

CS 2750 Machine Learning Lecture 5

Density estimation

Milos Hauskrecht

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**

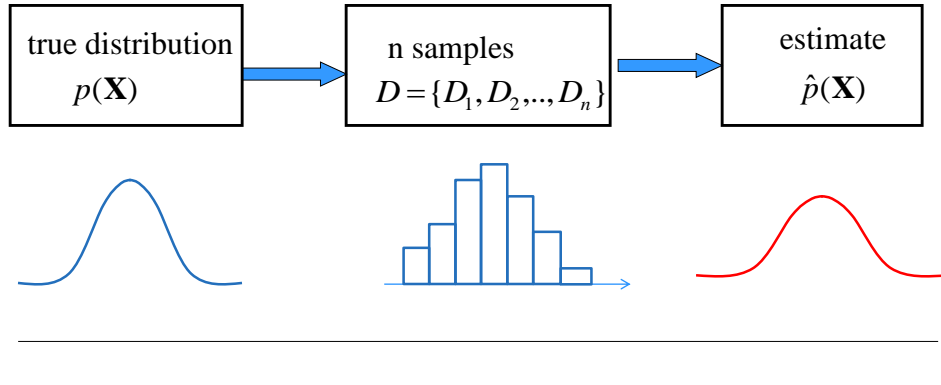
Density estimation: learn an underlying probability

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

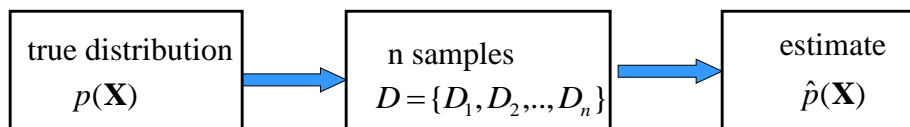
Density estimation

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

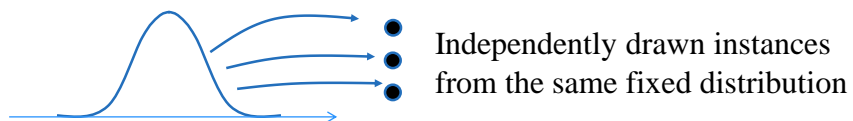


Density estimation



Standard (iid) assumptions: Samples

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

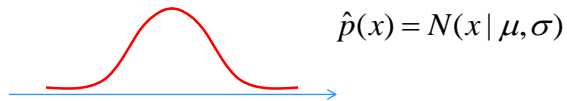


Density estimation

Types of density estimation:

(1) Parametric

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



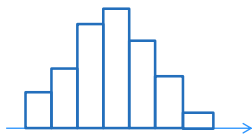
Density estimation

Types of density estimation:

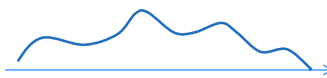
(2) Non-parametric

- The model of the distribution utilizes all examples in D
- As if all examples were parameters of the distribution
- $$\hat{p}(\mathbf{X}) = p(\mathbf{X} | D)$$
- **Examples:**

histogram



Kernel density estimation



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 $\Theta : \hat{p}(\mathbf{X} | \Theta)$
- **Data** $D = \{D_1, D_2, \dots, D_n\}$
- **Objective:** find parameters Θ such that $p(\mathbf{X} | \Theta)$ fits data D the best
- How to measure the goodness of fit or alternative the error?

Bayesian parameter estimation

The ML estimate picks just one value of the parameter

- **Problem:** if there are two different parameter values that are close in terms of the likelihood, using only one of them may introduce a strong bias, if we use it, for example, for predictions.

Bayesian parameter estimation

- Remedies the limitation of one choice
- Uses the posterior distribution for parameters Θ
- Posterior ‘covers’ all possible parameter values (and their “weights”)

$$p(\Theta | D, \xi) = \frac{p(D | \Theta, \xi) p(\Theta | \xi)}{p(D | \xi)}$$

Parameter posterior Data Likelihood Parameter prior

ML Parameter estimation

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

- **Maximum likelihood (ML)**

$$\max_{\Theta} p(D | \Theta, \xi)$$

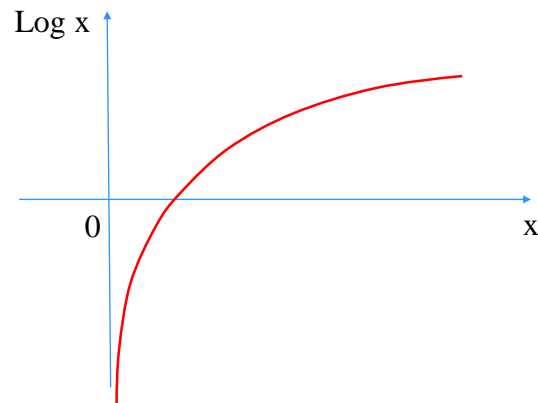
– Find Θ that maximizes likelihood $p(D | \Theta, \xi)$

