

## CS 2750 Machine Learning Lecture 5

### 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

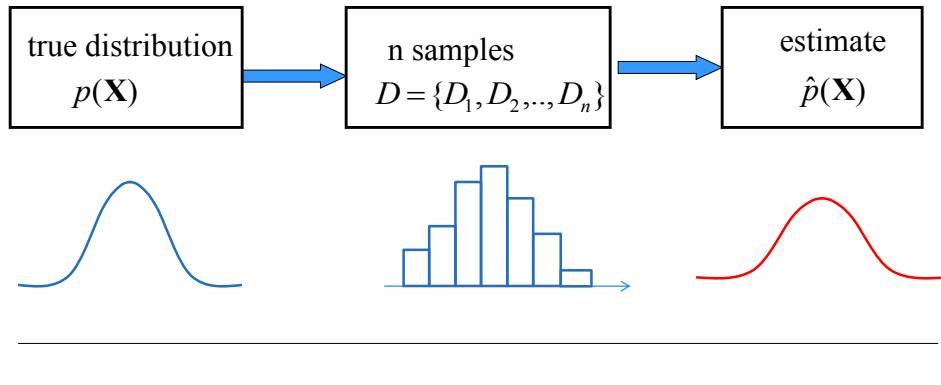
**distribution model :**  $p(\mathbf{X}) = p(X_1, X_2, \dots, X_d)$  from  $\mathbf{D}$

---

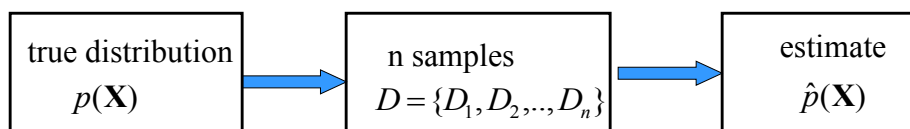
## Density estimation

**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  $D$

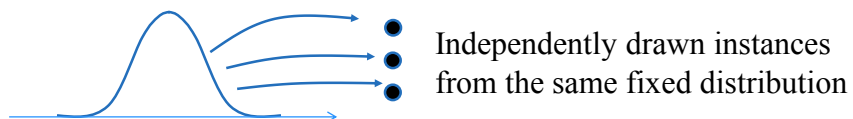


## Density estimation



**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

**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  $\Theta$
- Posterior ‘covers’ all possible parameter values (and their “weights”)

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


---

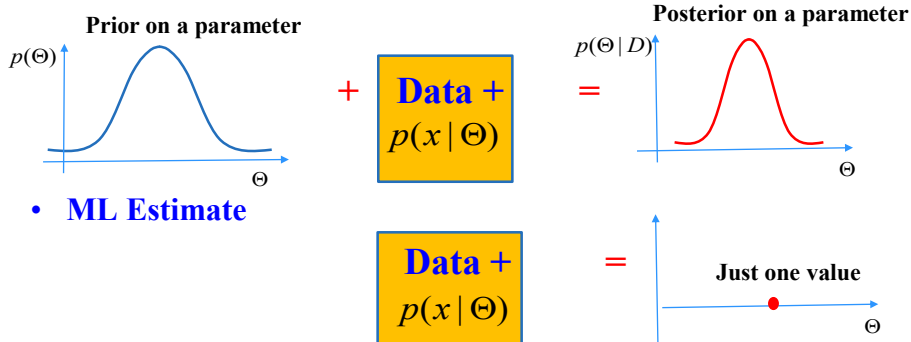
## Bayesian parameter estimation

### 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  $\Theta$

- **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”)

Parameter posterior Data Likelihood

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

← Parameter prior

- **How to use the posterior for modeling  $p(X)$ ?**

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

## 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:

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

## 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  $\theta$

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  $\theta$



- **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  $\theta$



- **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  $\theta$

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

**Assume:** we know the probability  $\theta$

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

**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  $\theta$   
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  $\theta$   
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  $\theta$   
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  $\theta$

**Our learning goal:**

- Find the parameter  $\theta$  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)**

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

$N_1$  - 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**

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

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

## Maximum likelihood estimate. Example

- **Assume** the unknown and possibly biased coin
- Probability of the head is  $\theta$



- **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  $\theta$



- **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 ?

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

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

## Bayesian parameter estimation

Uses the distributions (prior and posterior) over all possible values of the parameter  $\theta$  of the sampling distribution  $p(x|\theta)$  (Bernoulli):

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

Labels in the diagram:  
- **Likelihood of data** points to  $P(D|\theta, \xi)$   
- **Prior** points to  $p(\theta|\xi)$   
- **Posterior** points to  $p(\theta|D, \xi)$   
- **Normalizing factor** points to  $P(D|\xi)$

We know that the likelihood is:

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

How to choose the prior probability?

