

# CS 1675 Introduction to Machine Learning

## Lecture 7

### Density estimation

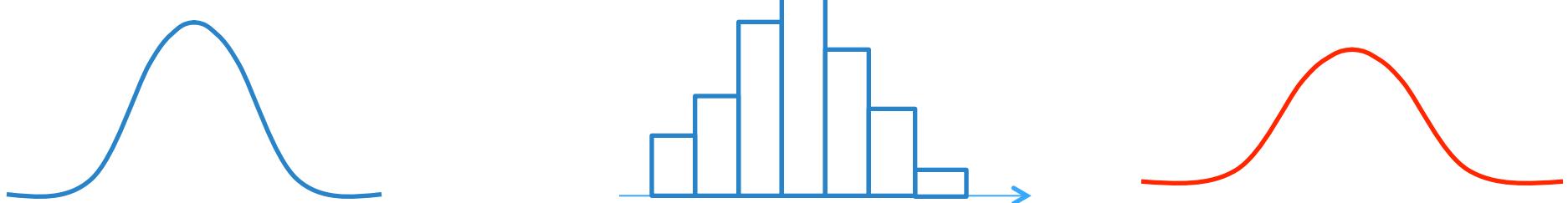
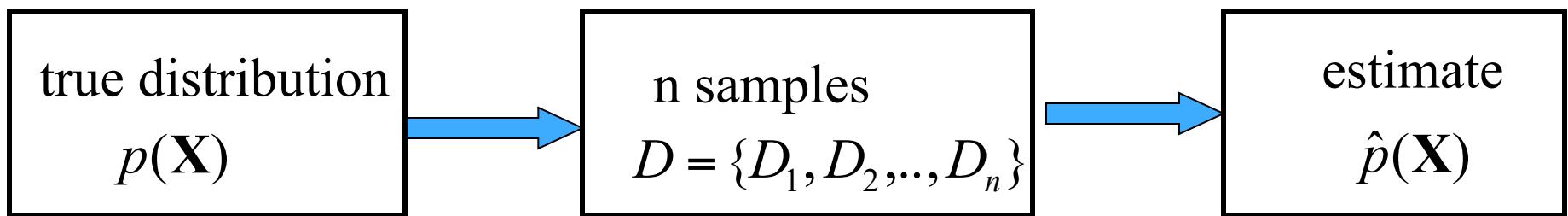
Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

# 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  $D$



# 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  $\Theta$  that maximizes  $p(D | \Theta, \xi)$

---

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

Independent examples

**Log likelihood – has the same maximum as likelihood**

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

---

# 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  $\Theta$  that maximizes  $p(D | \Theta, \xi)$
- 

**Log likelihood – has the same maximum as likelihood**

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

**The optimum satisfies:**

$$\frac{d}{d\Theta} \log p(D | \Theta, \xi) = 0$$

**It can be often solved analytically**

---

# MAP Parameter estimation

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

- **Maximum a posterior probability**  $\max_{\Theta} p(\Theta | D, \xi)$ 
  - Find  $\Theta$  that maximizes  $p(\Theta | D, \xi)$

Likelihood of data



Prior distribution



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

(via Bayes theorem)

$P(D | \xi)$  ← Normalizing factor

Conjugate choices:

- Prior distribution on the parameters **matches** the data distribution
- Posterior is the same type of the distribution as the prior

# MAP Parameter estimation

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

- **Maximum a posterior probability** 
$$\max_{\Theta} p(\Theta | D, \xi)$$
    - Find  $\Theta$  that maximizes  $p(\Theta | D, \xi)$
- 

- **The optimum satisfies:**

$$\frac{d}{d\Theta} \log p(\Theta | D, \xi) = 0 \quad \text{or} \quad \frac{d}{d\Theta} p(\Theta | D, \xi) = 0$$

- **It can be often solved analytically**
-

# Bernoulli distribution

- Model for random variable with two outcomes

**Random variable:**  $x$

**Two outcomes:** 0 or 1

**Distribution:**

$$P(x | \theta) = \theta^x (1 - \theta)^{(1-x)}$$

where  $\theta$  is the probability of  $x=1$

**Example:** Coin toss

**Outcomes:**

- Head  $\rightarrow x=1$
- Tail  $\rightarrow x=0$
- $\theta \rightarrow$  probability of a Head