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

Independent
examples

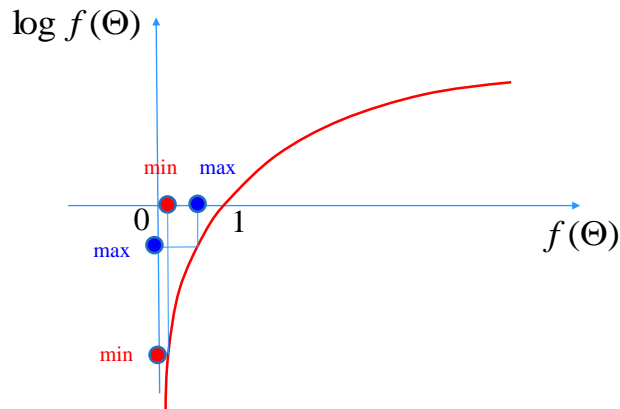
$$\Theta_{ML} = \arg \max_{\Theta} p(D | \Theta, \xi)$$

Log-likelihood



Properties of log function: ?

Log-likelihood



$$\Theta^* = \arg \max_{\Theta} f(\Theta) = \arg \max_{\Theta} \log f(\Theta)$$

ML Parameter estimation

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

• **Maximum likelihood (ML)** $\max_{\Theta} p(D | \Theta, \xi)$

– Find Θ that maximizes likelihood $p(D | \Theta, \xi)$

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

Independent examples

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

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

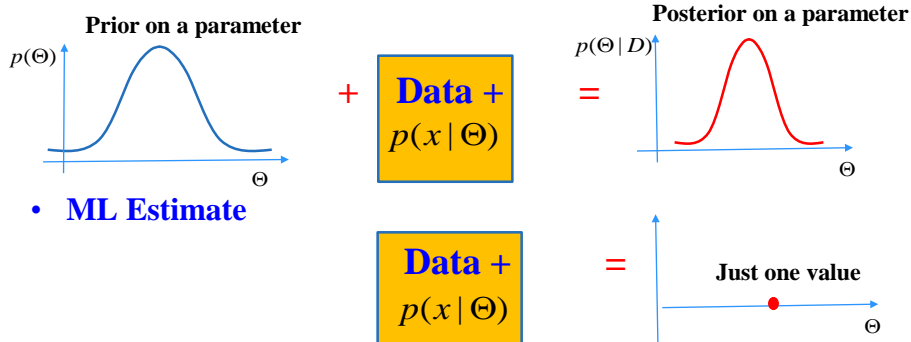
Bayesian parameter estimation

What does it do?

- Prior and Posterior ‘covers’ all possible parameter values (and their “weights”)

Assume: we have a model of $p(x | \Theta)$ with a parameter Θ

- **Bayesian parameter estimation:**



- **ML Estimate**

Bayesian parameter estimation

Bayesian parameter estimation

- Uses the posterior distribution for parameters
- Posterior ‘covers’ all possible parameter values (and their “weights”)

$$\text{Parameter posterior} \quad p(\Theta | D, \xi) = \frac{p(D | \Theta, \xi) p(\Theta | \xi)}{p(D | \xi)}$$

← Data Likelihood ← Parameter prior

- **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 Θ_{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:

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

Parameter estimation. Coin example.

Coin example: we have a coin that can be biased

Outcomes: two possible values -- head or tail

Data: D a sequence of outcomes x_i such that

- **head** $x_i = 1$

- **tail** $x_i = 0$

Model: probability of a head θ

probability of a tail $(1 - \theta)$

Objective:

We would like to estimate the probability of a **head** $\hat{\theta}$
from data



Parameter estimation. Example.

- **Assume** the unknown and possibly biased coin
- Probability of the head is θ



- **Data:**

H H T T H H T H T H T T T H T H H H H T H H H H T

- **Heads:** 15
- **Tails:** 10

What would be your estimate of the probability of a head ?

$$\tilde{\theta} = ?$$

Parameter estimation. Example

- **Assume** the unknown and possibly biased coin
- Probability of the head is θ



- **Data:**

H H T T H H T H T H T T T H T H H H H T H H H H T

- **Heads:** 15
- **Tails:** 10

What would be your choice of the probability of a head ?

Solution: use frequencies of occurrences to do the estimate

$$\tilde{\theta} = \frac{15}{25} = 0.6$$

This is **the maximum likelihood estimate** of the parameter θ

Probability of an outcome

Data: D a sequence of outcomes x_i such that

- **head** $x_i = 1$
- **tail** $x_i = 0$



Model: probability of a head θ
probability of a tail $(1-\theta)$

Assume: we know the probability θ

Probability of an outcome of a coin flip x_i

$$P(x_i | \theta) = \theta^{x_i} (1-\theta)^{(1-x_i)} \leftarrow \text{Bernoulli distribution}$$

- Combines the probability of a head and a tail
 - So that x_i is going to pick its correct probability
 - Gives θ for $x_i = 1$
 - Gives $(1-\theta)$ for $x_i = 0$
-

Probability of a sequence of outcomes.