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

## 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?

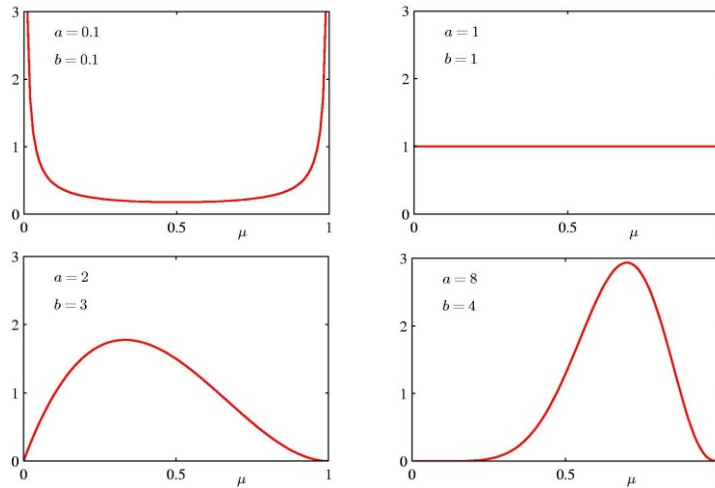
Beta distribution “fits” Bernoulli sample - **conjugate choices**

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

Posterior distribution is again a Beta distribution

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

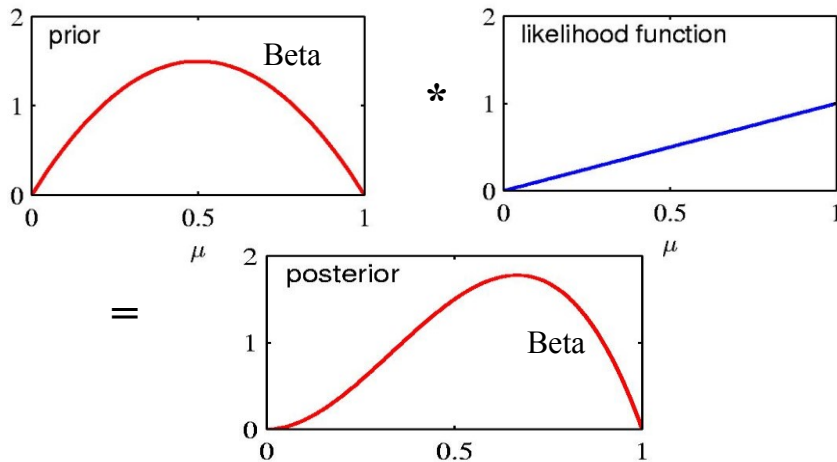
## Beta distribution



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

CS 2750 Machine Learning

## Posterior distribution



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

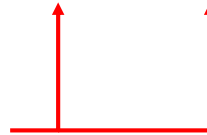
## Posterior distribution

### Beta posterior

– A conjugate prior to Bernoulli sample

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) \text{Beta}(\theta | \alpha_1, \alpha_2)}{P(D | \xi)} = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$
$$= \frac{\Gamma(\alpha_1 + \alpha_2 + N_1 + N_2)}{\Gamma(\alpha_1 + N_1) \Gamma(\alpha_2 + N_2)} \theta^{N_1 + \alpha_1 - 1} (1 - \theta)^{N_2 + \alpha_2 - 1}$$

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



## Maximum a posteriori probability (MAP)

### Maximum a posteriori estimate

– Selects **the mode of the posterior distribution**

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

**Likelihood of data**

**prior**

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

**Normalizing factor**

- Selects the model of the posterior represented as a Beta distribution

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

## Maximum posterior probability

### Maximum a posteriori estimate

- Selects the mode of the **posterior distribution**
- **Assumes conjugate prior to Bernoulli sample**

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) \text{Beta}(\theta | \alpha_1, \alpha_2)}{P(D | \xi)} = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$
$$= \frac{\Gamma(\alpha_1 + \alpha_2 + N_1 + N_2)}{\Gamma(\alpha_1 + N_1) \Gamma(\alpha_2 + N_2)} \theta^{N_1 + \alpha_1 - 1} (1 - \theta)^{N_2 + \alpha_2 - 1}$$

Mode of the posterior satisfies :  $\frac{\partial \log p(\theta | D, \xi)}{\partial \theta} = 0$

|  |
|--|
| <b>MAP Solution:</b> $\theta_{MAP} = \frac{\alpha_1 + N_1 - 1}{\alpha_1 + \alpha_2 + N_1 + N_2 - 2}$ |
|--|

## MAP estimate example

- Assume the unknown and possibly biased coin
- Probability of the head is  $\theta$

- **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?



