# Learning theory: Maximum Likelihood, Bayesian Learning, Model Selection

Shikui Tu

Shanghai Jiao Tong University

2021-03-30
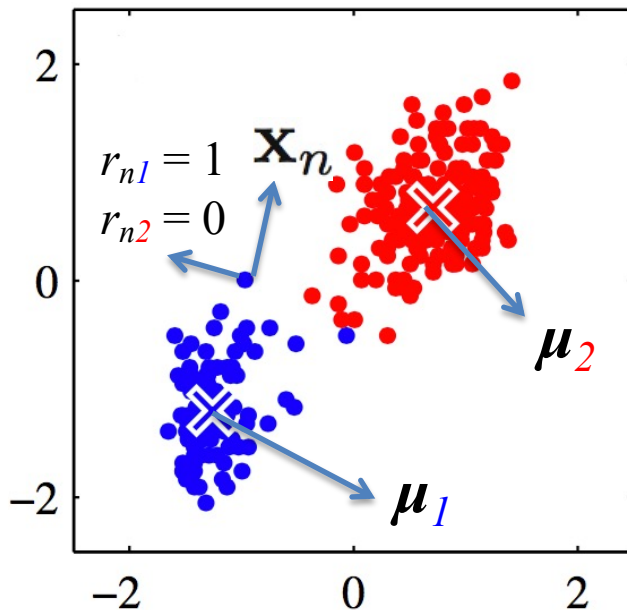
# Outline

- Recall from the previous lectures

- Maximum Likelihood (ML) learning

- Bayesian learning, Maximum A Posterior (MAP)

- Model Selection

# From minimizing sum of square distances to finding maximum likelihood

minimize

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \| \mathbf{x}_n - \boldsymbol{\mu}_k \|^2$$
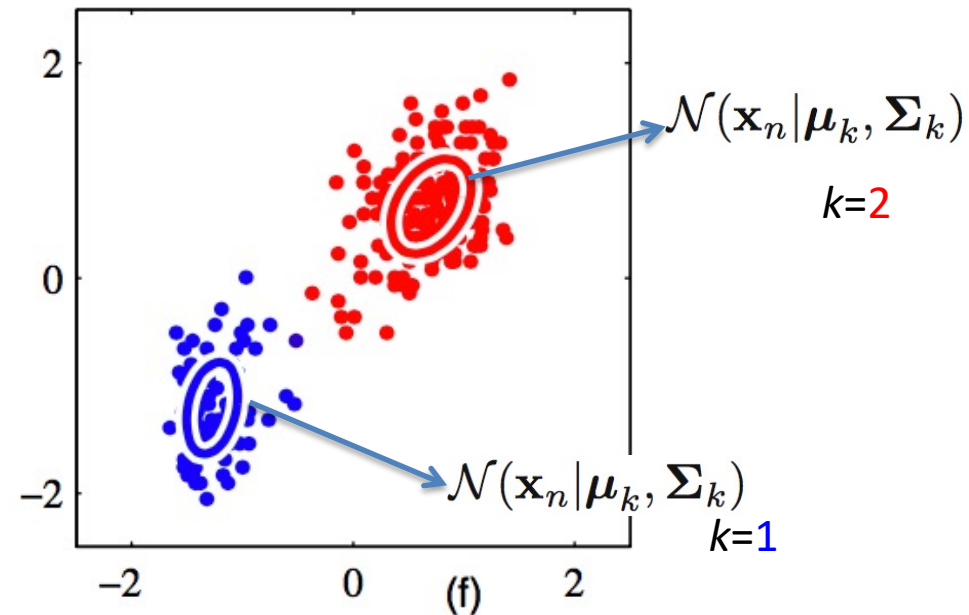
maximize likelihood

$$p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})$$

$$X = \{x_1, ..., x_N\}$$
$$\pi = \{\pi_1, ..., \pi_K\}$$
$$\mu = \{\mu_1, ..., \mu_K\}$$
$$\Sigma = \{\Sigma_1, ..., \Sigma_K\}$$

$r_{n1} = 1$
$r_{n2} = 0$

$\mathbf{x}_n$

$\boldsymbol{\mu}_2$

$\boldsymbol{\mu}_1$

$\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

$k=2$

$\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

$k=1$

(f)

Remember: **The closer the distance, the more likely the probability.**

# The General EM Algorithm

Given a joint distribution $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ over observed variables $\mathbf{X}$ and latent variables $\mathbf{Z}$, governed by parameters $\boldsymbol{\theta}$, the goal is to maximize the likelihood function $p(\mathbf{X}|\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$.

