

# Point estimations

# Agenda

- Maximum Likelihood Estimation (MLE)
- Method of Moments (MoM)
- Bayesian estimation
  - Maximum a posteriori estimation (MAP)
  - Posterior mean

# Maximum Likelihood Estimation

# Likelihood function

The likelihood of unknown parameter  $\theta$  given your data.

$X_1, \dots, X_n$  : random sample from  $f(x|\theta)$

$$L(\theta | X_1, \dots, X_n) = \frac{f(X_1, \dots, X_n | \theta)}{\text{joint probability/density}}$$

special case : if  $X_1, \dots, X_n$  <sup>iid</sup>  $f(x|\theta)$

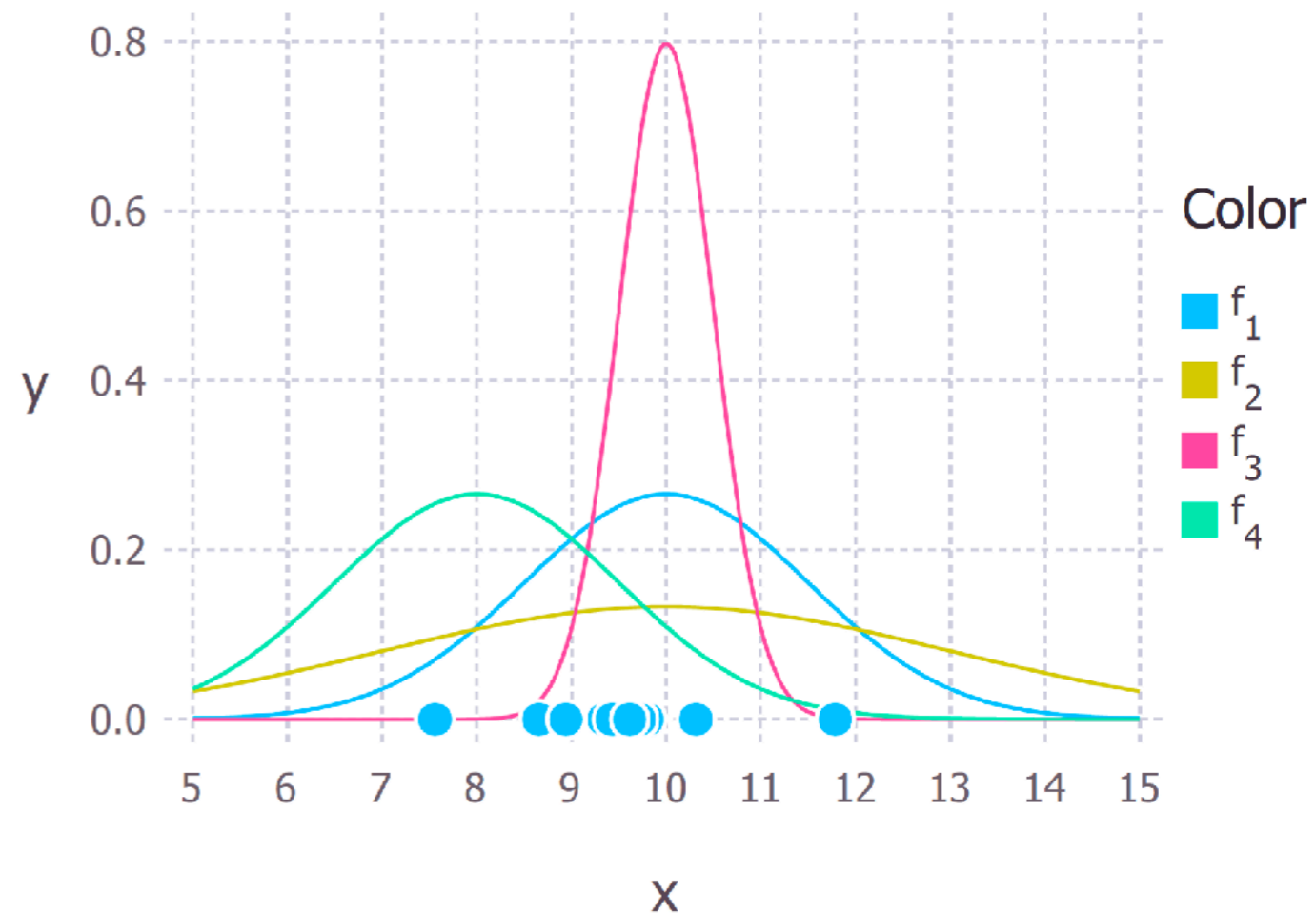
$$L(\theta | X_1, \dots, X_n) = \prod_{i=1}^n f(X_i | \theta) \quad (\text{if iid})$$

- The likelihood function is a function of  $\theta$
- It is **not** a probability density function
- It measures the “support” (i.e. likelihood) provided by the data for each possible value of the parameter.

