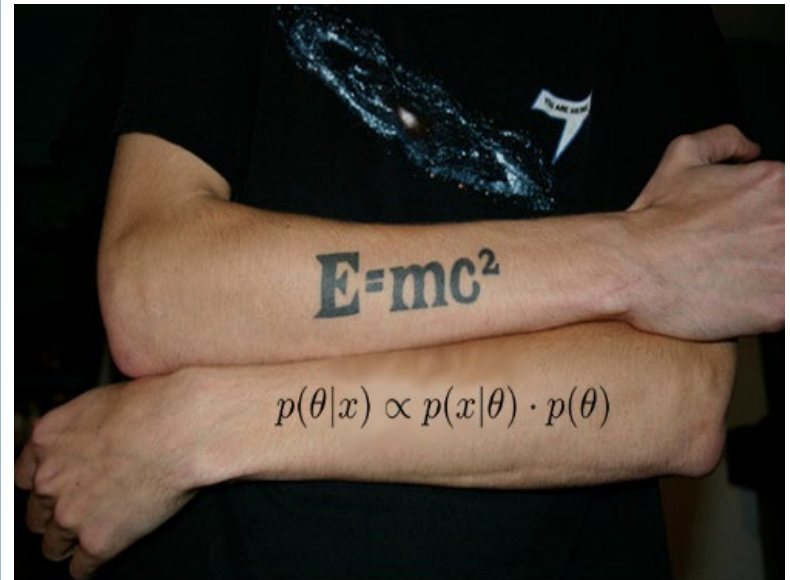
# Bayes theorem

For random variables:

**Bayes Theorem**

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

$$p(Y|X) \propto p(X|Y)p(Y)$$

$$p(Y|X) = \frac{p(X|Y)p(Y)}{\int p(X|Y)p(Y)dY}$$



$$E=mc^2$$

$$p(\theta|x) \propto p(x|\theta) \cdot p(\theta)$$

# Some conventional distributions

Bernoulli distribution

- Events: Success (X=1) and Failure (X=0)
- P(X=1)=p, P(X=0)=1-p

- $E(X) = p$
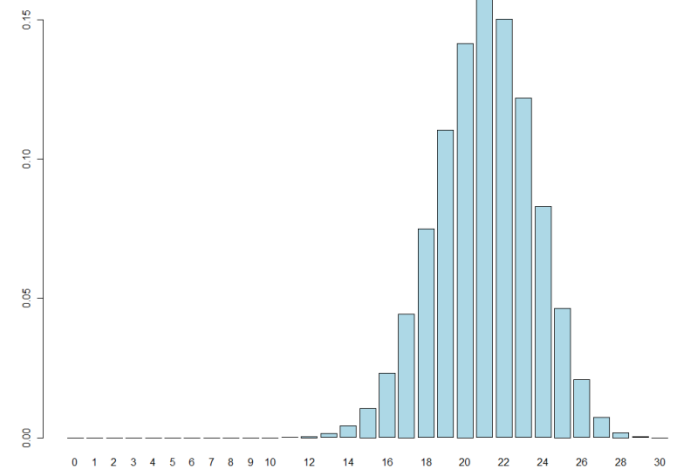- $Var(X) = p(1 - p)$

Examples: Tossing coin, vinning a lottery,..

