

CS 2750 Machine Learning Lecture 4

Density estimation II.

Milos Hauskrecht
milos@cs.pitt.edu
5329 Sennott Square

CS 2750 Machine Learning

Outline

Outline:

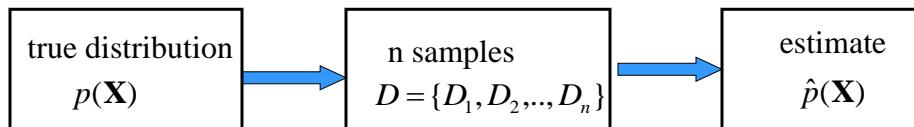
- **Density estimation:**
 - Maximum likelihood (ML)
 - Bayesian parameter estimates
 - MAP
- **Bernoulli distribution.** ✓
- **Binomial distribution**
- **Multinomial distribution**
- **Normal distribution**
- **Exponential family**

CS 2750 Machine Learning

Density estimation

Data: $D = \{D_1, D_2, \dots, D_n\}$
 $D_i = \mathbf{x}_i$ a vector of attribute values

Objective: try to estimate the underlying ‘true’ probability distribution over variables \mathbf{X} , $p(\mathbf{X})$, using examples in D



Standard (iid) assumptions: ?

CS 2750 Machine Learning

Parameter estimation (learning)

- **Maximum likelihood (ML)**
 $\Theta_{ML} = \arg \max_{\Theta} p(D | \Theta, \xi)$
- **Bayesian parameter estimation**
keep the posterior density $p(\Theta | D, \xi)$
- **Maximum a posteriori probability (MAP)**
 $\Theta_{MAP} = \arg \max_{\Theta} p(\Theta | D, \xi)$
- **Expected value**

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

CS 2750 Machine Learning

Bernoulli trials

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 $(1-\theta)$

Probability of an outcome of a coin flip

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

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

N_1, N_2 - Number of heads and tails respectively

CS 2750 Machine Learning

Posterior distribution

Posterior density

Likelihood of data \downarrow prior \downarrow

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

Normalizing factor

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

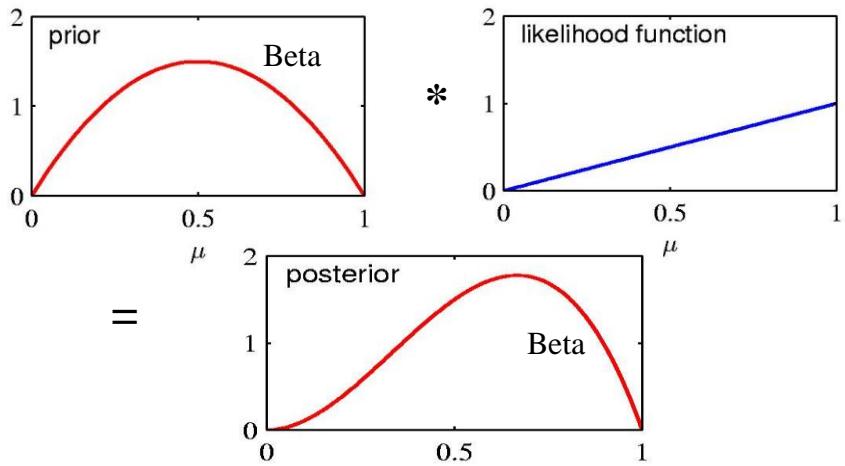
$p(\theta | \xi)$ - is the prior probability on θ

Conjugate choice of prior: Beta

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

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

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

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

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

Bayesian framework

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

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

Posterior density

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

- 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 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 α

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: a biased coin

Outcomes: two possible values -- head or tail

Data: D a set of order-independent outcomes

We treat D as a multi-set !!!