# MLE

- Find the parameter that is “most likely” to observe your data.
- Maximizing the likelihood function is equivalent to maximizing the log-likelihood function (for computational issues)

$$Q(\theta | x_1, \dots, x_n) = \log L(\theta | x_1, \dots, x_n)$$



<https://towardsdatascience.com/probability-concepts-explained-maximum-likelihood-estimation-c7b4342fdbb1>

How to maximize a function?

# Example: coin tossing

Assume  $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$

$$L(p | x_1, \dots, x_n) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

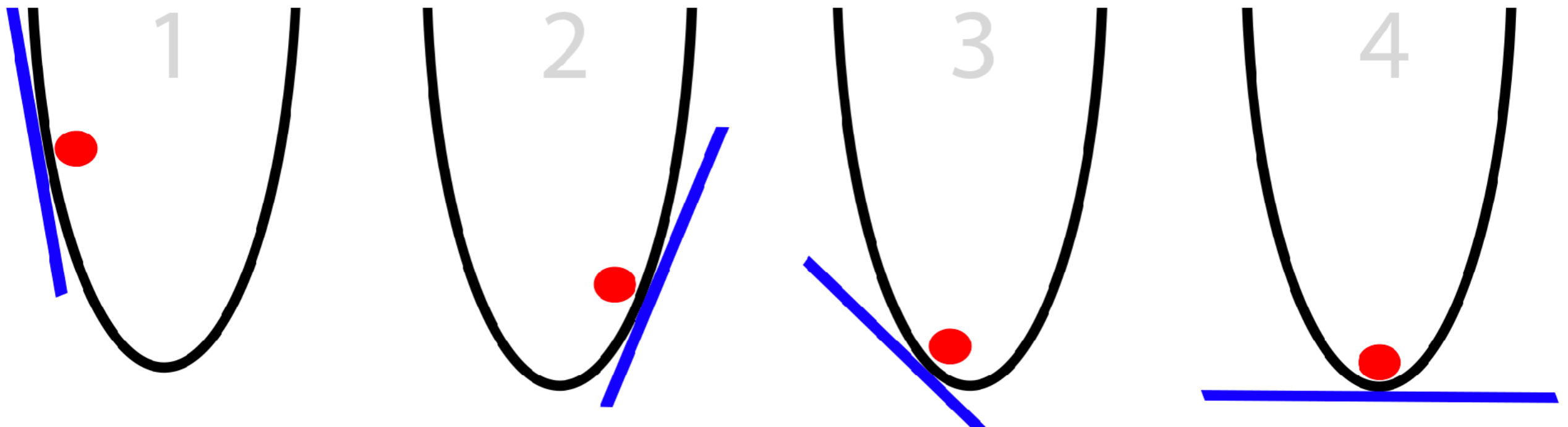
$$\log L(p | x_1, \dots, x_n) = \log L = \sum_{i=1}^n x_i \cdot \log p + (1-x_i) \log(1-p)$$

$$\frac{dL}{dp} = \sum_{i=1}^n \frac{x_i}{p} - \frac{1-x_i}{1-p} = \dots = \frac{\sum_{i=1}^n x_i - np}{p(1-p)} \stackrel{\Delta}{=} 0$$

$$\Rightarrow \hat{p} = \frac{1}{n} \sum_{i=1}^n x_i \stackrel{\Delta}{=} \bar{x}$$

