

RWTH AACHEN  
UNIVERSITY

# Advanced Machine Learning Lecture 6

## Probability Distributions

16.11.2015

Bastian Leibe  
RWTH Aachen  
<http://www.vision.rwth-aachen.de/>  
leibe@vision.rwth-aachen.de

Advanced Machine Learning Winter'15

RWTH AACHEN  
UNIVERSITY

## This Lecture: Advanced Machine Learning

- Regression Approaches
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Gaussian Processes
- Learning with Latent Variables
  - Probability Distributions & Mixture Models
  - Approximate Inference
  - EM and Generalizations
- Deep Learning
  - Neural Networks
  - CNNs, RNNs, RBMs, etc.

B. Leibe

RWTH AACHEN  
UNIVERSITY

## Recap: GPs with Noise-free Observations

- Assume our observations are noise-free:
 
$$\{(x_n, f_n) \mid n = 1, \dots, N\}$$
  - Joint distribution of the training outputs  $\mathbf{f}$  and test outputs  $\mathbf{f}_*$  according to the prior:
 
$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right)$$
  - Calculation of posterior corresponds to conditioning the joint Gaussian prior distribution on the observations:
 
$$f_* | X_*, X, \mathbf{f} \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}[\mathbf{f}_*]) \quad \bar{\mathbf{f}}_* = \mathbb{E}[\mathbf{f}_* | X, X_*, \mathbf{f}]$$
  - with:
 
$$\bar{\mathbf{f}}_* = K(X_*, X) K(X, X)^{-1} \mathbf{f}$$

$$\text{cov}[\mathbf{f}_*] = K(X_*, X_*) - K(X_*, X) K(X, X)^{-1} K(X, X_*)$$

4

Slide adapted from Bernt Schiele B. Leibe

RWTH AACHEN  
UNIVERSITY

## Recap: GPs with Noisy Observations

- Joint distribution of the observed values and the test locations under the prior:
 
$$\begin{bmatrix} \mathbf{t} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right)$$
  - Calculation of posterior corresponds to conditioning the joint Gaussian prior distribution on the observations:
 
$$f_* | X_*, X, \mathbf{t} \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}[\mathbf{f}_*]) \quad \bar{\mathbf{f}}_* = \mathbb{E}[\mathbf{f}_* | X, X_*, \mathbf{t}]$$
  - with:
 
$$\bar{\mathbf{f}}_* = K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} \mathbf{t}$$

$$\text{cov}[\mathbf{f}_*] = K(X_*, X_*) - K(X_*, X) (K(X, X) + \sigma_n^2 I)^{-1} K(X, X_*)$$
- ⇒ This is the key result that defines Gaussian process regression!
  - Predictive distribution is Gaussian whose mean and variance depend on test points  $X_*$  and on the kernel  $k(x, x')$ , evaluated on  $X$ .

5

Slide adapted from Bernt Schiele B. Leibe

RWTH AACHEN  
UNIVERSITY

## Recap: Bayesian Model Selection for GPs

- Goal
  - Determine/learn different parameters of Gaussian Processes
- Hierarchy of parameters
  - Lowest level
    - w - e.g. parameters of a linear model.
  - Mid-level (hyperparameters)
    - $\theta$  - e.g. controlling prior distribution of w.
  - Top level
    - Typically discrete set of model structures  $\mathcal{H}_i$ .
- Approach
  - Inference takes place one level at a time.