Data: D a sequence of outcomes x_i such that

- **head** $x_i = 1$
- **tail** $x_i = 0$



Model: probability of a head θ
probability of a tail $(1-\theta)$

Assume: a sequence of independent coin flips

D = H H T H T H (encoded as **D= 110101**)

What is the probability of observing the data sequence **D**:

$$P(D | \theta) = ?$$

Probability of a sequence of outcomes.

Data: D a sequence of outcomes x_i such that

- head $x_i = 1$
- tail $x_i = 0$



Model: probability of a head θ
probability of a tail $(1-\theta)$

Assume: a sequence of coin flips $D = H H T H T H$
encoded as $D = 110101$

What is the probability of observing a data sequence D :

$$P(D | \theta) = \theta\theta(1-\theta)\theta(1-\theta)\theta$$

Probability of a sequence of outcomes.

Data: D a sequence of outcomes x_i such that

- head $x_i = 1$
- tail $x_i = 0$



Model: probability of a head θ
probability of a tail $(1-\theta)$

Assume: a sequence of coin flips $D = H H T H T H$
encoded as $D = 110101$

What is the probability of observing a data sequence D :

$$P(D | \theta) = \theta\theta(1-\theta)\theta(1-\theta)\theta$$

 **likelihood of the data**

Probability of a sequence of outcomes.

Data: D a sequence of outcomes x_i such that

- head $x_i = 1$
- tail $x_i = 0$



Model: probability of a head θ
probability of a tail $(1-\theta)$

Assume: a sequence of coin flips $D = H H T H T H$
encoded as $D = 110101$

What is the probability of observing a data sequence D :

$$P(D | \theta) = \theta\theta(1-\theta)\theta(1-\theta)\theta$$

$$P(D | \theta) = \prod_{i=1}^6 \theta^{x_i} (1-\theta)^{(1-x_i)}$$

Can be rewritten using the Bernoulli distribution:

The goodness of fit to the data

Learning: we do not know the value of the parameter θ

Our learning goal:

- Find the parameter θ that fits the data D the best?

Criterion for the best fit: Maximize the likelihood

$$P(D | \theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{(1-x_i)}$$

Intuition:

- more likely are the data given the model, the better is the fit

Note: Instead of an error function that measures how bad the data fit the model we have a measure that tells us how well the data fit :

$$Error(D, \theta) = -P(D | \theta)$$



Maximum likelihood (ML) estimate.

Likelihood of data:

$$P(D | \theta, \xi) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{(1-x_i)}$$



Maximum likelihood estimate

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

Optimize log-likelihood (the same as maximizing likelihood)

$$l(D, \theta) = \log P(D | \theta, \xi) = \log \prod_{i=1}^n \theta^{x_i} (1-\theta)^{(1-x_i)} =$$
$$\sum_{i=1}^n x_i \log \theta + (1-x_i) \log(1-\theta) = \log \theta \sum_{i=1}^n x_i + \log(1-\theta) \sum_{i=1}^n (1-x_i)$$

N_1 - number of heads seen N_2 - number of tails seen

Maximum likelihood (ML) estimate.

Optimize log-likelihood

$$l(D, \theta) = N_1 \log \theta + N_2 \log(1-\theta)$$



Set derivative to zero

$$\frac{\partial l(D, \theta)}{\partial \theta} = \frac{N_1}{\theta} - \frac{N_2}{(1-\theta)} = 0$$

Solving

$$\theta = \frac{N_1}{N_1 + N_2}$$

$$\text{ML Solution: } \theta_{ML} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$$

Maximum likelihood estimate. Example

- **Assume** the unknown and possibly biased coin
- Probability of the head is θ



- **Data:**

H H T T H H T H T H T T T H T H H H H T H H H H T

- **Heads:** 15
- **Tails:** 10

What is the ML estimate of the probability of a head and a tail?

Maximum likelihood estimate. Example

- Assume the unknown and possibly biased coin
- Probability of the head is θ



- **Data:**

H H T T H H T H T H T T T H T H H H H T H H H H T

- **Heads:** 15
- **Tails:** 10

What is the ML estimate of the probability of head and tail ?

$$\text{Head: } \theta_{ML} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2} = \frac{15}{25} = 0.6$$

$$\text{Tail: } (1 - \theta_{ML}) = \frac{N_2}{N} = \frac{N_2}{N_1 + N_2} = \frac{10}{25} = 0.4$$