

# CS 2750 Machine Learning

## Lecture 3

### Density estimation

Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

---

CS 2750 Machine Learning

### Outline

**Outline:**

- **Density estimation:**
  - Maximum likelihood (ML)
  - Bayesian parameter estimates
  - MAP
- **Bernoulli distribution**
- **Binomial distribution**
- **Multinomial distribution**
- **Normal distribution**

---

CS 2750 Machine Learning

## Density estimation

**Density estimation:** is an unsupervised learning

- Learn relations among attributes in the data

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

$D_i = \mathbf{x}_i$  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 attempts to learn the underlying**

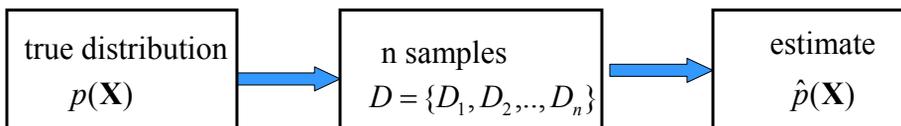
**probability distribution:**  $p(\mathbf{X}) = p(X_1, X_2, \dots, X_d)$

## Density estimation

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

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

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



**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  $\Theta$

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

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

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

## Parameter estimation

- **Maximum likelihood (ML)**

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

- yields: one set of parameters  $\Theta_{ML}$
- the target distribution is approximated as:

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

- **Bayesian parameter estimation**

- uses the posterior distribution over possible parameters

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

- Yields: all possible settings of  $\Theta$  (and their “weights”)
- The target distribution is approximated as:

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

## Parameter estimation

### Other possible 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 = 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  $D = \text{H H T H T H}$   
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  $D = \text{H H T H T H}$   
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 :

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

## 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**  $x_i$

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

## 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

CS 2750 Machine Learning

## 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**

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

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

CS 2750 Machine Learning

## 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 ?

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

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

## 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  $\theta$

### 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**

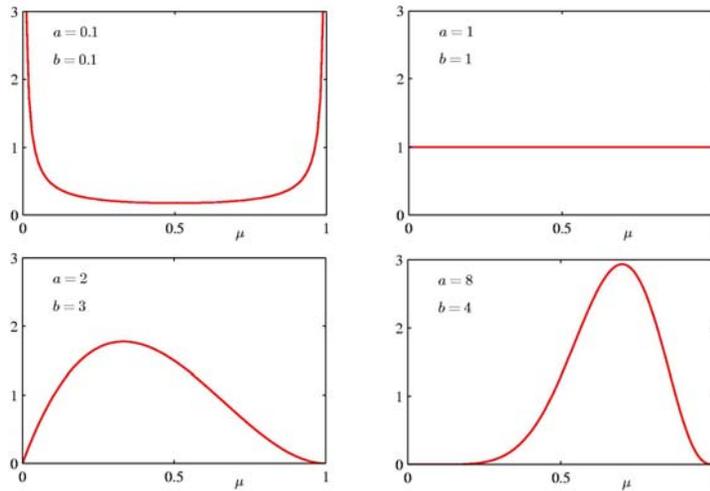
$$P(D | \theta, \xi) = \theta^{N_1} (1 - \theta)^{N_2}$$

### 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

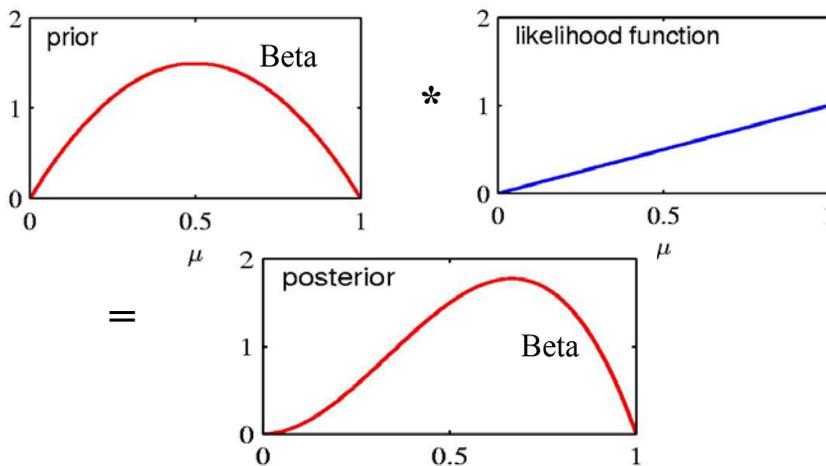
## Beta distribution



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

CS 2750 Machine Learning

## Posterior distribution



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

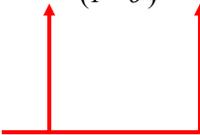
CS 2750 Machine Learning

## 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  $\theta$
- **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?

CS 2750 Machine Learning

## 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 ?

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

## 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

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

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

## Bayesian framework

Both ML or MAP estimates pick one value of the parameter

- **Assume:** there are two different parameter settings that are close in terms of their probability values. Using only one of them may introduce a strong bias, if we use them, for example, for predictions.

**Bayesian parameter estimate**

- Remedies the limitation of one choice
- Keeps all possible parameter values
- Where  $p(\theta | D, \xi) \approx \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$

- **The posterior can be used to define  $p(A | D)$ :**

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

## 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

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

## Expected value of the parameter

### How to obtain the expected value?

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

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

## 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:

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

## 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**

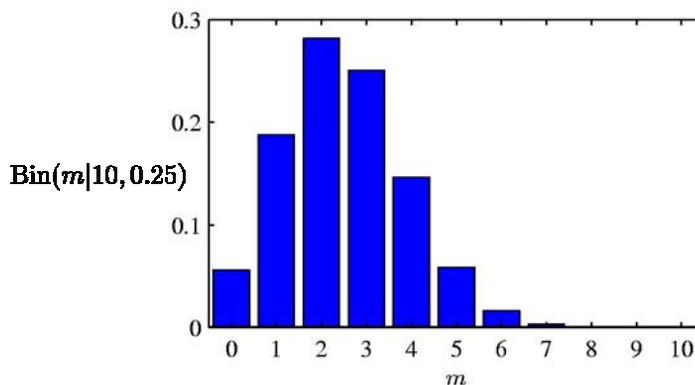
$$P(N_1 | N, \theta) = \binom{N}{N_1} \theta^{N_1} (1 - \theta)^{N - N_1} \quad \text{Binomial distribution}$$

**Objective:**

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

## Binomial distribution

**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

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

## Expected value of the parameter

The result is the same as for Bernoulli distribution

$$E(\theta) = \int_0^1 \theta \text{Beta}(\theta | \eta_1, \eta_2) d\theta = \frac{\eta_1}{\eta_1 + \eta_2}$$

Expected value of the parameter

$$E(\theta) = \frac{\alpha_1 + N_1}{\alpha_1 + N_1 + \alpha_2 + N_2}$$

Predictive probability of event  $x=1$

$$P(x = 1 | \theta, \xi) = E(\theta) = \frac{\alpha_1 + N_1}{\alpha_1 + N_1 + \alpha_2 + N_2}$$