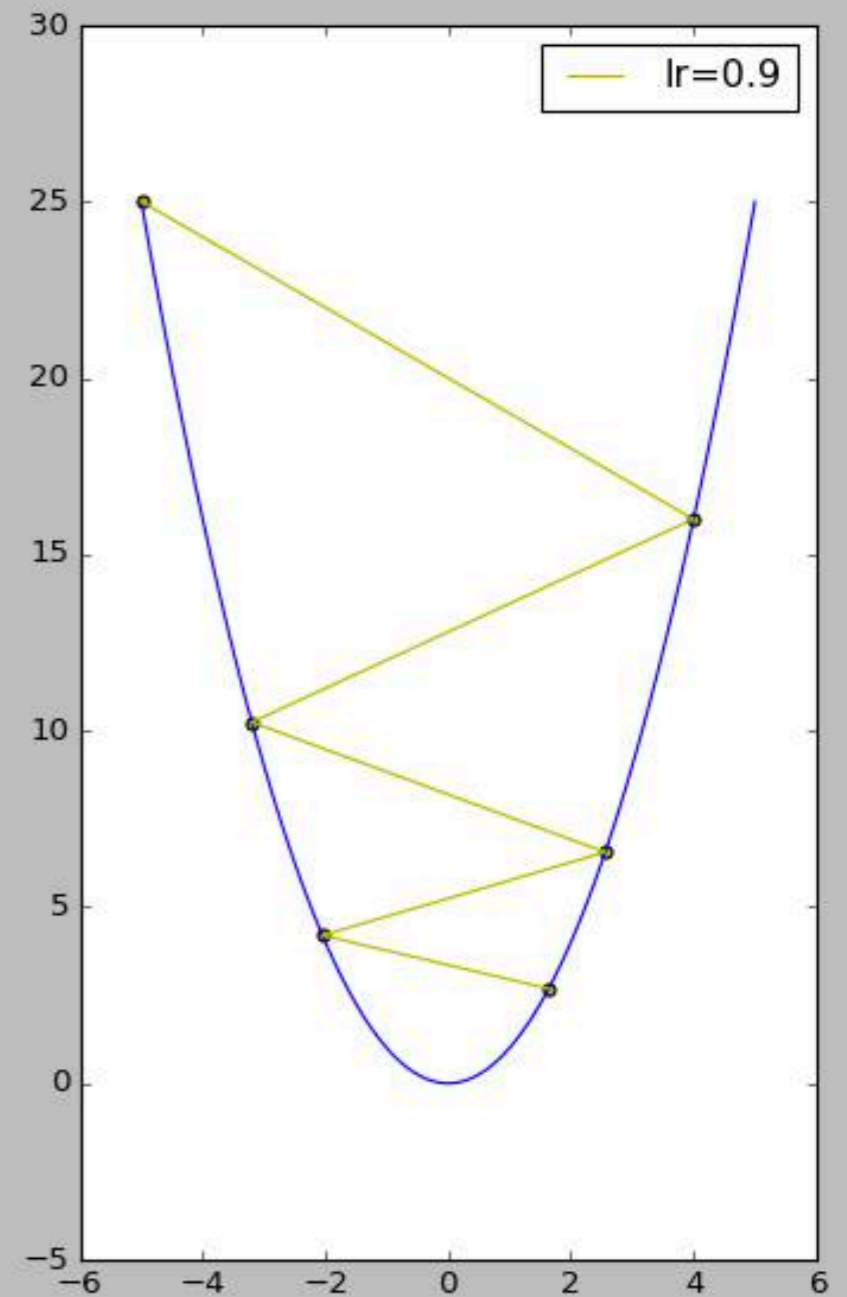
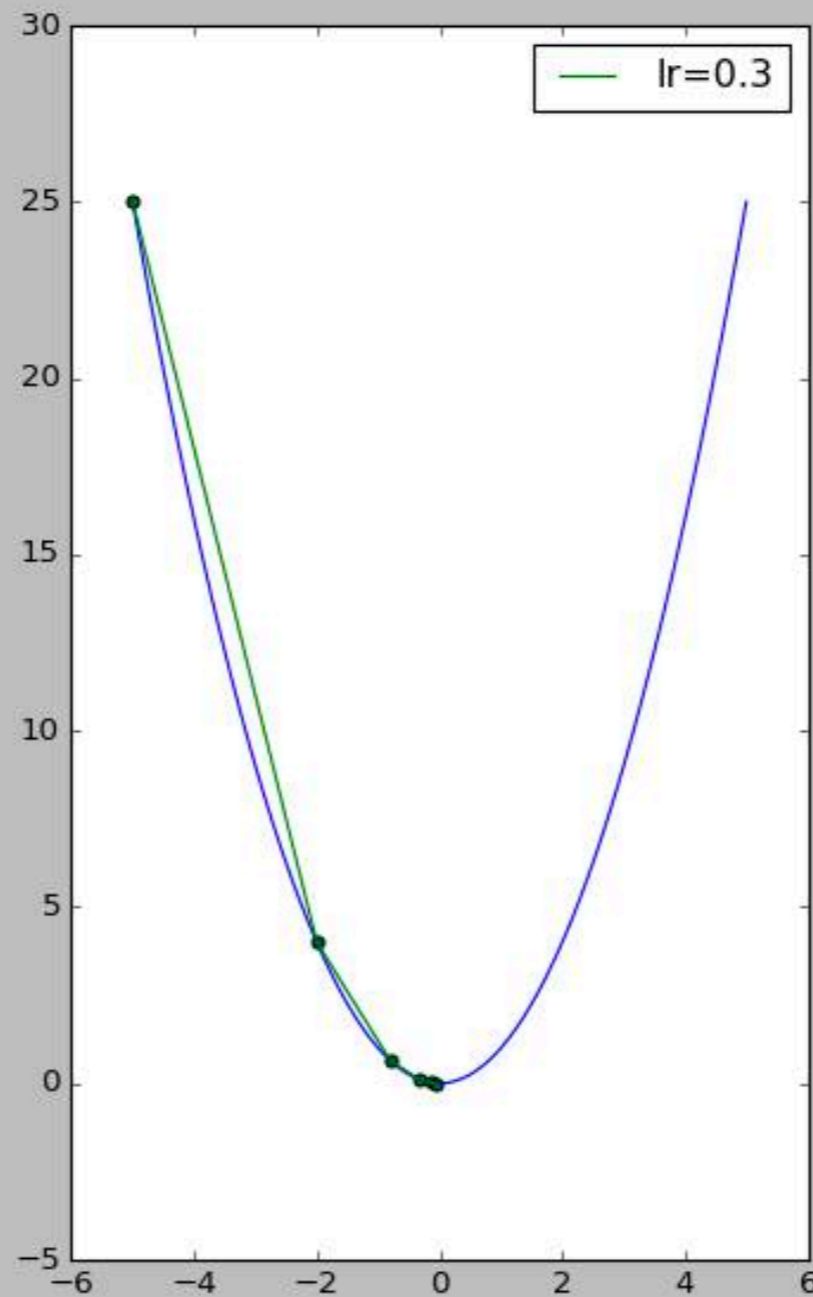
