

CS 1675 Introduction to Machine Learning

Lecture 6

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

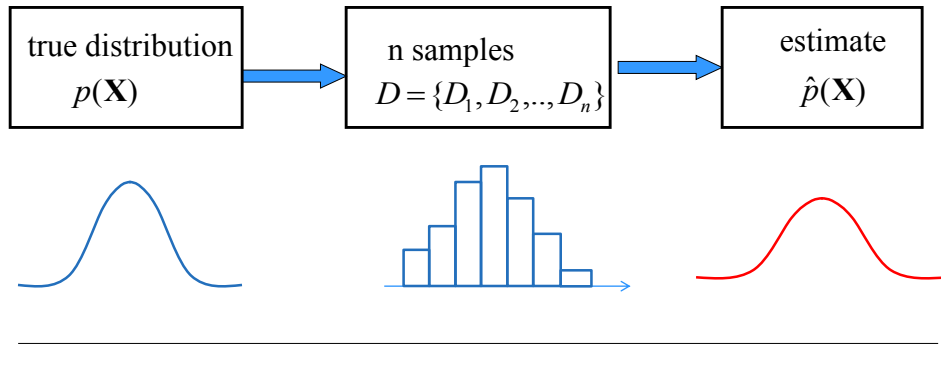
Density estimation: learn an underlying probability

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

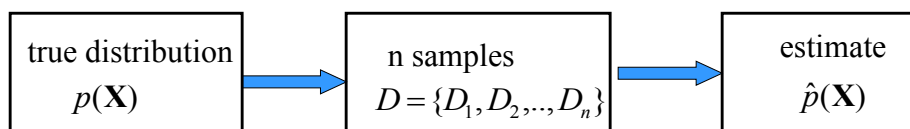
Density estimation

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

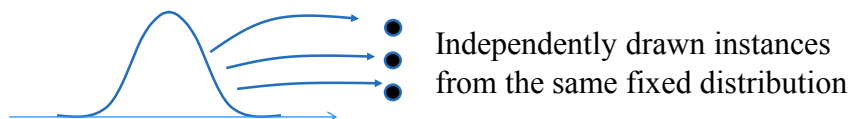


Density estimation



Standard (iid) assumptions: Samples

- are **independent** of each other
- come from the same **(i)dentical (d)istribution** (fixed $p(\mathbf{X})$)



Density estimation

Types of density estimation:

Parametric

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

Next 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 Θ : $\hat{p}(\mathbf{X} | \Theta)$

Example: Gaussian distribution with mean and variance parameters

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

Objective: find parameters Θ 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)**

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

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

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

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 outcomes 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 (0.6)
probability of a tail $(1-\theta)$ (0.4)

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 or 0.6 for $x_i = 1$
- Gives $(1-\theta)$ or 0.4 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 (0.6)
probability of a tail $(1-\theta)$ (0.4)

Assume: a sequence of independent coin flips

$D = \text{H H T H T H}$ (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 (0.6)
probability of a tail $(1-\theta)$ (0.4)

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 :

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

$$P(D | \theta) \equiv 0.6 * 0.6 * 0.4 * 0.6 * 0.4 * 0.6 = 0.6^4 * 0.4^2$$

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?

One solution to the “best”: 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

$$\theta_{ML} = \arg \max_{\theta} P(D | \theta, \xi)$$

Optimize log-likelihood (the same as maximizing likelihood)

$$\begin{aligned} l(D, \theta) &= \log P(D | \theta, \xi) = \log \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)} = \\ &= \sum_{i=1}^n x_i \log \theta + (1 - x_i) \log(1 - \theta) = \log \theta \underbrace{\sum_{i=1}^n x_i}_{N_1} + \log(1 - \theta) \underbrace{\sum_{i=1}^n (1 - x_i)}_{N_2} \end{aligned}$$

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

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

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 ?

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

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

Maximum a posteriori estimate

Maximum a posteriori estimate

- Selects the mode of the **posterior distribution**

$$\theta_{MAP} = \arg \max_{\theta} p(\theta | D, \xi)$$

Likelihood of data

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) p(\theta | \xi)}{P(D | \xi)} \quad (\text{via Bayes rule})$$

prior

Normalizing factor

$$P(D | \theta, \xi) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)} = \theta^{N_1} (1 - \theta)^{N_2}$$

$p(\theta | \xi)$ - is the prior probability on θ

How to choose the prior probability?