1. Choose an initial setting for the parameters $\boldsymbol{\theta}^{\text{old}}$.

2. **E step** Evaluate $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$.

3. **M step** Evaluate $\boldsymbol{\theta}^{\text{new}}$ given by

$$\boldsymbol{\theta}^{\text{new}} = \arg\max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$$

where

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}).$$

4. Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let

$$\boldsymbol{\theta}^{\text{old}} \leftarrow \boldsymbol{\theta}^{\text{new}}$$

and return to step 2.

# Summary for the EM algorithm for GMM

- Does it find the global optimum?
  - No, like K-means, EM only finds the nearest local optimum and the optimum depends on the initialization

- GMM is more general then K-means by considering mixing weights, covariance matrices, and soft assignments.

- Like K-means, it does not tell you the best K.

# Outline

- Recall from the previous lectures


- **Maximum Likelihood (ML) learning**


- Bayesian learning, Maximum A Posterior (MAP)


- Model Selection

# An example

- If flipping a coin a few times, and get



- What is the probability it will fall with the head up?

You may say:    3/5

Because ….

# Bernoulli distribution

The dataset $D = \{ x_t \}$, $t=1,...,N,$ $x_t \in \{H, T\}$

$D =$ 

$$P(x = Head) = \theta$$

$$P(x = Tail) = 1 - \theta$$

Flipping coins are **i.i.d.**, i.e., independent identically distributed according to Bernoulli distribution

**Question**: What is the parameter $\theta$ that maximizes the probability of observed data?

# Maximum Likelihood Estimation

- Choose parameter $\theta$ that <u>maximizes the probability of observed data</u>

$$\hat{\theta}_{MLE} = \arg\max_{\theta} \; P(D \mid \theta)$$

$$= \arg\max_{\theta} \prod_{i=1}^{n} P(X_i \mid \theta) \qquad \text{Independent draws}$$

$$= \arg\max_{\theta} \prod_{i:X_i=H} \theta \prod_{i:X_i=T} (1-\theta) \qquad \text{Identically distributed}$$

$$= \arg\max_{\theta} \underbrace{\theta^{\alpha_H}(1-\theta)^{\alpha_T}}_{J(\theta)}$$

$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T} \qquad \textbf{= 3/5} \quad \text{"Frequency of heads"}$$

Number of heads     Number of tails

9

# Outline

- Recall from the previous lectures

- Maximum Likelihood (ML) learning

- **Bayesian learning, Maximum A Posterior (MAP)**

- Model Selection

# Bayesian Learning

- Bayes rule

$$P(\Theta|X) = \frac{P(X|\Theta)P(\Theta)}{P(X)}$$

$P(X|\Theta)$: likelihood of data $X$ given parameter $\Theta$

$P(\Theta)$: prior distribution over the parameter $\Theta$

$P(X)$: marginal distribution of data $X$

- Prior distribution
  - Represents expert knowledge
  - Uninformative priors: Uniform distribution
  - Conjugate priors: Closed-form representation of posterior, $P(\theta)$ and $P(\theta|D)$ have the same form

**Bayes, Thomas (1763)** An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, **53:370-418**

11

# Bayesian learning

- Maximum A Posteriori (MAP)

$$\max_{\Theta} p(\Theta|X)$$

Equivalent to:

$$\log p(X, \Theta) = \log p(X|\Theta) + \log p(\Theta)$$

Consider a simple example:

$$p(x|\Theta) = G(x|\mu, \Sigma)$$

$$p(\mu) = G(\mu|\mu_0, \sigma_0^2)$$

# When is MAP the same as MLE?

- Maximum Likelihood estimation (MLE)

  Choose value that maximizes the probability of observed data

  $$\hat{\theta}_{MLE} = \arg\max_{\theta} P(D|\theta)$$

- Maximum *a posteriori* (MAP) estimation

  Choose value that is most probable given observed data and prior belief

  $$\hat{\theta}_{MAP} = \arg\max_{\theta} P(\theta|D)$$
  $$= \arg\max_{\theta} P(D|\theta)P(\theta)$$

# Bayesians vs Frequentists

# Thank you!