# Bernoulli distribution

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

**Likelihood of data:**

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

$N_1$  - number of 1s

$N_2$  - number of 0s

**Loglikelihood of data:**

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

**ML estimate:**

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

# 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:**

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

**Posterior:**

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

# Binomial distribution

- models counts of occurrences of binary outcomes in order-independent sequence of trials

**Model:** probability of an outcome 1 (head)  $\theta$   
probability of an outcome 0 (tail)  $(1 - \theta)$

**Probability distribution function:**

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

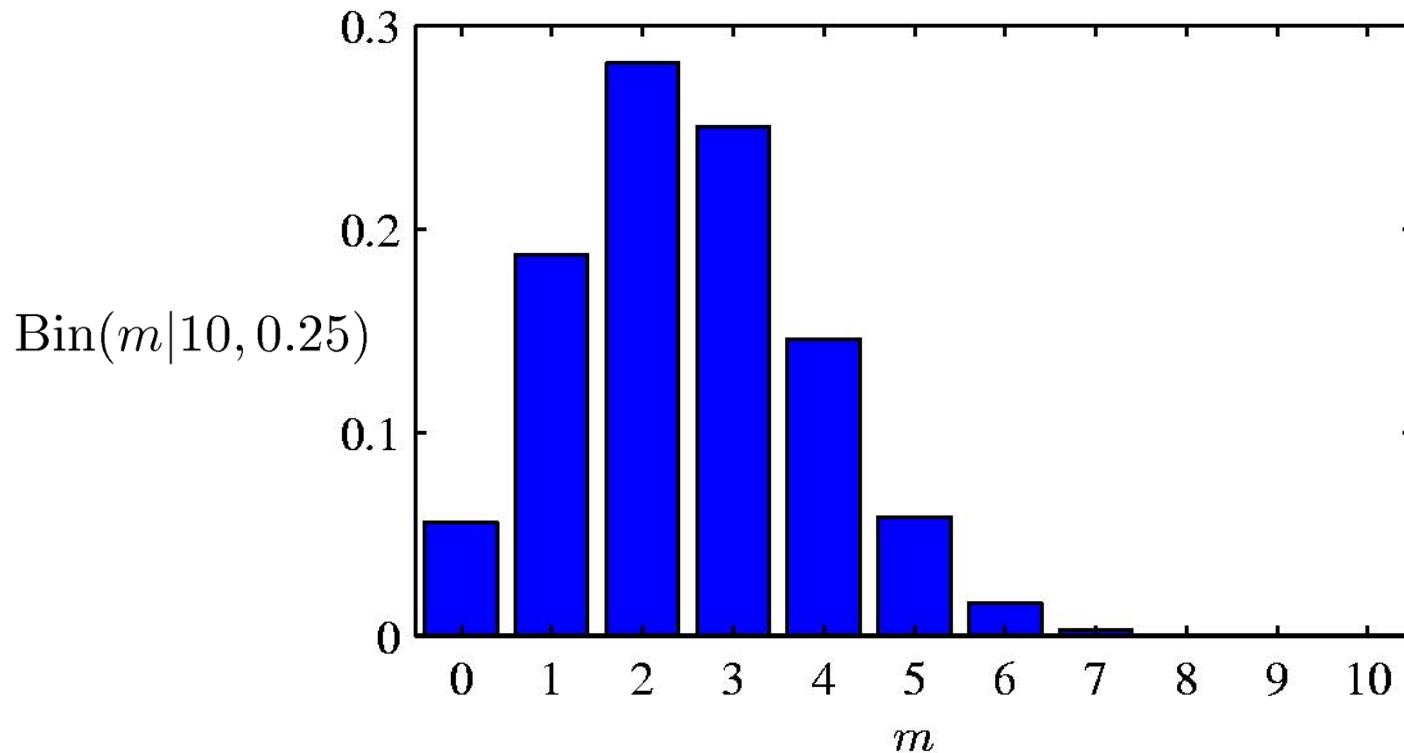
$N_1$  - number of outcomes 1     $N_2$  - number of outcomes 2