## MAP estimate example

- Assume the unknown and possibly biased coin
- Probability of the head is  $\theta$

- **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

## Bayesian framework

- **Predictive probability of an outcome  $x=1$  in the next trial**

$$P(x=1 | D, \xi)$$

$$\begin{aligned}
 P(x=1 | D, \xi) &= \int_0^1 P(x=1 | \theta, \xi) \overbrace{p(\theta | D, \xi)}^{\text{Posterior density}} d\theta \\
 &= \int_0^1 \theta p(\theta | D, \xi) d\theta = E(\theta)
 \end{aligned}$$

- **Equivalent to the expected value of the parameter**
  - expectation is taken with respect to the posterior distribution

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

CS 2750 Machine Learning

## Expected value of the parameter

**How to calculate the expected value of Beta?**

$$\begin{aligned}
 E(\theta) &= \int_0^1 \theta \text{Beta}(\theta | \eta_1, \eta_2) d\theta = \int_0^1 \theta \frac{\Gamma(\eta_1 + \eta_2)}{\Gamma(\eta_1)\Gamma(\eta_2)} \theta^{\eta_1-1} (1-\theta)^{\eta_2-1} d\theta \\
 &= \frac{\Gamma(\eta_1 + \eta_2)}{\Gamma(\eta_1)\Gamma(\eta_2)} \int_0^1 \theta^{\eta_1} (1-\theta)^{\eta_2-1} d\theta \\
 &= \frac{\Gamma(\eta_1 + \eta_2)}{\Gamma(\eta_1)\Gamma(\eta_2)} \frac{\Gamma(\eta_1 + 1)\Gamma(\eta_2)}{\Gamma(\eta_1 + \eta_2 + 1)} \underbrace{\int_0^1 \text{Beta}(\eta_1 + 1, \eta_2) d\theta}_1 \\
 &= \frac{\eta_1}{\eta_1 + \eta_2}
 \end{aligned}$$

**Note:**  $\Gamma(\alpha + 1) = \alpha\Gamma(\alpha)$  for integer values of  $\alpha$

CS 2750 Machine Learning

## Expected value of the parameter

- **Substituting the results for the posterior:**

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

- **We get**  $E(\theta) = \frac{\alpha_1 + N_1}{\alpha_1 + N_1 + \alpha_2 + N_2}$

- **Note that the mean of the posterior is yet another** “reasonable” parameter choice:

$$\hat{\theta} = E(\theta)$$

---

CS 2750 Machine Learning

## Binomial distribution

**Example problem:** a biased coin

**Outcomes:** two possible values -- head or tail

**Data:** a set of order-independent outcomes for N trials

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

**can be calculated from the trial data !!!**

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

**Probability of an outcome**

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

**Objective:**

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

---

## Binomial distribution



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

**Probability of an outcome:**

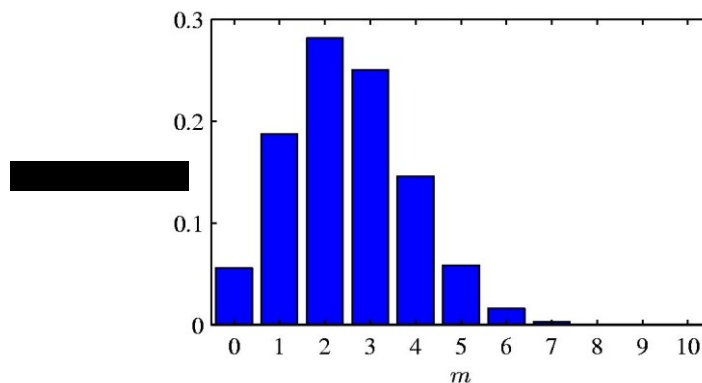
$$P(N_1 | N, \theta) = \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N-N_1} \quad \text{Binomial distribution}$$

**Binomial distribution:**

- models order independent sequence of Bernoulli trials

## Binomial distribution

**Binomial distribution:**



## 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 !!!

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

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

## 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

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

## Multinomial distribution



**Example:** multiple rolls of a dice with 6 results

**Outcome:** counts of occurrences of  $k$  possible outcomes of  $N$  trials:

$$\sum_{i=1}^k N_i = N$$

**Model parameters:**  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$  s.t.  $\sum_{i=1}^k \theta_i = 1$   
 $\theta_i$  - probability of an outcome  $i$

**Probability distribution:**

$$P(N_1, N_2, \dots, N_k | \boldsymbol{\theta}, \xi) = \frac{N!}{N_1! N_2! \dots N_k!} \theta_1^{N_1} \theta_2^{N_2} \dots \theta_k^{N_k} \quad \text{Multinomial distribution}$$