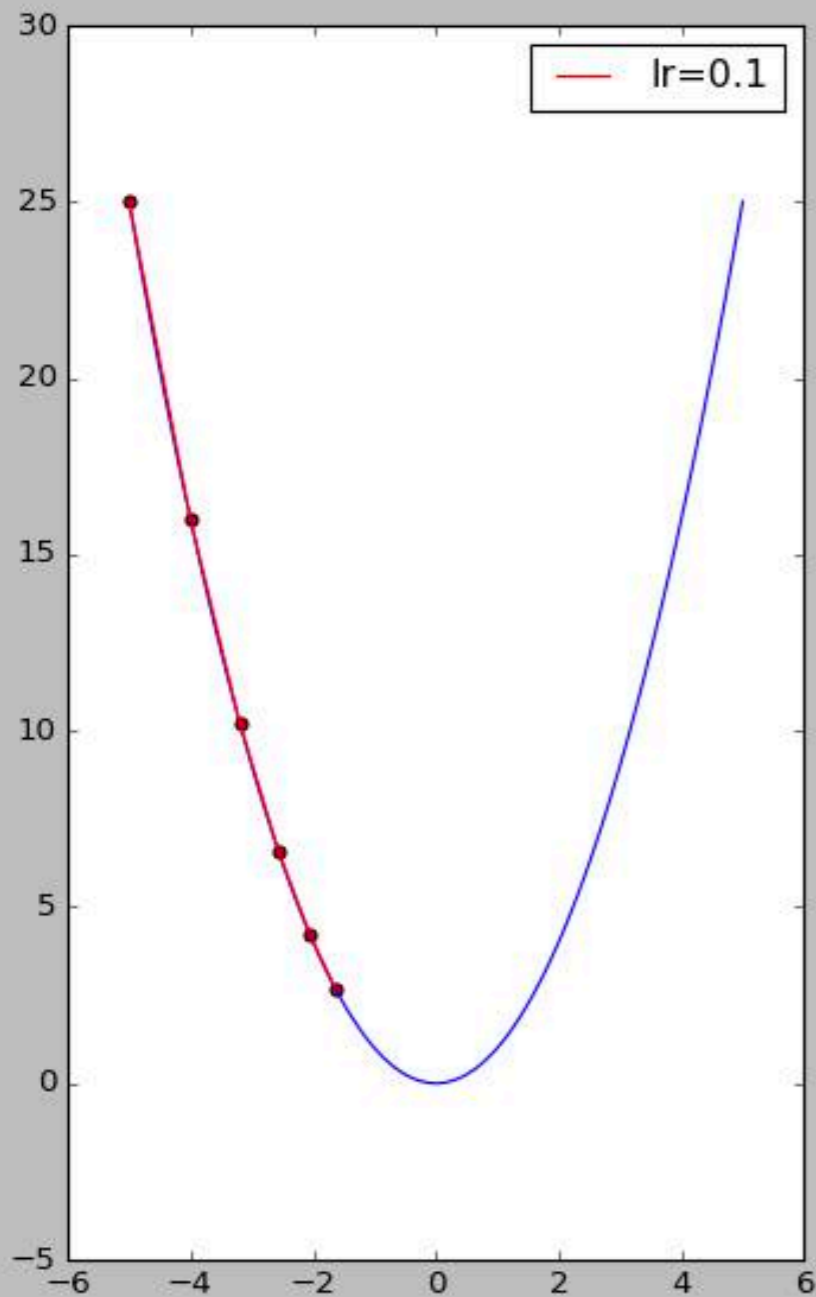