6

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN  
UNIVERSITY

## Recap: Model Selection at Lowest Level

- Posterior of the parameters  $\mathbf{w}$  is given by Bayes' rule
 
$$p(\mathbf{w} | \mathbf{t}, X, \theta, \mathcal{H}_i) = \frac{p(\mathbf{t} | X, \mathbf{w}, \theta, \mathcal{H}_i) p(\mathbf{w} | \theta, X, \mathcal{H}_i)}{p(\mathbf{t} | X, \theta, \mathcal{H}_i)}$$

$$= \frac{p(\mathbf{t} | X, \mathbf{w}, \theta, \mathcal{H}_i) p(\mathbf{w} | \theta, \mathcal{H}_i)}{p(\mathbf{t} | X, \theta, \mathcal{H}_i)}$$
- with
  - $p(\mathbf{t} | X, \mathbf{w}, \theta, \mathcal{H}_i)$  likelihood and
  - $p(\mathbf{w} | \theta, \mathcal{H}_i)$  prior parameters w,
  - Denominator (normalizing constant) is independent of the parameters and is called **marginal likelihood**.
$$p(\mathbf{t} | X, \theta, \mathcal{H}_i) = \int p(\mathbf{t} | X, \mathbf{w}, \theta, \mathcal{H}_i) p(\mathbf{w} | \theta, \mathcal{H}_i) d\mathbf{w}$$

7

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

## Recap: Model Selection at Mid Level

- Posterior of parameters  $\theta$  is again given by Bayes' rule

$$p(\theta|t, X, \mathcal{H}_i) = \frac{p(t|X, \theta, \mathcal{H}_i)p(\theta|X, \mathcal{H}_i)}{p(t|X, \mathcal{H}_i)}$$

$$= \frac{p(t|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i)}{p(t|X, \mathcal{H}_i)}$$

- where
  - The marginal likelihood of the previous level  $p(t|X, \theta, \mathcal{H}_i)$  plays the role of the likelihood of this level.
  - $p(\theta|\mathcal{H}_i)$  is the **hyperprior** (prior of the hyperparameters)
  - Denominator (normalizing constant) is given by:

$$p(t|X, \mathcal{H}_i) = \int p(t|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i)d\theta$$

Slide credit: Bernt Schiele B. Leibe 8

RWTH AACHEN UNIVERSITY

## Recap: Model Selection at Top Level

- At the top level, we calculate the posterior of the model

$$p(\mathcal{H}_i|t, X) = \frac{p(t|X, \mathcal{H}_i)p(\mathcal{H}_i)}{p(t|X)}$$

- where
  - Again, the denominator of the previous level  $p(t|X, \mathcal{H}_i)$  plays the role of the likelihood.
  - $p(\mathcal{H}_i)$  is the prior of the model structure.
  - Denominator (normalizing constant) is given by:

$$p(t|X) = \sum_i p(t|X, \mathcal{H}_i)p(\mathcal{H}_i)$$

Slide credit: Bernt Schiele B. Leibe 9

RWTH AACHEN UNIVERSITY

## Recap: Bayesian Model Selection

- Discussion
  - Marginal likelihood is main difference to non-Bayesian methods
  - It automatically incorporates a trade-off between the model fit and the model complexity:
    - A simple model can only account for a limited range of possible sets of target values - if a simple model fits well, it obtains a high posterior.
    - A complex model can account for a large range of possible sets of target values - therefore, it can never attain a very high posterior.

Slide credit: Bernt Schiele B. Leibe Image source: Rasmussen & Williams, 2006 10

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Probability Distributions
  - Bayesian Estimation Reloaded
- Binary Variables
  - Bernoulli distribution
  - Binomial distribution
  - Beta distribution
- Multinomial Variables
  - Multinomial distribution
  - Dirichlet distribution
- Continuous Variables
  - Gaussian distribution
  - Gamma distribution
  - Student's t distribution
  - Exponential Family

Slide credit: Bernt Schiele B. Leibe 11

RWTH AACHEN UNIVERSITY

## Motivation

- Recall: Bayesian estimation

$$p(x|X) = \int p(x|\theta) \frac{p(X|\theta)p(\theta)}{\int p(X|\theta')p(\theta')d\theta'} d\theta$$

- So far, we have only done this for Gaussian distributions, where the integrals could be solved analytically.
- Now, let's also examine other distributions...