CS 2750 Machine Learning

Prior distribution

Choice of prior: Beta distribution

$$p(\theta | \xi) = \text{Beta}(\theta | \alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \theta^{\alpha_1-1} (1-\theta)^{\alpha_2-1}$$

$\Gamma(x)$ - a Gamma function $\Gamma(x) = (x-1)\Gamma(x-1)$

For integer values of x $\Gamma(n) = (n-1)!$

Why to use Beta distribution?

Beta distribution “fits” Bernoulli trials - **conjugate choices**

$$P(D | \theta, \xi) = \theta^{N_1} (1 - \theta)^{N_2}$$

Posterior distribution is again a Beta distribution

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) \text{Beta}(\theta | \alpha_1, \alpha_2)}{P(D | \xi)} = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$

CS 2750 Machine Learning

Bernoulli distribution

Data D: iid sample of n outcomes (coin flips)

Posterior of data:

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

Likelihood

$$P(D | \theta, \xi) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{(1-x_i)} = \theta^{N_1} (1 - \theta)^{N_2}$$

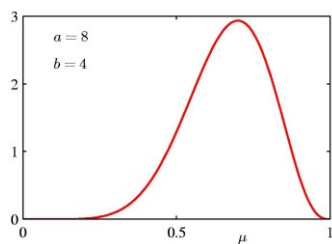
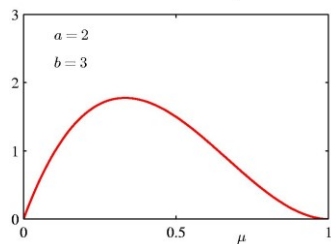
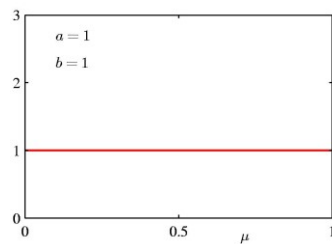
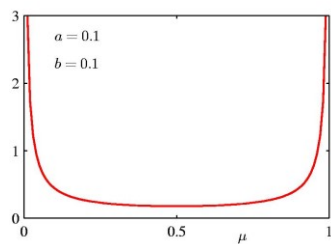
Conjugate prior:

$$p(\theta | \xi) = \text{Beta}(\theta | \alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \theta^{\alpha_1-1} (1 - \theta)^{\alpha_2-1}$$

Posterior:

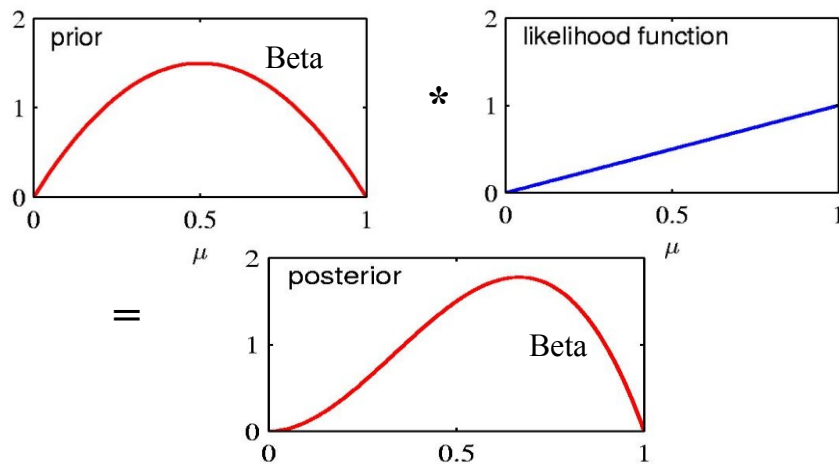
$$p(\theta | D, \xi) = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$

Beta distribution



$$p(\theta | \xi) = \text{Beta}(\theta | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1 - \theta)^{b-1}$$

Posterior distribution



$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) \text{Beta}(\theta | \alpha_1, \alpha_2)}{P(D | \xi)} = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$

CS 2750 Machine Learning

Maximum a posterior probability

Maximum a posteriori estimate

- Selects the mode of the **posterior distribution**

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) \text{Beta}(\theta | \alpha_1, \alpha_2)}{P(D | \xi)} = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$

$$= \frac{\Gamma(\alpha_1 + \alpha_2 + N_1 + N_2)}{\Gamma(\alpha_1 + N_1) \Gamma(\alpha_2 + N_2)} \theta^{N_1 + \alpha_1 - 1} (1 - \theta)^{N_2 + \alpha_2 - 1}$$

Notice that parameters of the prior
act like counts of heads and tails
(sometimes they are also referred to as **prior counts**)

MAP Solution:

$$\theta_{MAP} = \frac{\alpha_1 + N_1 - 1}{\alpha_1 + \alpha_2 + N_1 + N_2 - 2}$$

CS 2750 Machine Learning

MAP 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

- Assume $p(\theta | \xi) = \text{Beta}(\theta | 5, 5)$

What is the MAP estimate?