# Gradient descent



<https://iamtrask.github.io/2015/07/27/python-network-part2/>

# Learning rate



# Local optimums



<https://iamtrask.github.io/2015/07/27/python-network-part2/>

# Method of Moments

# Expectation

- If  $X$  is a discrete random variable with p.m.f.  $f_X(x)$ ,

$$E[X] = \sum_x xf_X(x)$$

- If  $X$  is a continuous random variable with p.d.f.  $f_X(x)$ ,

$$E[X] = \int_x xf_X(x)$$

- $E[X^k]$  are called moments with  $k \geq 1$

# Law of large number

- If  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_X(x)$ , then  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  converges to  $E[X_i]$  (in probability) as  $n \rightarrow \infty$
- LLN holds for every moments

# Method of moments

- The values of moments depend on the values of unknown parameters  $\theta$  (moment conditions)
- By LLN, we can estimate moments by sample means

# Example: normal distribution

Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

- moment conditions:

$$\frac{1}{n} \sum_{i=1}^n X_i \approx E[X_1] = \mu$$

$$\frac{1}{n} \sum_{i=1}^n X_i^2 \approx E[X_1^2] = \mu^2 + \sigma^2$$

- By solving the above moment conditions we have

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i, \text{ and } \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \hat{\mu}^2$$



# Example: linear regression

Let  $Y_i \stackrel{iid}{\sim} N(\mu(\mathbf{x}_i), \sigma^2)$  with  $\mathbf{x}_i = [x_{i1}, \dots, x_{ip}]^T \in \mathbb{R}^p$  and

$$\mu(\mathbf{x}_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

- Moment conditions:

$$E[Y - \mu(\mathbf{x})] = 0$$

$$E\left[x_j (Y - \mu(\mathbf{x}))\right] = 0$$

$$E\left[(Y - \mu(\mathbf{x}))^2\right] = \sigma^2$$

- Plug the sample moments into the moment conditions, we obtain

$$\frac{1}{n} \sum_{i=1}^n \left( Y_i - \left( \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} \right) \right) \approx 0$$

$$\frac{1}{n} \sum_{i=1}^n x_{ij} \left( Y_i - \left( \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} \right) \right) \approx 0$$

$$\frac{1}{n} \sum_{i=1}^n \left( Y_i - \left( \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} \right) \right)^2 \approx \sigma^2$$