Slide credit: Bernt Schiele B. Leibe Image created with Tagxedo.com 12

RWTH AACHEN UNIVERSITY

## Teaser: Conjugate Priors

- Problem: How to evaluate the integrals?
  - We will see that if likelihood and prior have the same functional form  $c \cdot f(x)$ , then the analysis will be greatly simplified and the integrals will be solvable in closed form.
 
$$p(X|\theta)p(\theta) = \prod_{x_n} c_1 f(x_n, \theta) c_2 f(\theta, \alpha)$$

$$= \prod_{x_n} c f(x_n, \theta, \alpha)$$
  - Such an algebraically convenient choice is called a **conjugate prior**. Whenever possible, we should use it.
  - To do this, we need to know for each probability distribution what is its conjugate prior.  $\Rightarrow$  *Topic of this lecture.*
- What to do when we cannot use the conjugate prior?
  - $\Rightarrow$  Use approximate inference methods. Next lecture...

Slide credit: Bernt Schiele B. Leibe 13

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Probability Distributions
  - Bayesian Estimation Reloaded
- Binary Variables
  - Bernoulli distribution
  - Binomial distribution
  - Beta distribution
- Multinomial Variables
  - Multinomial distribution
  - Dirichlet distribution
- Continuous Variables
  - Gaussian distribution
  - Gamma distribution
  - Student's t distribution
  - Exponential Family

B. Leibe 14

RWTH AACHEN UNIVERSITY

## Binary Variables

- Example: Flipping a coin
  - Binary random variable  $x \in \{0,1\}$
  - Outcome heads:  $x = 1$
  - Outcome tails:  $x = 0$
  - Denote probability of landing heads by parameter  $\mu$ 

$$p(x = 1|\mu) = \mu$$
- Bernoulli distribution
  - Probability distribution over  $x$ :
 
$$\text{Bern}(x|\mu) = \mu^x(1-\mu)^{1-x}$$

$$\mathbb{E}[x] = \mu$$

$$\text{var}[x] = \mu(1-\mu)$$

B. Leibe 15

RWTH AACHEN UNIVERSITY

## The Binomial Distribution

- Now consider  $N$  coin flips
  - Probability of landing  $m$  heads:  $p(m \text{ heads}|N, \mu)$
- Binomial distribution
 
$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m}$$
  - Properties
 
$$\mathbb{E}[m] \equiv \sum_{m=0}^N m \text{Bin}(m|N, \mu) = N\mu$$

$$\text{var}[m] \equiv \sum_{m=0}^N (m - \mathbb{E}[m])^2 \text{Bin}(m|N, \mu) = N\mu(1-\mu)$$
  - Note: Bernoulli is a special case of the Binomial for  $n = 1$ .

B. Leibe 16

RWTH AACHEN UNIVERSITY

## Binomial Distribution: Visualization

B. Leibe 17

RWTH AACHEN UNIVERSITY

## Parameter Estimation: Maximum Likelihood

- Maximum Likelihood for Bernoulli
  - Given a data set  $\mathcal{D} = \{x_1, \dots, x_N\}$  of observed values for  $x$ .
  - Likelihood
 
$$p(\mathcal{D}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n}$$

$$\log p(\mathcal{D}|\mu) = \sum_{n=1}^N \log p(x_n|\mu) = \sum_{n=1}^N \{x_n \log \mu + (1-x_n) \log(1-\mu)\}$$
  - Observation
    - The log-likelihood depends on the observations  $x_n$  only through their sum.
    - ➔  $\sum_n x_n$  is a **sufficient statistic** for the Bernoulli distribution.