# Some conventional distributions

## Binomial distribution

- Sequence of *n* Bernoulli events

- X={Amount of successes among these events}, X=0,…,n

$$P(X = r) = \frac{n!}{(n-r)!\, r!} p^r (1-p)^{n-r}$$

- $EX = np$

- $Var(X) = np(1-p)$
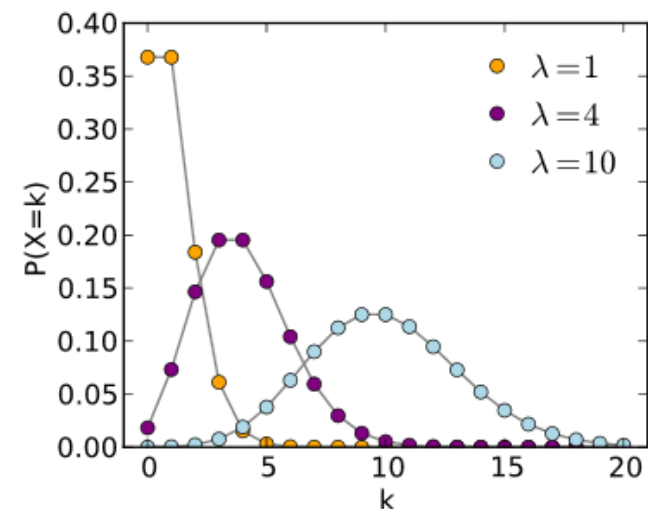
# Poisson distribution

- Customers of a bank **n** (in theory, endless population)
- Probability that a specific person will make a call to the bank between 13.00 and 14.00 a certain day is **_p_**
  - **_p_** can be very small if population is large (rare event)
  - Still, some people will make calls between 13.00 and 14.00 that day, and their amount may be quite big
  - A known quantity **λ=_np_** is mean amount of persons that call between 13.00 and 14.00
  - **X**={amount of persons that have called between 13.00 and 14.00}

# Poisson distribution

- $P(X = r) = \lim\limits_{n \to \infty} \dfrac{n!}{(n-r)!\,r!} p^r (1-p)^{n-r}$

- It can be shown that

$$P(X = r) = \frac{\lambda^r e^{-\lambda}}{r!}$$

- $E(X) = \lambda$

- $Var(X) = \lambda$

# Poisson distribution

- Further properties:
  - Poisson distribution is a good approximation of the binomial distribution if n >20 and $p < 0.05$
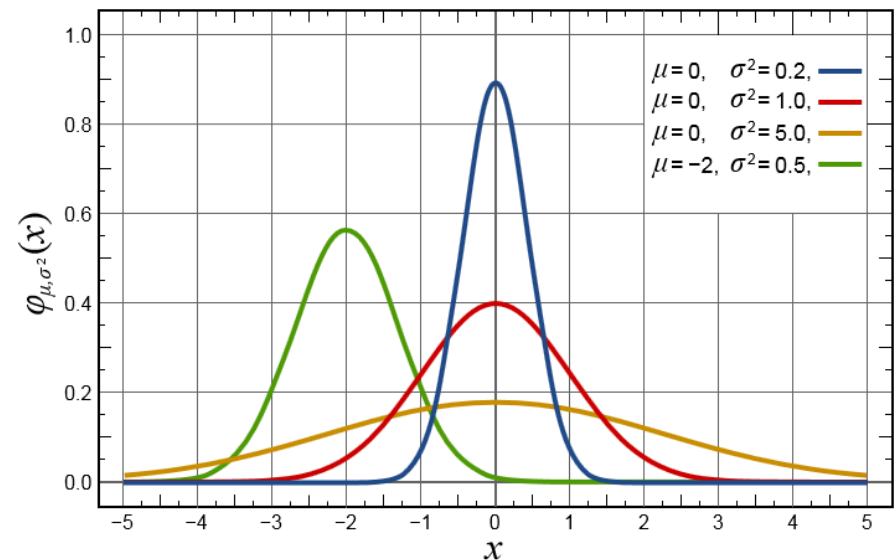  - Excellent approximation if $n \geq 100$ and $np \leq 10$