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

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

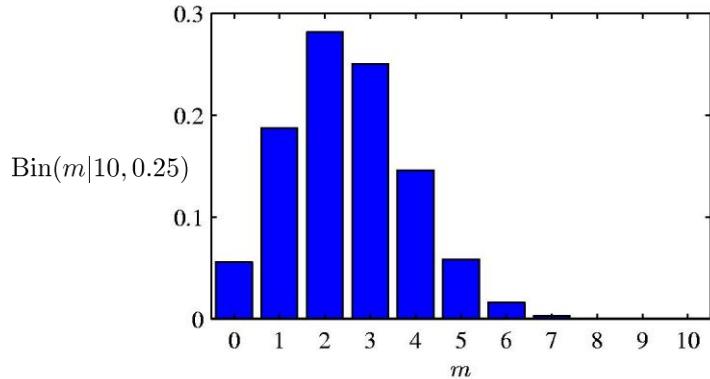
Probability of an outcome

$$P(D | \theta) = \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N_2} \quad \text{Binomial distribution}$$

CS 2750 Machine Learning

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 an 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

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

Likelihood

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

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

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

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

CS 2750 Machine Learning

Expected value of the parameter

The result is the same as for Bernoulli distribution

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

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

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

CS 2750 Machine Learning

Binomial distribution

Example: a biased coin

Outcomes: two possible values -- head or tail

Data: D a set of order-independent outcomes

We treat D as a multi-set !!!

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

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

Probability of an outcome

$$P(D | \theta) = \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N_2} \quad \text{Binomial distribution}$$

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

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

CS 2750 Machine Learning

Multinomial distribution

Example: Multi-way coin toss, roll of dice

- Data:** a set of N outcomes (multi-set)

N_i - a number of times an outcome i has been seen

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

Probability of data (likelihood)

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

ML estimate:

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

CS 2750 Machine Learning

Posterior density and MAP estimate

Choice of the prior: Dirichlet distribution

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

Dirichlet is the **conjugate choice for multinomial**

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

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

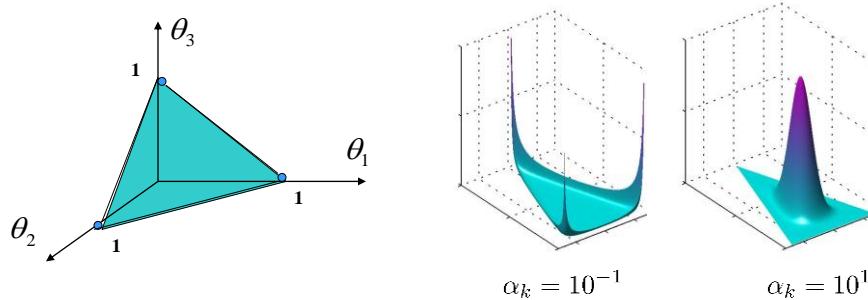
CS 2750 Machine Learning

Dirichlet distribution

Dirichlet distribution:

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

Assume: $k=3$



CS 2750 Machine Learning

Expected value

The result is analogous to the result for binomial

$$E(\boldsymbol{\theta}) = \int_{\sum \theta_i \leq 1, \sum \theta_i = 1} \boldsymbol{\theta} Dir(\boldsymbol{\theta} | \boldsymbol{\eta}) d\boldsymbol{\theta} = \left(\frac{\eta_1}{\eta_1 + \eta_2 + \eta_k}, \dots, \frac{\eta_i}{\eta_1 + \eta_2 + \eta_k}, \dots, \frac{\eta_k}{\eta_1 + \eta_2 + \eta_k} \right)$$

Expectation based parameter estimate

$$E(\boldsymbol{\theta}) = \left(\frac{\alpha_1 + N_1}{\alpha_1 + N_1 + \dots + \alpha_k + N_k}, \dots, \frac{\alpha_i + N_i}{\alpha_1 + N_1 + \dots + \alpha_k + N_k}, \dots, \frac{\alpha_k + N_k}{\alpha_1 + N_1 + \dots + \alpha_k + N_k} \right)$$

Represents the predictive probability of an event $x=i$

$$P(x=i | \boldsymbol{\theta}, \xi) = \frac{\alpha_i + N_i}{\alpha_1 + N_1 + \dots + \alpha_k + N_k}$$