B. Leibe 18

RWTH AACHEN UNIVERSITY

## ML for Bernoulli Distribution

$$\log p(\mathcal{D}|\mu) = \sum_{n=1}^N \{x_n \log \mu + (1-x_n) \log(1-\mu)\}$$

$$\nabla_{\mu} \log p(\mathcal{D}|\mu) = \frac{1}{\mu} \sum_{n=1}^N x_n - \frac{1}{1-\mu} \sum_{n=1}^N (1-x_n) \stackrel{!}{=} 0$$

$$(1-\mu) \sum_{n=1}^N x_n = \mu \sum_{n=1}^N (1-x_n)$$

$$\sum_{n=1}^N x_n - \cancel{\mu \sum_{n=1}^N x_n} = N\mu - \cancel{\mu \sum_{n=1}^N x_n}$$

- ML estimate:
 
$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$

B. Leibe 19

RWTH AACHEN UNIVERSITY

## ML for Bernoulli Distribution

- Maximum Likelihood estimate
 
$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n = \frac{m}{N} \quad \text{for } m \text{ heads } (x_n = 1)$$
- Discussion
  - Consider a data set  $\mathcal{D} = \{1,1,1\}$ .  $\rightarrow \mu_{ML} = \frac{3}{3} = 1$
  - $\Rightarrow$  Prediction: *all* future tosses will land head up!
  - $\Rightarrow$  Overfitting to  $\mathcal{D}$ !

Advanced Machine Learning Winter'12 | Slide adapted from C. Bishop | B. Leibe | 20

RWTH AACHEN UNIVERSITY

## Bayesian Bernoulli: First Try

- Bayesian estimation
  - We can improve the ML estimate by incorporating a prior for  $\mu$ .
  - How should such a prior look like?
  - Consider the Bernoulli/Binomial form
 
$$p(\mathcal{D}|\mu) \propto \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n}$$
  - If we choose a prior with the same functional form, then we will get a closed-form expression for the posterior; otherwise, a difficult numerical integration may be necessary.
  - Most general form here:
 
$$p(\mu|a, b) \propto \mu^a (1-\mu)^b$$
  - This algebraically convenient choice is called a **conjugate prior**.

Advanced Machine Learning Winter'12 | B. Leibe | 21

RWTH AACHEN UNIVERSITY

## The Beta Distribution

- Beta distribution
  - Distribution over  $\mu \in [0,1]$ :
 
$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$
  - Where  $\Gamma(x)$  is the **gamma function**

$$\Gamma(x) \equiv \int_0^\infty u^{x-1} e^{-u} du$$
 for which  $\Gamma(x+1) = x!$  iff  $x$  is an integer.
  $\Rightarrow \Gamma(x)$  is a continuous generalization of the factorial.
  - The Beta distribution generalizes the Binomial to arbitrary values of  $a$  and  $b$ , while keeping the same functional form.
  - It is therefore a **conjugate prior** for the Bernoulli and Binomial.

Advanced Machine Learning Winter'12 | B. Leibe | 22

RWTH AACHEN UNIVERSITY

## Beta Distribution

- Properties
  - In general, the Beta distribution is a suitable model for the random behavior of percentages and proportions.
  - Mean and variance
 
$$\mathbb{E}[\mu] = \frac{a}{a+b}$$

$$\text{var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)}$$
  - The parameters  $a$  and  $b$  are often called **hyperparameters**, because they control the distribution of the parameter  $\mu$ .
  - General observation: if a distribution has  $K$  parameters, then the conjugate prior typically has  $K+1$  hyperparameters.

Advanced Machine Learning Winter'12 | B. Leibe | 23

RWTH AACHEN UNIVERSITY

## Beta Distribution: Visualization

Advanced Machine Learning Winter'12 | Slide credit: C. Bishop | B. Leibe | Image source: C. Bishop, 2008 | 24