**ML estimate:**

$$\theta_{i,ML} = \frac{N_i}{N}$$

## Posterior and MAP estimate



**Choice of the prior: Dirichlet distribution**

$$Dir(\boldsymbol{\theta} | \alpha_1, \dots, \alpha_k) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \dots \theta_k^{\alpha_k-1}$$

**Dirichlet is the conjugate choice for the multinomial sampling**

$$P(D | \boldsymbol{\theta}, \xi) = P(N_1, N_2, \dots, N_k | \boldsymbol{\theta}, \xi) = \frac{N!}{N_1! N_2! \dots N_k!} \theta_1^{N_1} \theta_2^{N_2} \dots \theta_k^{N_k}$$

**Posterior density**

$$p(\boldsymbol{\theta} | D, \xi) = \frac{P(D | \boldsymbol{\theta}, \xi) Dir(\boldsymbol{\theta} | \alpha_1, \alpha_2, \dots, \alpha_k)}{P(D | \xi)} = Dir(\boldsymbol{\theta} | \alpha_1 + N_1, \dots, \alpha_k + N_k)$$

**MAP estimate:**

$$\theta_{i,MAP} = \frac{\alpha_i + N_i - 1}{\sum_{i=1, \dots, k} (\alpha_i + N_i) - k}$$

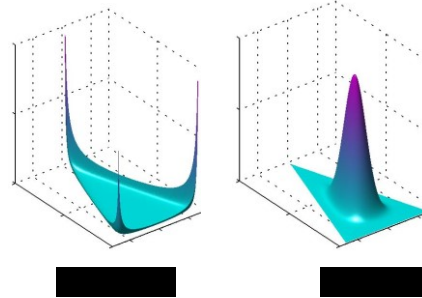
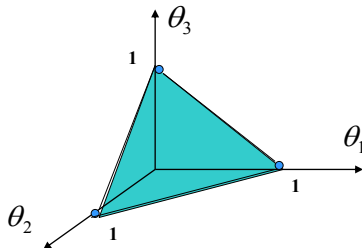
## Dirichlet distribution



**Dirichlet distribution:**

$$Dir(\boldsymbol{\theta} | \alpha_1, \dots, \alpha_k) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \dots \theta_k^{\alpha_k-1}$$

**Assume: k=3**



## Other distributions

**The same ideas can be applied to other distributions**

- Typically we choose distributions that behave well so that computations lead to “nice” solutions

- **Exponential family of distributions**

**Conjugate choices** for some of the distributions from the exponential family:

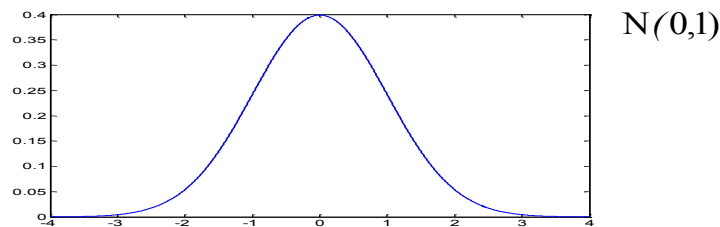
- **Binomial – Beta**
- **Multinomial - Dirichlet**
- **Exponential – Gamma**
- **Poisson – Inverse Gamma**
- **Gaussian - Gaussian (mean) and Wishart (covariance)**

## Gaussian (normal) distribution

- **Gaussian:**  $x \sim N(\mu, \sigma)$
- **Parameters:**  $\mu$  - mean  
 $\sigma$  - standard deviation
- **Density function:**

$$p(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right]$$

- **Example:**



CS 2750 Machine Learning

## Parameter estimates

- **Loglikelihood**  $l(D, \mu, \sigma) = \log \prod_{i=1}^n p(x_i | \mu, \sigma)$

- **ML estimates of the mean and variance:**

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

– ML variance estimate is biased

$$E_n(\hat{\sigma}^2) = E_n\left(\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2\right) = \frac{n-1}{n} \sigma^2 \neq \sigma^2$$

- **Unbiased estimate:**

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

CS 2750 Machine Learning

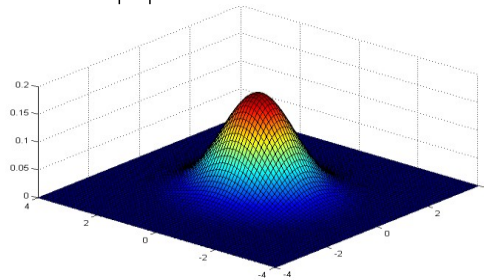


## Multivariate normal distribution

- **Multivariate normal:**  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- **Parameters:**  $\boldsymbol{\mu}$ - mean  
 $\boldsymbol{\Sigma}$ - covariance matrix
- **Density function:**

