

## CS 2750 Machine Learning Lecture 4

### Density estimation

Milos Hauskrecht

[milos@pitt.edu](mailto:milos@pitt.edu)

5329 Sennott Square

---

### Density estimation

**Density estimation:** is an unsupervised learning problem

- **Goal:** Learn a model that represent the relations among attributes in the data

$$D = \{D_1, D_2, \dots, D_n\}$$

**Data:**  $D_i = \mathbf{x}_i$  a vector of attribute values

**Attributes:**

- modeled by random variables  $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$  with
  - **Continuous or discrete valued variables**

**Density estimation:** learn an underlying probability

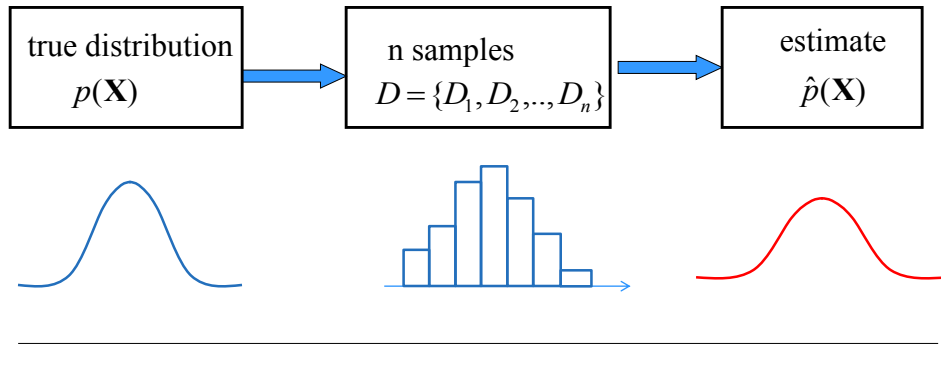
**distribution model :**  $p(\mathbf{X}) = p(X_1, X_2, \dots, X_d)$  from  $\mathbf{D}$

---

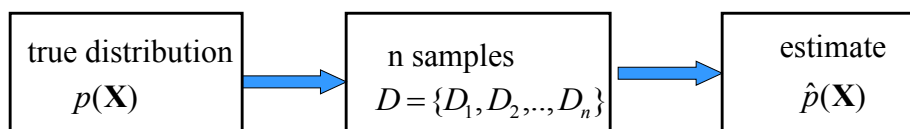
## Density estimation

**Data:**  $D = \{D_1, D_2, \dots, D_n\}$   
 $D_i = \mathbf{x}_i$  a vector of attribute values

**Objective:** estimate the model of the underlying probability distribution over variables  $\mathbf{X}$ ,  $p(\mathbf{X})$ , using examples in  $D$

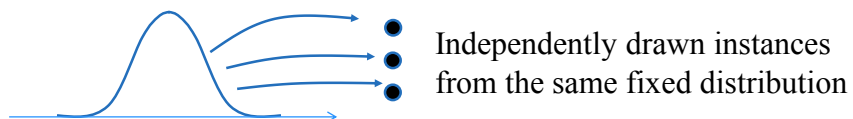


## Density estimation



**Standard (iid) assumptions: Samples**

- are **independent** of each other
- come from the same **(identical) distribution** (fixed  $p(\mathbf{X})$ )



## Density estimation

**Types of density estimation:**

### Parametric

- the distribution is modeled using a set of parameters  $\Theta$   
$$\hat{p}(\mathbf{X}) = p(\mathbf{X} | \Theta)$$
- **Example:** mean and covariances of a multivariate normal
- **Estimation:** find parameters  $\Theta$  describing data  $D$

### Non-parametric

- The model of the distribution utilizes all examples in  $D$
- As if all examples were parameters of the distribution
- **Examples:** Nearest-neighbor

## Learning via parameter estimation

In this lecture we consider **parametric density estimation**

### Basic settings:

- A set of random variables  $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$
- **A model of the distribution** over variables in  $\mathbf{X}$   
with parameters  $\Theta : \hat{p}(\mathbf{X} | \Theta)$

**Example:** Gaussian distribution with mean and variance parameters

- **Data**  $D = \{D_1, D_2, \dots, D_n\}$

**Objective:** find parameters  $\Theta$  such that  $p(\mathbf{X} | \Theta)$  fits data  $D$  the best

## ML Parameter estimation

**Model**  $\hat{p}(\mathbf{X}) = p(\mathbf{X} | \Theta)$      **Data**  $D = \{D_1, D_2, \dots, D_n\}$

- **Maximum likelihood (ML)**

- Find  $\Theta$  that maximizes the likelihood  $p(D | \Theta, \xi)$

$$\begin{aligned}
 P(D | \Theta, \xi) &= P(D_1, D_2, \dots, D_n | \Theta, \xi) \\
 &= P(D_1 | \Theta, \xi) P(D_2 | \Theta, \xi) \dots P(D_n | \Theta, \xi) \quad \text{Independent examples} \\
 &= \prod_{i=1}^n P(D_i | \Theta, \xi)
 \end{aligned}$$

**log-likelihood**  $\log p(D | \Theta, \xi) = \sum_{i=1}^n \log P(D_i | \Theta, \xi)$

$$\Theta_{ML} = \arg \max_{\Theta} p(D | \Theta, \xi) = \arg \max_{\Theta} \log p(D | \Theta, \xi)$$

$$\hat{p}(\mathbf{X}) = p(\mathbf{X} | \Theta_{ML})$$


---

## Bayesian parameter estimation

### Bayesian parameter estimation

- uses the posterior distribution over possible parameters
- Yields: all possible settings of (and their “weights”)
- The target distribution is approximated as:

$$\begin{aligned}
 &\text{Parameter posterior} \quad \swarrow \quad \text{Data Likelihood} \\
 p(\Theta | D, \xi) &= \frac{p(D | \Theta, \xi) p(\Theta | \xi)}{p(D | \xi)} \quad \longleftarrow \text{Parameter prior}
 \end{aligned}$$

- How to use the posterior for modeling  $p(\mathbf{X})$ ?

$$\hat{p}(\mathbf{X}) = p(\mathbf{X} | D) = \int_{\Theta} p(\mathbf{X} | \Theta) p(\Theta | D, \xi) d\Theta$$


---

## Parameter estimation

### Other criteria:

- **Maximum a posteriori probability (MAP)**

maximize  $p(\Theta | D, \xi)$  (mode of the posterior)

– Yields: one set of parameters  $\Theta_{MAP}$

– Approximation:

$$\hat{p}(\mathbf{X}) = p(\mathbf{X} | \Theta_{MAP})$$

- **Expected value of the parameter**

$\hat{\Theta} = E(\Theta)$  (mean of the posterior)

– Expectation taken with regard to posterior  $p(\Theta | D, \xi)$

– Yields: one set of parameters

– Approximation:

$$\hat{p}(\mathbf{X}) = p(\mathbf{X} | \hat{\Theta})$$