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


$$= 2 * \text{penny} + 3 * \text{penny}$$

# 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 !!!

**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

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

## Likelihood

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

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

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

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

# Multinomial distribution



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

**Outcome:** counts of occurrences of  $k$  possible outcomes of  $N$  trials:  $N_i$  - a number of times an outcome  $i$  has been seen

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

**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:**

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

**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:

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

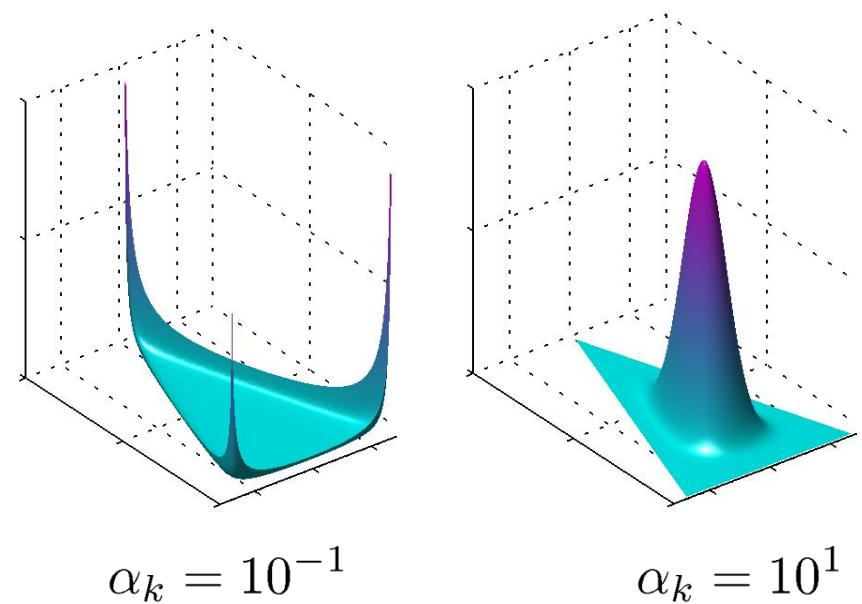
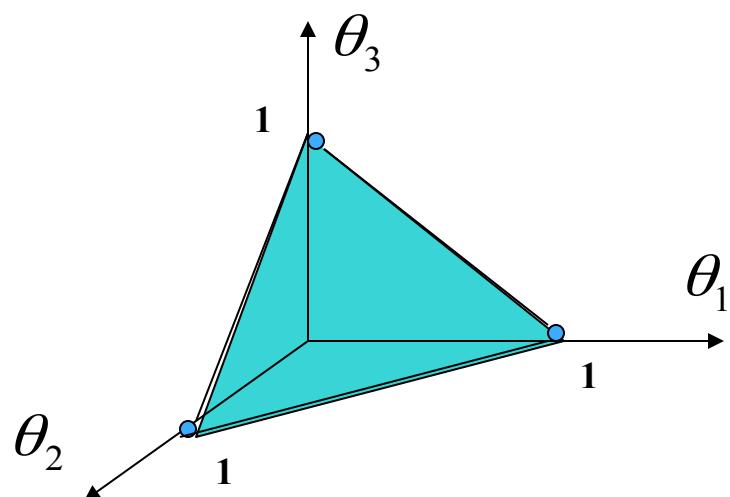
# Dirichlet distribution