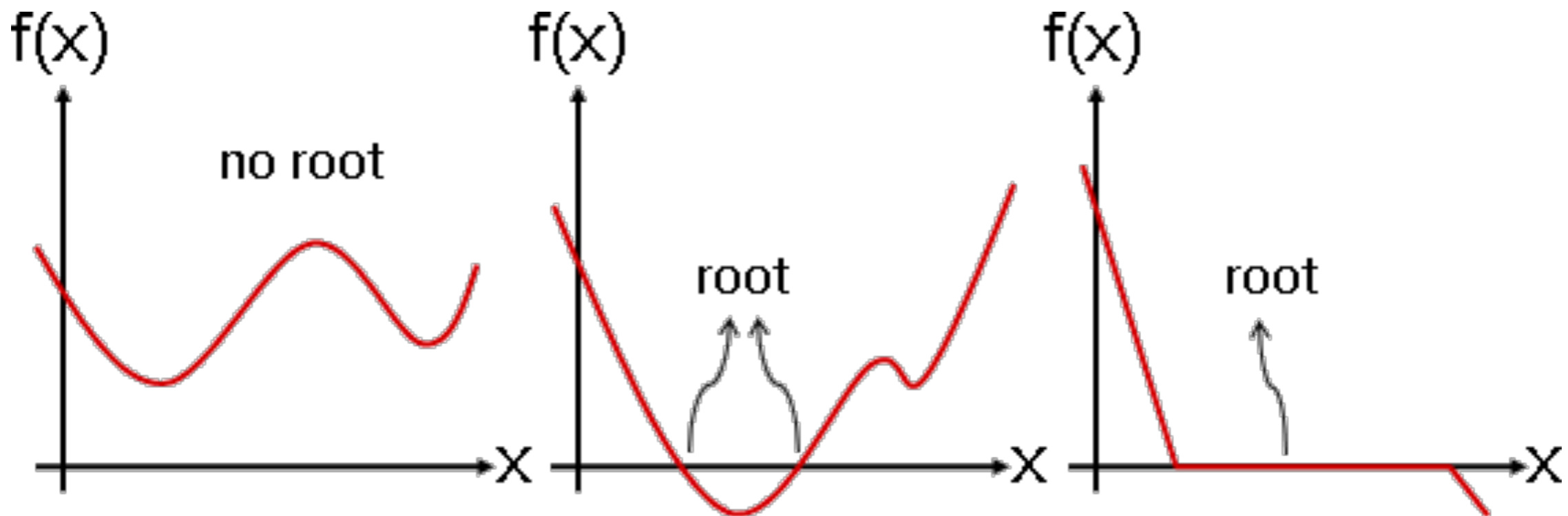
# Solving linear systems

The solution of a linear system  $\mathbf{Ax} = \mathbf{b}$  can be found (if it exists) by

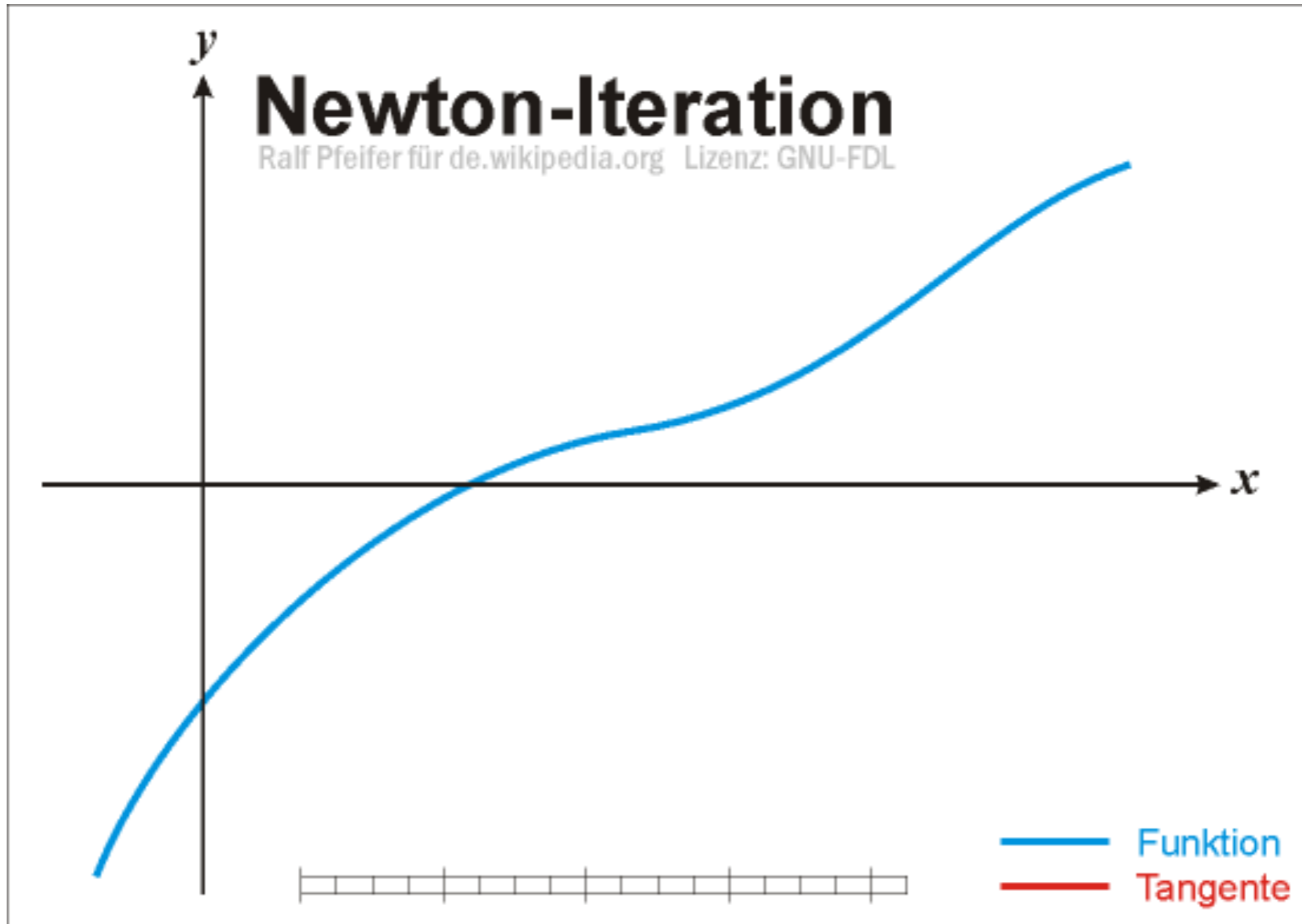
- Gaussian elimination (`numpy.linalg.solve`)
- Minimize  $\|\mathbf{Ax} - \mathbf{b}\|^2$  (`numpy.linalg.lstsq`)