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

- **Example:**



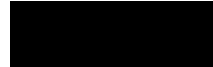
CS 2750 Machine Learning

## Partitioned Gaussian Distributions

- **Multivariate Gaussian:**



- **Example:**



Precision matrix

- **What are the distributions for marginals and conditionals?**

$$p(x_a)$$

$$p(x_a | x_b)$$

## Partitioned Conditionals and Marginals

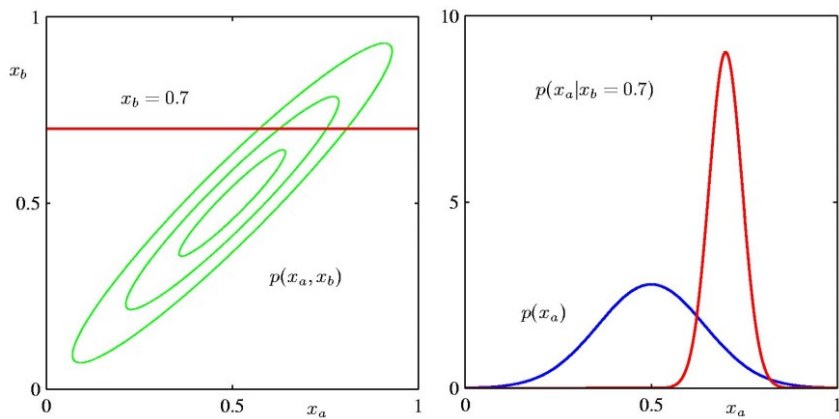
- Conditional density:



- Marginal Density:



## Partitioned Conditionals and Marginals



## Parameter estimates

- **Loglikelihood**  $l(D, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

- **ML estimates of the mean and covariances:**

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

– Covariance estimate is biased

$$E_n(\hat{\boldsymbol{\Sigma}}) = E_n\left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T\right) = \frac{n-1}{n} \boldsymbol{\Sigma} \neq \boldsymbol{\Sigma}$$

- **Unbiased estimate:**

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

CS 2750 Machine Learning

## Posterior of a multivariate normal

- **Assume a prior on the mean  $\boldsymbol{\mu}$  that is normally distributed:**

$$p(\boldsymbol{\mu}) \approx N(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$$

- **Then the posterior of  $\boldsymbol{\mu}$  is normally distributed**

$$\begin{aligned} p(\boldsymbol{\mu} | D) &\approx \left( \prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})\right] \right) \\ &\quad * \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_p|^{1/2}} \exp\left[-\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_p)^T \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_p)\right] \\ &= \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_n|^{1/2}} \exp\left[-\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)^T \boldsymbol{\Sigma}_n^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)\right] \end{aligned}$$

CS 2750 Machine Learning

## Posterior of a multivariate normal

- Then the posterior of  $\boldsymbol{\mu}$  is normally distributed

$$p(\boldsymbol{\mu} | D) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_n|^{1/2}} \exp\left[-\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)^T \boldsymbol{\Sigma}_n^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)\right]$$

$$\boldsymbol{\Sigma}_n^{-1} = n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}_p^{-1}$$

$$\boldsymbol{\mu}_n = \boldsymbol{\Sigma}_p \left( \boldsymbol{\Sigma}_p + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right) + \frac{1}{n} \boldsymbol{\Sigma} \left( \boldsymbol{\Sigma}_p + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\mu}_p$$

$$\boldsymbol{\Sigma}_n = \boldsymbol{\Sigma}_p \left( \boldsymbol{\Sigma}_p + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \frac{1}{n} \boldsymbol{\Sigma}$$

CS 2750 Machine Learning

## Other distributions

### Gamma distribution:

$$p(x | a, b) = \frac{1}{\Gamma(a)b^a} x^{a-1} e^{-\frac{x}{b}} \quad \text{for } x \in [0, \infty]$$

### Exponential distribution:

- A special case of Gamma for  $a=1$

$$p(x | b) = \left( \frac{1}{b} \right) e^{-\frac{x}{b}}$$

### Poisson distribution:

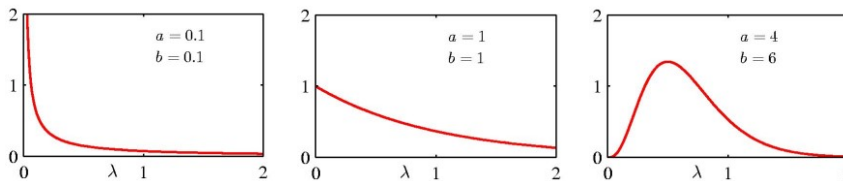
$$p(x | \lambda) = \frac{e^{-\lambda} \lambda^x}{x!} \quad \text{for } x \in \{0, 1, 2, \dots\}$$