CS 2750 Machine Learning

Other distributions

The same ideas can be applied to other distributions

- Typically we choose distributions that behave well so that computations lead to a 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)

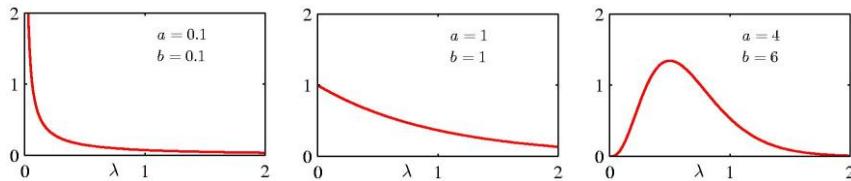
CS 2750 Machine Learning

Gamma distribution

- **Gamma distribution**

$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$

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



CS 2750 Machine Learning

Other distributions

Exponential distribution:

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

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

Poisson distribution:

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

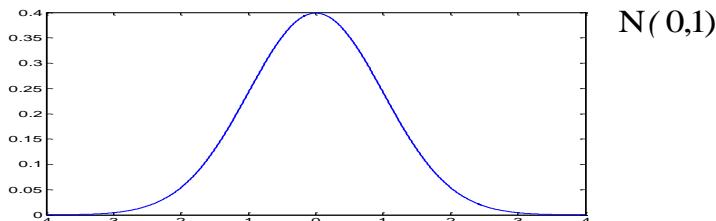
CS 2750 Machine Learning

Gaussian (normal) distribution

- **Gaussian:** $x \sim N(\mu, \sigma)$
- **Parameters:** μ - mean
 σ - 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

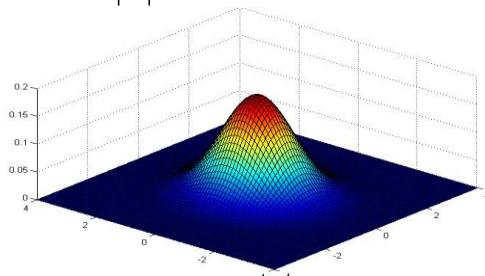
- Multivariate normal: $\mathbf{x} \sim N(\boldsymbol{\mu}, \Sigma)$

• Parameters: $\boldsymbol{\mu}$ - mean
 Σ - 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:



CS 2750 Machine Learning

Partitioned Gaussian Distributions

- Multivariate Gaussian:

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

- what are marginals and conditionals?

- Example:

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

$$\Lambda \equiv \Sigma^{-1} \quad \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

Precision matrix

CS 2750 Machine Learning

Partitioned Conditionals and Marginals

- Conditional density:

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

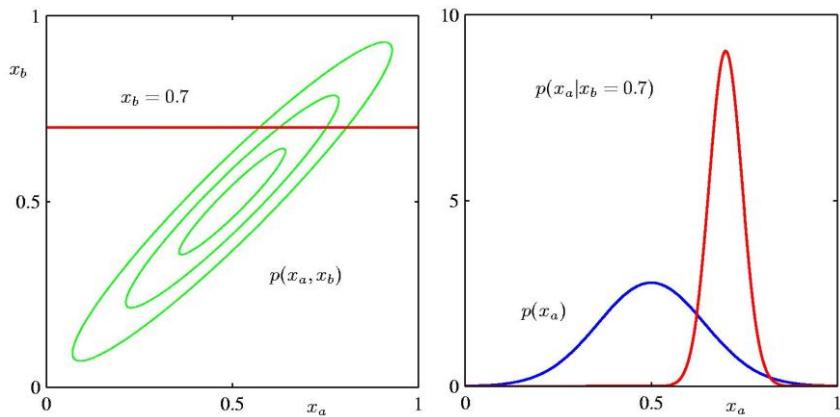
$$\begin{aligned}\boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba} \\ \boldsymbol{\mu}_{a|b} &= \boldsymbol{\Sigma}_{a|b} \{ \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}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

- Marginal Density:

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

CS 2750 Machine Learning

Partitioned Conditionals and Marginals



CS 2750 Machine Learning

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:**

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

CS 2750 Machine Learning