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

- **Model of a real-valued outcome**

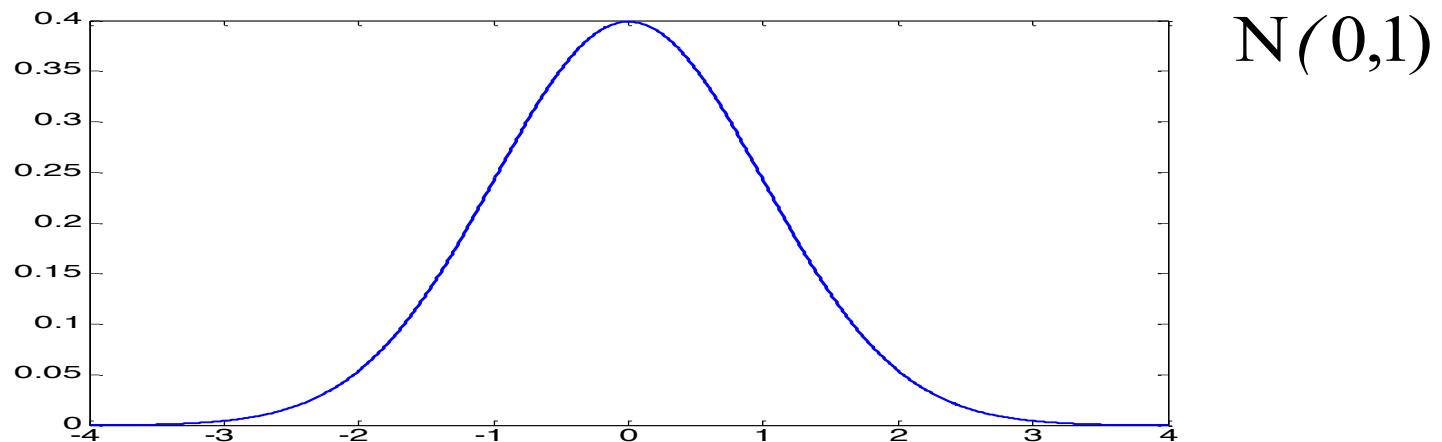
- **Gaussian:**  $x \sim N(\mu, \sigma)$

- **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:**



# Parameter estimates

- **Loglikelihood**

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

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

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

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

- ML variance estimate is biased

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

- **Unbiased estimate:**

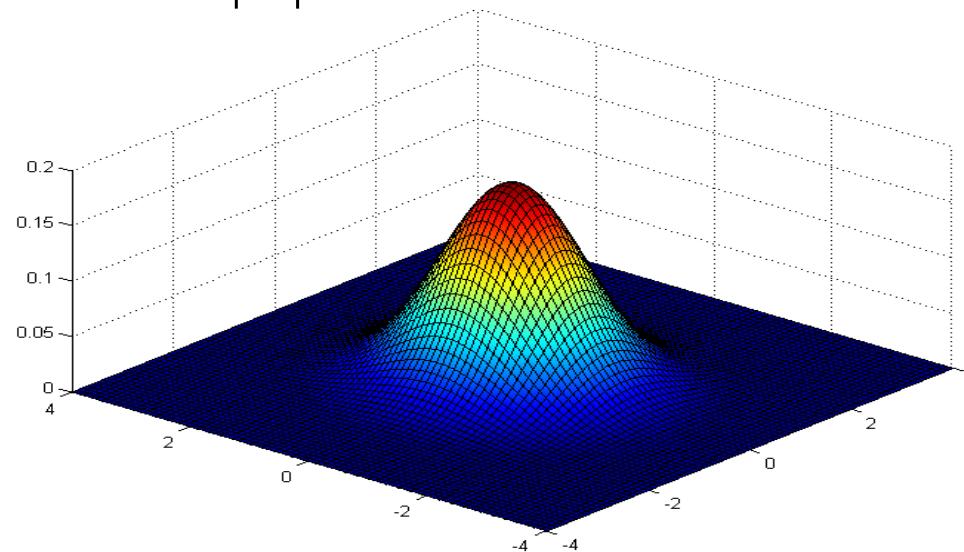
$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

# Multivariate normal distribution

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

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

- **Example:**



# Partitioned Gaussian Distributions

- Multivariate Gaussian:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Example:

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

$$\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$$

$$\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

Precision matrix

- What are the distributions for marginals and conditionals?