CS 2750 Machine Learning

## Other distributions

### Gamma distribution:

$$p(\lambda | a, b) = \frac{1}{\Gamma(a)b^a} \lambda^{a-1} e^{-\frac{\lambda}{b}} \quad \text{for } \lambda \in [0, \infty]$$



## Sequential Bayesian parameter estimation

- **Sequential Bayesian approach**

- Under the iid the estimates of the posterior can be computed incrementally for a sequence of data points

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

- If we use a conjugate prior we get back the same posterior
- Assume we split the data D in the last element  $\mathbf{x}$  and the rest  $p(D | \Theta) = P(\mathbf{x} | \Theta) P(D_{n-1} | \Theta)$

- **Then:** A “new” prior

$$p(\Theta | D, \xi) = \frac{P(\mathbf{x} | \Theta) \overbrace{P(D_{n-1} | \Theta) p(\Theta | \xi)}^{\text{A “new” prior}}}{\int_{\Theta} P(\mathbf{x} | \Theta) P(D_{n-1} | \Theta) p(\Theta | \xi) d\Theta}$$

## Exponential family

### Exponential family:

- all probability mass / density functions that can be written in the exponential normal form

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$$

- $\boldsymbol{\eta}$  a vector of natural (or canonical) parameters
- $t(\mathbf{x})$  a function referred to as a sufficient statistic
- $h(\mathbf{x})$  a function of  $\mathbf{x}$  (it is less important)
- $Z(\boldsymbol{\eta})$  a normalization constant (a partition function)

$$Z(\boldsymbol{\eta}) = \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}$$

- Other common form:

$$f(\mathbf{x} | \boldsymbol{\eta}) = h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x}) - A(\boldsymbol{\eta})] \quad \log Z(\boldsymbol{\eta}) = A(\boldsymbol{\eta})$$

CS 2750 Machine Learning

## Exponential family: examples

- Bernoulli distribution**

$$\begin{aligned} p(x | \pi) &= \pi^x (1 - \pi)^{1-x} \\ &= \exp\left\{ \log\left(\frac{\pi}{1-\pi}\right)x + \log(1-\pi) \right\} \\ &= \exp\{\log(1-\pi)\} \exp\left\{ \log\left(\frac{\pi}{1-\pi}\right)x \right\} \end{aligned}$$

- Exponential family**

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$$

- Parameters**

$$\boldsymbol{\eta} = ?$$

$$t(\mathbf{x}) = ?$$

$$Z(\boldsymbol{\eta}) = ?$$

$$h(\mathbf{x}) = ?$$

CS 2750 Machine Learning

## Exponential family: examples

- **Bernoulli distribution**

$$\begin{aligned} p(x | \pi) &= \pi^x (1 - \pi)^{1-x} \\ &= \exp \left\{ \log \left( \frac{\pi}{1 - \pi} \right) x + \log(1 - \pi) \right\} \\ &= \exp \{ \log(1 - \pi) \} \exp \left\{ \log \left( \frac{\pi}{1 - \pi} \right) x \right\} \end{aligned}$$

- **Exponential family**

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$$

- **Parameters**

$$\boldsymbol{\eta} = \log \frac{\pi}{1 - \pi} \quad (\text{note } \pi = \frac{1}{1 + e^{-\eta}}) \quad t(\mathbf{x}) = x$$

$$Z(\boldsymbol{\eta}) = \frac{1}{1 - \pi} = 1 + e^{\eta} \quad h(\mathbf{x}) = 1$$

CS 2750 Machine Learning

## Exponential family: examples

- **Univariate Gaussian distribution**

$$\begin{aligned} p(x | \mu, \sigma) &= \frac{1}{\sigma \sqrt{2\pi}} \exp \left[ -\frac{1}{2\sigma^2} (x - \mu)^2 \right] \\ &= \frac{1}{2\pi} \exp \left( -\frac{\mu}{2\sigma^2} - \log \sigma \right) \exp \left\{ \frac{\mu}{\sigma^2} x - \frac{1}{2\sigma^2} x^2 \right\} \end{aligned}$$

- **Exponential family**

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(x) \exp[\boldsymbol{\eta}^T t(x)]$$