# Solving nonlinear systems

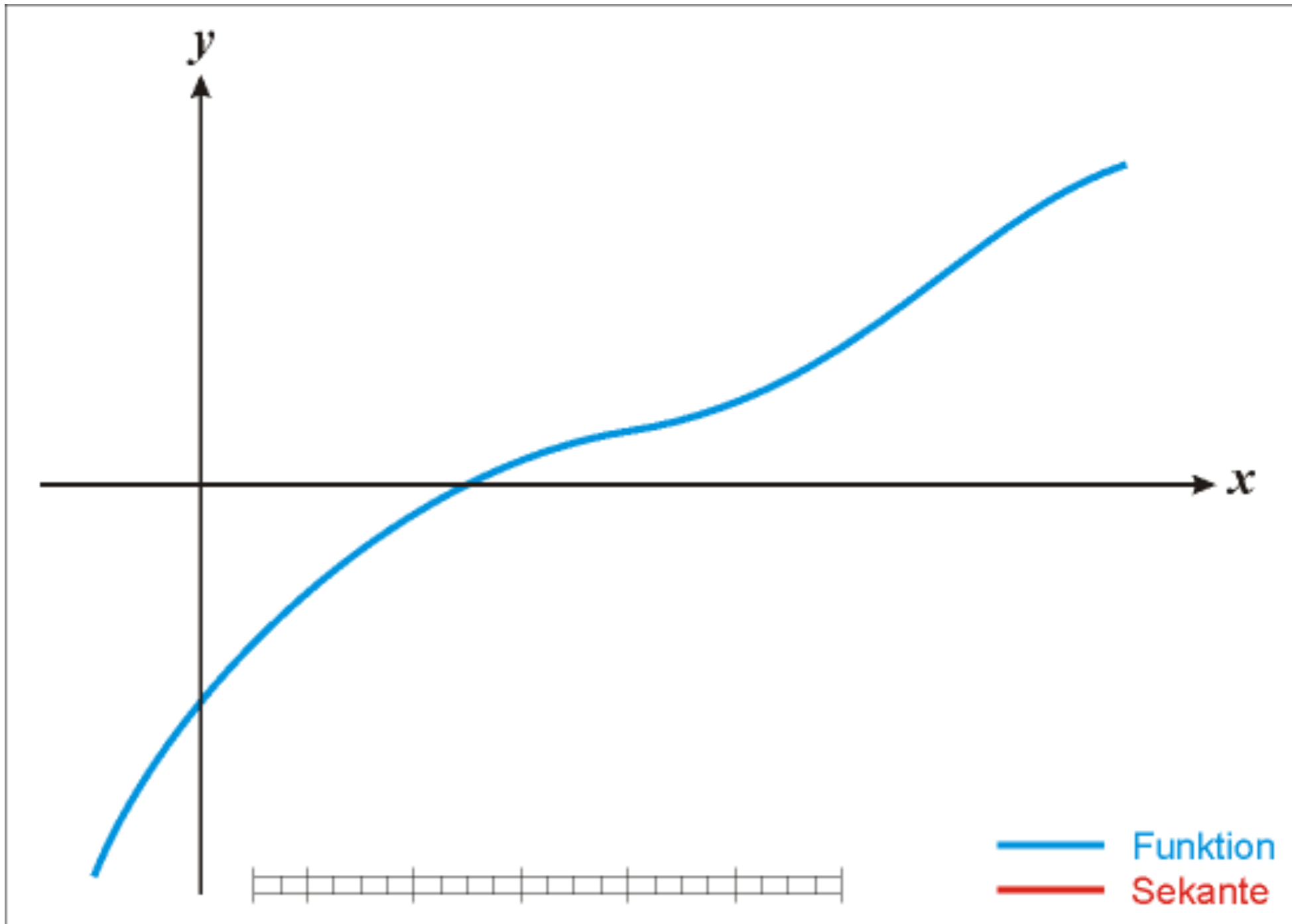
The solution of a nonlinear system  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$  can be found (if it exists) by various root-finding algorithms (scipy.optimize.root)