$$p(x_a)$$

$$p(x_a | x_b)$$

# Partitioned Conditionals and Marginals

- **Conditional density:**

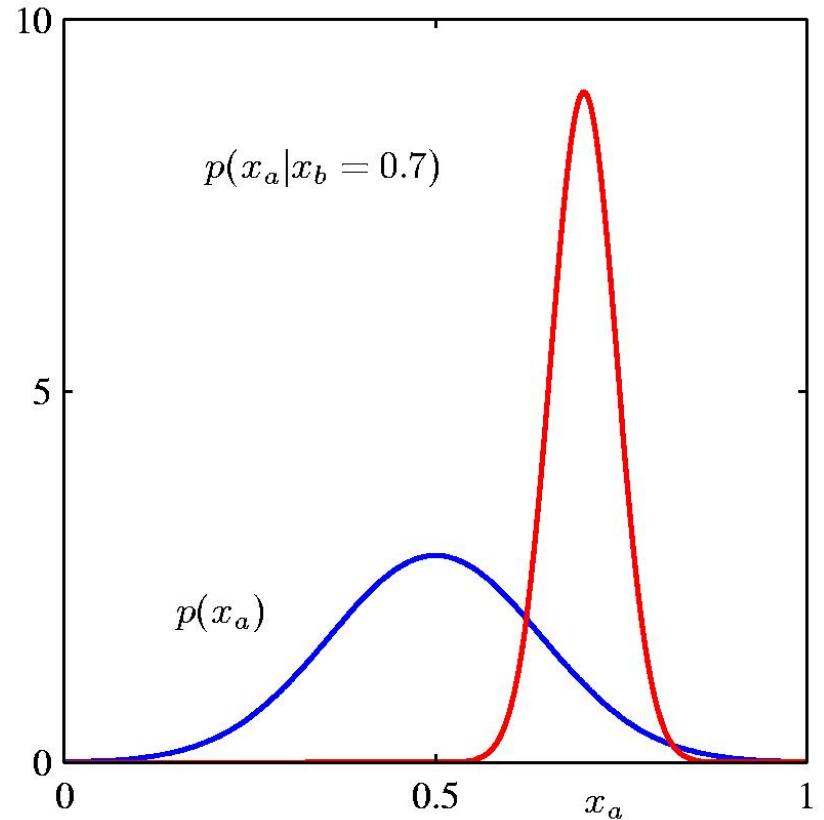
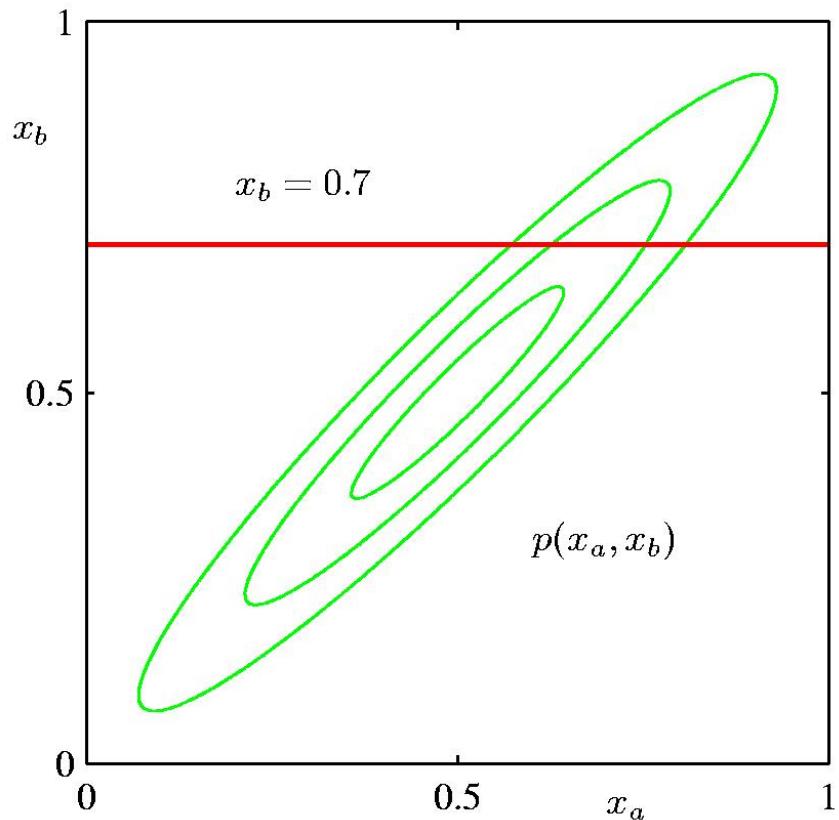
$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$$

$$\begin{aligned}\boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba} \\ \boldsymbol{\mu}_{a|b} &= \boldsymbol{\Sigma}_{a|b} \{ \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \} \\ &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b)\end{aligned}$$