# Normal distribution

- Appears in almost all applications
  - Difference between the times required to download two specific documents to a specific computer
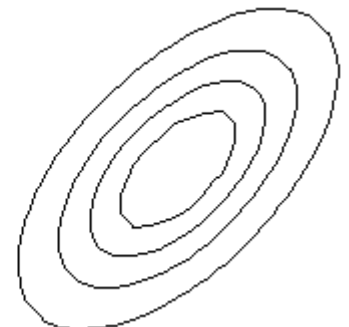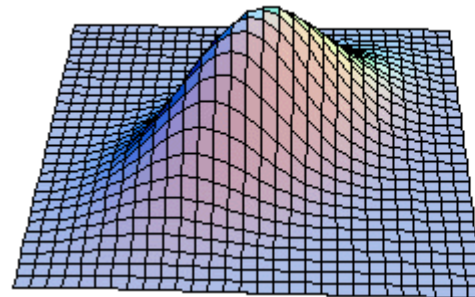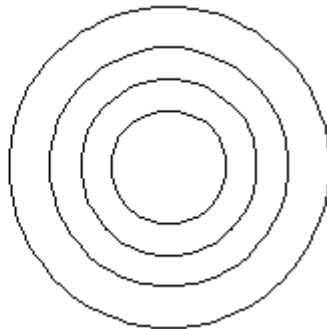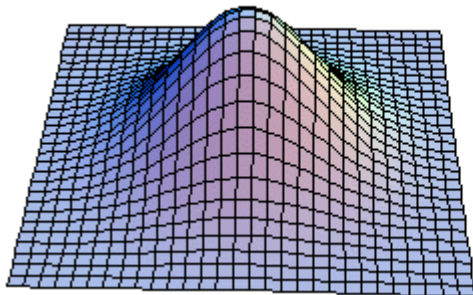
$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \sigma > 0$$

- $E(X) = \mu$
- $Var(X) = \sigma^2$

# Multivariate distributions

- Probability of two variables having certain values at the same time
  - P.D.F. p(x,y)
  - Correlation

# Basic ML ingridients

- Data $T$: observations
  - Features $x_1, .. x_p$
  - Targets $y_1, ..., y_r$

| Case | $x_1$ | $x_2$ | $y$ |
|------|-------|-------|-----|
| 1 | | | |
| 2 | | | |
| ... | | | |

- Mathematical Model $P(x|w_1, ... w_k)$ or $P(y|x, w_1, ... w_k)$
  - Example: Linear regression $p(y|x, w) = N(w_0 + w_1 x, \sigma^2)$

- Learning algorithm (data→get parameters $\hat{w}$ or $p(w|T)$ )
  - Maximum likelihood, Bayesian estimation
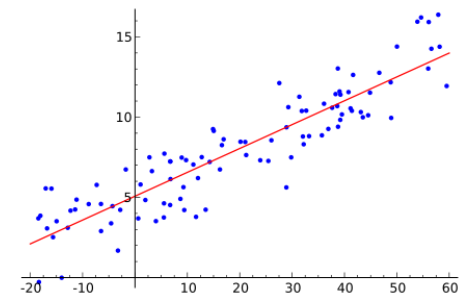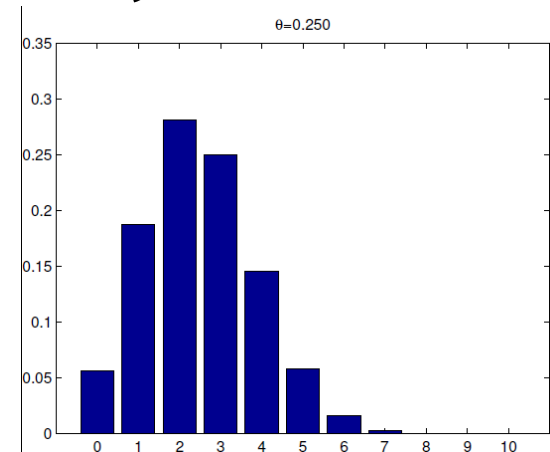
- Predict  new data $x_*$ by using the fitted model

# Probabilistic models

- A distribution $p(x|w)$ or $p(y|x,w)$

- Example:

  - $x \sim Bin(n, \theta)$

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

  - $y \sim N(\alpha_0 + \alpha_1 x, \sigma^2)$



Source: Wikipedia

Learn basic distributions and their properties

# Fitting a model

- Given dataset $T$ and model $p(x|w)$ or $p(y|x, w)$

  - Frequentist approach: which combination of parameter values fits my data best?