# Newton's method



# Secant method



# Pros

- Easy to compute and always work
- MoM is consistent; i.e.  $\hat{\theta}_{MoM} \rightarrow \theta$  as  $n \rightarrow \infty$

# Cons

- MoM **may not be unique**: different moment conditions yields to different results!
- Not the most efficient (i.e. achieving minimum mean squared error, MSE) estimators
- Sometimes MoM may be meaningless



# MoM may be meaningless

- Suppose we observe 3, 5, 6, 18 from  $U(0,\theta)$
- Since  $E[X] = \theta/2$  for  $X \sim U(0,\theta)$ , the MoM of  $\theta$  is

$$\hat{\theta}_{MoM} = 2\bar{X} = 2 \times \frac{3 + 5 + 6 + 18}{4} = 16$$

- This estimation is not acceptable since we have already observed a 18.

# Extensions

- Generalized MoM
  - Generalized moment conditions
  - The number of conditions may exceed the number of parameters
- Method of simulated moments
  - Approximate the theoretical moments when they are not available

# Bayesian estimation

# Key concepts

- Since  $\hat{\theta}$  (derived from some random sample) is random, we can treat  $\theta$  as random and specify its probability distribution as  $\pi(\theta)$
- The probability  $\pi(\theta)$  is usually specified by a data scientist to express one's beliefs. Thus, we call it a “prior”.

# Posterior

- By Bayes' theorem, we can derive the “posterior” distribution of  $\theta$ :

$$\begin{aligned} p(\theta | X_1, \dots, X_n) &= \frac{f(X_1, \dots, X_n | \theta)\pi(\theta)}{f(X_1, \dots, X_n)} \\ &= \frac{f(X_1, \dots, X_n | \theta)\pi(\theta)}{\int f(X_1, \dots, X_n | \theta)\pi(\theta)d\theta} \\ &\propto f(X_1, \dots, X_n | \theta)\pi(\theta) \end{aligned}$$

# Maximum a posteriori estimation

- The posterior distribution can be interpreted as the conditional probability of  $\theta$  given observational data
- Thus, similar to MLE, we may find the mode of the posterior since it is the most likely value of  $\theta$

$$\hat{\theta}_{MAP} = \arg \max f(X_1, \dots, X_n | \theta) \pi(\theta)$$

# Posterior mean

- The “mean” of the posterior is another frequently used Bayesian estimator,

$$\hat{\theta} = E[\theta | X_1, \dots, X_n] = \int \theta p(\theta | X_1, \dots, X_n) d\theta$$

- The expectation is often approximated by Markov chain Monte Carlo (MCMC) method since the above integration is usually difficult

# Example: normal distribution

Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$  and we assume that  $\mu \sim N(\mu_0, \tau^2)$ . Then

$$p(\mu | X_1, \dots, X_n) \propto \frac{1}{\sqrt{2\pi\tau}} e^{-\frac{1}{2}\left(\frac{\mu - \mu_0}{\tau}\right)^2} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2}\left(\frac{X_i - \mu}{\sigma}\right)^2}$$

and

$$\hat{\mu} = \frac{n\tau^2}{n\tau^2 + \sigma^2} \bar{X} + \frac{1}{n\tau^2 + \sigma^2} \mu_0$$



# (My unfair) suggestions

- Use Bayesian estimations when you have a domain expert; otherwise, use MLE
- Use MoM only for computational issues
  - The posterior (or likelihood function) is not convex
  - Big data

# Homework: logistic regression

- Breast Cancer Wisconsin (Diagnostic) Data Set (also available in scikit-learn)

- Assume that  $Y_i \in \{0,1\} \stackrel{iid}{\sim} \text{Bernoulli}(p(\mathbf{x}_i))$  with

$$p(x_i) = \frac{\exp \left[ \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} \right]}{1 + \exp \left[ \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} \right]}$$

Estimate the unknown coefficients  $\theta = [\beta_0, \beta_1, \dots, \beta_p]'$  by either MLE or MoM. Compare your results with the the ones provided by scikit-learn (example) with very large  $C$ .

- Moment conditions:

$$E [Y - p(\mathbf{x})] = 0$$

$$E [x_j (Y - p(\mathbf{x}))] = 0$$

- Plug the sample moments into the moment conditions, we obtain

$$\frac{1}{n} \sum_{i=1}^n (Y_i - p(\mathbf{x}_i)) = 0$$

$$\frac{1}{n} \sum_{i=1}^n x_{ij} (Y_i - p(\mathbf{x}_i)) = 0$$

# Readings

- Chapters 10.1–10.3 and 12.1–12.2 of “All of statistics”