- **Marginal Density:**

$$\begin{aligned}p(\mathbf{x}_a) &= \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b \\ &= \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})\end{aligned}$$

# Partitioned Conditionals and Marginals



# Parameter estimates

- **Loglikelihood**

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

- **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:**

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

# Posterior of a multivariate normal

- Assume a prior on the mean  $\mu$  that is normally distributed:

$$p(\mu) \approx N(\mu_p, \Sigma_p)$$

- Then the posterior of  $\mu$  is normally distributed

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

# Posterior of a multivariate normal

- Then the posterior of  $\mu$  is normally distributed

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

$$\Sigma_n^{-1} = n\Sigma^{-1} + \Sigma_p^{-1}$$

$$\mu_n = \Sigma_p \left( \Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n x_i \right) + \frac{1}{n} \Sigma \left( \Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \mu_p$$

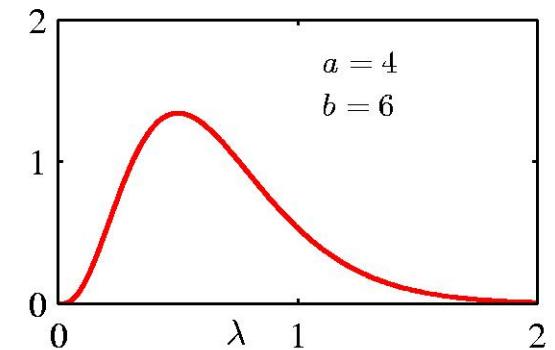
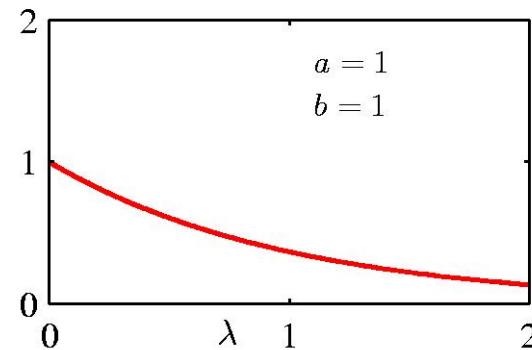
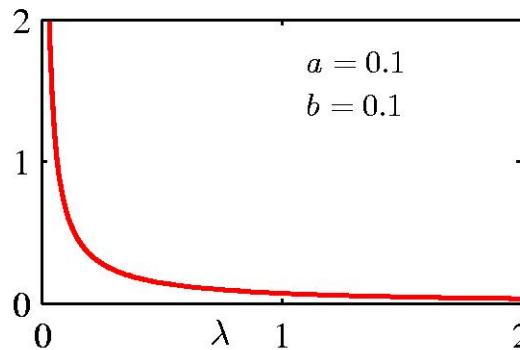
$$\Sigma_n = \Sigma_p \left( \Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \frac{1}{n} \Sigma$$

# Other distributions

**Gamma distribution (models density over non-negative numbers):**

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

$$\mathbb{E}[\lambda] = \frac{a}{b} \qquad \qquad \text{var}[\lambda] = \frac{a}{b^2}$$



# Other distributions

## Exponential distribution:

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

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

## Uniform distribution:

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

**Poisson distribution:** models a number of events occurring in a specific time interval

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

# 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

$$p(D | \Theta) = P(x | \Theta)P(D_{n-1} | \Theta)$$

A “new” prior

- **Then:**

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

# 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:

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

# 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

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

- Parameters

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

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

# 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**

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

- **Parameters**

$$\boldsymbol{\eta} = \begin{bmatrix} \mu / 2\sigma^2 \\ -1 / 2\sigma^2 \end{bmatrix} \quad t(\mathbf{x}) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$

$$Z(\boldsymbol{\eta}) = \exp\left\{\frac{\mu}{2\sigma^2} + \log\sigma\right\} = \exp\left\{-\frac{\eta_1^2}{4\eta_2} - \frac{1}{2}\log(-2\eta_2)\right\}$$

$$h(x) = 1/\sqrt{2\pi}$$

---