  - Bayesian approach: parameters are random variables, all feasible values are acceptable
    - Different parameter values have different probabilities

# Fitting a model

- Frequenist principle: **Maximum likelihood** principle

  - Compute likelihood $p(\boldsymbol{T}|\, w)$

$$p(\boldsymbol{T}|\, w) = \prod_{i=1}^{n} p(X_i|w)$$

$$p(\boldsymbol{T}|\, w) = \prod_{i=1}^{n} p(Y_i|X_i, w)$$



  - Maximize the likelihood and find the optimal $w^*$
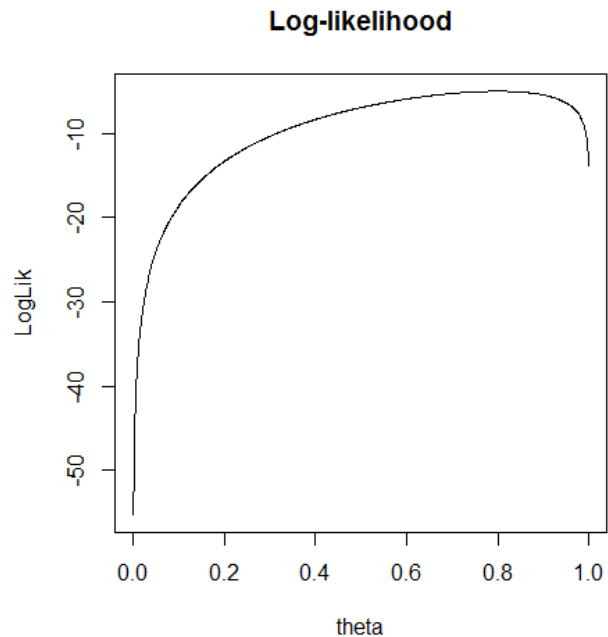
# Fitting a model

**Remarks:**

- Likelihood shows how much the chosen parameter value is proper for a specific model and the given data

- Normally **log-likelihood** is used in computations instead

- Other alternatives to ML exist…

# Fitting a model

Example: tossing a coin.

$$T = \{0,1,1,0,1,1,1,1,1,1,1,1\},$$
$$p(x = 1|\theta) = \theta, p(x = 0|\theta) = 1 - \theta$$

**Log-likelihood**

# Bayesian probabilities

- Probability reflects your knowledge (uncertainty) about a phenomenon → **subjective probabilities**
  - **Prior probability** $p(w)$, can be uninformative $p(w) \propto 1$
  - Formulate a model, compute **likelihood** $p(T|w)$
  - **Posterior probability** $p(w|T)$, after observing data
    - $p(w|T) \propto p(T|w)p(w)$

- Model parameters are considered as random variables
  - In real life, do not need to be random, but we model as random

# Fitting a model

- Bayesian principle
  - Compute $p(w|T)$ and then decide yourself what to do with this (for ex. MAP, mean, median)
- Use bayes theorem

$$p(w|T) = \frac{p(T|w)p(w)}{p(T)} \propto p(T|w)p(w)$$

- $p(T)$ is **marginal likelihood**
  - $p(T) = \int p(T|w)p(w)dw$ or
  - $p(T) = \sum_i p(T|w_i)p(w_i)$

Example: tossing a coin. Find $p(\theta|T)$, estimate posterior mean $\theta^*$

# Fitting a model

- How to chose the prior?
  - Expert knowledge about the phenomenon
  - Forcing a model to have a certain structure
    - Example: decision trees: prior prefers smaller trees

      http://en.wikipedia.org/wiki/Conjugate_prior
  - Conjugacy
    - Distribution of the posterior is the same type as the distribution of the likelihood or prior

- Prior is the most controversial about Bayesian methods, but
  - When $N \rightarrow \infty$, data overwhelms the prior

# Measuring uncertainty

- **Confidence interval** (frequentist)

- **Credible interval** (Bayes)

- **Prediction interval** (models)



- Example: Prediction interval for $Y \sim N(2x + 4, 1)$ at $x = 5$