- **Parameters**

$$\boldsymbol{\eta} = ? \quad t(\mathbf{x}) = ?$$

$$Z(\boldsymbol{\eta}) = ? \quad h(\mathbf{x}) = ?$$

CS 2750 Machine Learning

## Exponential family: examples

- **Univariate Gaussian distribution**

$$\begin{aligned} p(x | \mu, \sigma) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2} (x - \mu)^2\right] \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\mu}{2\sigma^2} - \log \sigma\right) \exp\left\{\frac{\mu}{\sigma^2} x - \frac{1}{2\sigma^2} x^2\right\} \end{aligned}$$

- **Exponential family**  $f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$

- **Parameters**

$$\boldsymbol{\eta} = \begin{bmatrix} \mu / 2\sigma^2 \\ -1 / 2\sigma^2 \end{bmatrix} \quad t(\mathbf{x}) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$
$$Z(\boldsymbol{\eta}) = \exp\left\{\frac{\mu}{2\sigma^2} + \log \sigma\right\} = \exp\left\{-\frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2)\right\}$$
$$h(\mathbf{x}) = 1/\sqrt{2\pi}$$

CS 2750 Machine Learning

## Exponential family

- **For iid samples, the likelihood of data is**

$$\begin{aligned} P(D | \boldsymbol{\eta}) &= \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\eta}) = \prod_{i=1}^n h(\mathbf{x}_i) \exp[\boldsymbol{\eta}^T t(\mathbf{x}_i) - A(\boldsymbol{\eta})] \\ &= \left[ \prod_{i=1}^n h(\mathbf{x}_i) \right] \exp\left[ \sum_{i=1}^n \boldsymbol{\eta}^T t(\mathbf{x}_i) - nA(\boldsymbol{\eta}) \right] \\ &= \left[ \prod_{i=1}^n h(\mathbf{x}_i) \right] \exp\left[ \boldsymbol{\eta}^T \left( \sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right] \end{aligned}$$

- **Important:**

- the dimensionality of the sufficient statistic remains the same with the number of samples

CS 2750 Machine Learning



## Exponential family

- The log likelihood of data is

$$\begin{aligned}l(D, \boldsymbol{\eta}) &= \log \left[ \prod_{i=1}^n h(\mathbf{x}_i) \right] \exp \left[ \boldsymbol{\eta}^T \left( \sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right] \\ &= \log \left[ \prod_{i=1}^n h(\mathbf{x}_i) \right] + \left[ \boldsymbol{\eta}^T \left( \sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right]\end{aligned}$$

- Optimizing the loglikelihood

$$\nabla_{\boldsymbol{\eta}} l(D, \boldsymbol{\eta}) = \left( \sum_{i=1}^n t(\mathbf{x}_i) \right) - n \nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \mathbf{0}$$

- For the ML estimate it must hold

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \frac{1}{n} \left( \sum_{i=1}^n t(\mathbf{x}_i) \right)$$

---

CS 2750 Machine Learning

## Exponential family

- Rewriting the gradient:

---

CS 2750 Machine Learning

## Exponential family

- **Rewriting the gradient:**

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \log Z(\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \log \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \frac{\int t(\mathbf{x}) h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}}{\int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \int t(\mathbf{x}) h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x}) - A(\boldsymbol{\eta})\} d\mathbf{x}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = E(t(\mathbf{x}))$$

- **Result:** 
$$E(t(\mathbf{x})) = \frac{1}{n} \left( \sum_{i=1}^n t(\mathbf{x}_i) \right)$$
- **For the ML estimate the parameters  $\boldsymbol{\eta}$  should be adjusted such that the expectation of the statistic  $t(\mathbf{x})$  is equal to the observed sample statistics**

CS 2750 Machine Learning

## Moments of the distribution

- **For the exponential family**

- The k-th moment of the statistic corresponds to the k-th derivative of  $A(\boldsymbol{\eta})$
- If  $x$  is a component of  $t(\mathbf{x})$  then we get the moments of the distribution by differentiating its corresponding natural parameter

- **Example: Bernoulli**  $p(x | \pi) = \exp\left\{ \log\left(\frac{\pi}{1-\pi}\right)x + \log(1-\pi) \right\}$

$$A(\boldsymbol{\eta}) = \log \frac{1}{1-\pi} = \log(1 + e^{\boldsymbol{\eta}})$$

- **Derivatives:**

$$\frac{\partial A(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \frac{\partial}{\partial \boldsymbol{\eta}} \log(1 + e^{\boldsymbol{\eta}}) = \frac{e^{\boldsymbol{\eta}}}{(1 + e^{\boldsymbol{\eta}})} = \frac{1}{(1 + e^{-\boldsymbol{\eta}})} = \pi$$

$$\frac{\partial A(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^2} = \frac{\partial}{\partial \boldsymbol{\eta}} \frac{1}{(1 + e^{-\boldsymbol{\eta}})} = \pi(1 - \pi)$$