CS 2750 Machine Learning

MAP 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

- Assume $p(\theta | \xi) = \text{Beta}(\theta | 5, 5)$

What is the MAP estimate ?

$$\theta_{MAP} = \frac{N_1 + \alpha_1 - 1}{N - 2} = \frac{N_1 + \alpha_1 - 1}{N_1 + N_2 + \alpha_1 + \alpha_2 - 2} = \frac{19}{33}$$

CS 2750 Machine Learning

MAP estimate example

- Note that the prior and data fit (data likelihood) are combined
- **The MAP can be biased with large prior counts**
- **It is hard to overturn it with a smaller sample size**
- **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

- Assume

$$p(\theta | \xi) = \text{Beta}(\theta | 5, 5) \quad \theta_{MAP} = \frac{19}{33}$$

$$p(\theta | \xi) = \text{Beta}(\theta | 5, 20) \quad \theta_{MAP} = \frac{19}{48}$$

CS 2750 Machine Learning

Binomial distribution

$$\text{5 coins} = 2 * \text{heads} + 3 * \text{tails}$$

Example problem: N coin flips, where each coin flip can have two results: head or tail

Outcome: N_1 - number of heads seen N_2 - number of tails seen in N trials

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

Probability of an outcome:

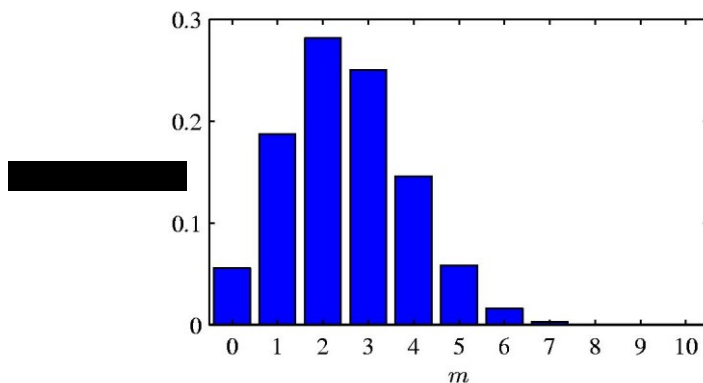
$$P(N_1 | N, \theta) = \binom{N}{N_1} \theta^{N_1} (1 - \theta)^{N - N_1} \quad \text{Binomial distribution}$$

Binomial distribution:

- **models order independent sequence of independent Bernoulli trials**

Binomial distribution

Binomial distribution:



CS 2750 Machine Learning

Maximum likelihood (ML) estimate.

Likelihood of data:

$$P(D|\theta) = \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N_2} = \frac{N!}{N_1! N_2!} \theta^{N_1} (1-\theta)^{N_2}$$

Log-likelihood

$$l(D, \theta) = \log \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N_2} = \log \frac{N!}{N_1! N_2!} + N_1 \log \theta + N_2 \log(1-\theta)$$

Constant from the point of optimization !!!

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

The same as for Bernoulli and D with iid sequence of examples

CS 2750 Machine Learning

Posterior density

Posterior density

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) p(\theta | \xi)}{P(D | \xi)} \quad (\text{via Bayes rule})$$

Prior choice

$$p(\theta | \xi) = \text{Beta}(\theta | \alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \theta^{\alpha_1-1} (1-\theta)^{\alpha_2-1}$$

Likelihood

$$P(D | \theta) = \frac{\Gamma(N_1 + N_2)}{\Gamma(N_1)\Gamma(N_2)} \theta^{N_1} (1-\theta)^{N_2}$$

Posterior

$$p(\theta | D, \xi) = \text{Beta}(\alpha_1 + N_1, \alpha_2 + N_2)$$

MAP estimate

$$\theta_{MAP} = \arg \max_{\theta} p(\theta | D, \xi)$$
$$\theta_{MAP} = \frac{\alpha_1 + N_1 - 1}{\alpha_1 + \alpha_2 + N_1 + N_2 - 2}$$