CS 2750 Machine Learning

**End**

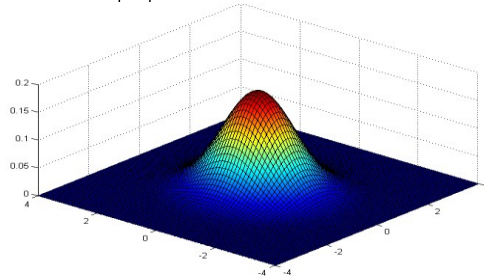
CS 2750 Machine Learning

## Multivariate normal distribution

- **Multivariate normal:**  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- **Parameters:**  $\boldsymbol{\mu}$ - mean  
 $\boldsymbol{\Sigma}$ - covariance matrix
- **Density function:**

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right]$$

- **Example:**



CS 2750 Machine Learning

## Parameter estimates

- **Loglikelihood**  $l(D, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

- **ML estimates of the mean and covariances:**

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

– Covariance estimate is biased

$$E_n(\hat{\boldsymbol{\Sigma}}) = E_n\left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T\right) = \frac{n-1}{n} \boldsymbol{\Sigma} \neq \boldsymbol{\Sigma}$$

- **Unbiased estimate:**

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

CS 2750 Machine Learning

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

**Objective:** find parameters  $\hat{\Theta}$  that fit the data the best

What is the best set of parameters?

There are various criteria one can apply here ...

CS 2750 Machine Learning

## Parameter estimation.

- **Maximum likelihood (ML)**

maximize  $p(D | \Theta, \xi)$

$\xi$  - represents prior (background) knowledge

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

maximize  $p(\Theta | D, \xi)$

**Selects the mode of the posterior**

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

- **Bayesian framework**

- use a posterior density
- no optimization

CS 2750 Machine Learning

## Posterior of a multivariate normal

- Assume that we use only a prior on the mean:  $\mu$

- A prior  $\mu \approx N(\mu_p, \Sigma_p)$

- Then the posterior is:

- Normally

$$p(\mathbf{x} | \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right]$$

- ML estimates of the mean and covariances:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T$$

- Covariance estimate is biased

$$E_n(\hat{\Sigma}) = E_n\left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T\right) = \frac{n-1}{n} \Sigma \neq \Sigma$$

- Unbiased estimate:  $\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T$

CS 2750 Machine Learning

## Parameter estimates

- **Loglikelihood**  $l(D, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

- **ML estimates of the mean and covariances:**

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

- **Unbiased estimate:**

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

---

CS 2750 Machine Learning

## Unsupervised learning

- **Data:**  $D = \{D_1, D_2, \dots, D_n\}$   
 $D_i = \mathbf{x}_i$  a vector of attribute values
  - e.g. the description of a patient
  - no specific target attribute we want to predict (no output y)
- **Objective:**
  - learn (describe) relations between attributes, examples

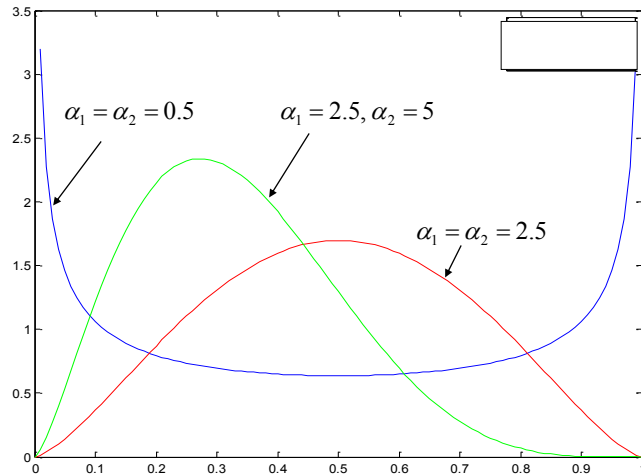
### Types of problems:

- **Clustering**
  - Group together “similar” examples
- **Density estimation**
  - Model probabilistically the population of examples

---

CS 2750 Machine Learning

## Beta distribution



CS 2750 Machine Learning

## Exponential family

- Exponential family of distributions

$$f(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\varphi}) = \exp \left\{ \frac{(\boldsymbol{\theta} \mathbf{x} - b(\boldsymbol{\theta}))}{a(\boldsymbol{\varphi})} + c(\mathbf{x}, \boldsymbol{\varphi}) \right\}$$

- Parameters:

$\boldsymbol{\theta}$  - location parameters

$\boldsymbol{\varphi}$  - scaling parameters

- Example:

- $$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

CS 2750 Machine Learning

## Example: Bernoulli distribution

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

**Objective:**

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

**Probability of an outcome**  $x_i$

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