RWTH AACHEN UNIVERSITY

## Bayesian Bernoulli

- Bayesian estimate
 
$$p(\mu|a_0, b_0, \mathcal{D}) \propto p(\mathcal{D}|\mu)p(\mu|a_0, b_0)$$

$$= \left( \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n} \right) \text{Beta}(\mu|a_0, b_0)$$

$$\propto \mu^{m+a_0-1} (1-\mu)^{(N-m)+b_0-1}$$

$$\propto \text{Beta}(\mu|a_N, b_N)$$
- This is again a Beta distribution with the parameters
 
$$a_N = a_0 + m \quad b_N = b_0 + (N - m)$$
- $\Rightarrow$  We can interpret the hyperparameters  $a$  and  $b$  as an **effective number of observations** for  $x = 1$  and  $x = 0$ , respectively.
- Note:  $a$  and  $b$  need not be integers!

Advanced Machine Learning Winter'12 | Slide adapted from C. Bishop | B. Leibe | 25

RWTH AACHEN UNIVERSITY

## Sequential Estimation

- Prior · Likelihood = Posterior
  - The posterior can act as a prior if we observe additional data.
  - The number of effective observations increases accordingly.
- Example: Taking observations one at a time

Beta( $\mu|a = 2, b = 2$ )   Bin( $m = 1|N = 1, \mu$ )   Beta( $\mu|a = 3, b = 2$ )

⇒ This sequential approach to learning naturally arises when we take a Bayesian viewpoint.

B. Leibe 26  
Image source: C. Bishop, 2004

RWTH AACHEN UNIVERSITY

## Properties of the Posterior

- Behavior in the limit of infinite data
  - As the size of the data set,  $N$ , increases

$$a_N = a_0 + m \rightarrow m$$

$$b_N = b_0 + N - m \rightarrow N - m$$

$$\mathbb{E}[\mu] = \frac{a_N}{a_N + b_N} \rightarrow \frac{m}{N} = \mu_{ML}$$

$$\text{var}[\mu] = \frac{a_N b_N}{(a_N + b_N)^2 (a_N + b_N + 1)} \rightarrow 0$$

⇒ As expected, the Bayesian result reduces to the ML result.

Advanced Machine Learning Winter'12 27  
Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## Prediction under the Posterior

- Predict the outcome of the next trial
  - “What is the probability that the next coin toss will land heads up?”
- ⇒ Evaluate the predictive distribution of  $x$  given the observed data set  $\mathcal{D}$ :

$$p(x = 1|a_0, b_0, \mathcal{D}) = \int_0^1 p(x = 1|\mu)p(\mu|a_0, b_0, \mathcal{D}) d\mu$$

$$= \int_0^1 \mu p(\mu|a_0, b_0, \mathcal{D}) d\mu$$

$$= \mathbb{E}[\mu|a_0, b_0, \mathcal{D}] = \frac{a_N}{a_N + b_N}$$

- Simple interpretation: total fraction of observations that correspond to  $x = 1$ .

Advanced Machine Learning Winter'12 28  
Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Probability Distributions
  - Bayesian Estimation Reloaded
- Binary Variables
  - Bernoulli distribution
  - Binomial distribution
  - Beta distribution
- Multinomial Variables
  - Multinomial distribution
  - Dirichlet distribution
- Continuous Variables
  - Gaussian distribution
  - Gamma distribution
  - Student's t distribution
  - Exponential Family

Advanced Machine Learning Winter'12 29  
B. Leibe

RWTH AACHEN UNIVERSITY

## Multinomial Variables

- Multinomial variables
  - Variables that can take one of  $K$  possible distinct states
  - Convenient: 1-of- $K$  coding scheme:  $\mathbf{x} = (0, 0, 1, 0, 0, 0)^T$
- Generalization of the Bernoulli distribution
  - Distribution of  $\mathbf{x}$  with  $K$  outcomes

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k}$$

with the constraints

$$\forall k : \mu_k \geq 0 \quad \text{and} \quad \sum_{k=1}^K \mu_k = 1$$

Advanced Machine Learning Winter'12 30  
Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## Multinomial Variables

- Properties
  - Distribution is normalized
 
$$\sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) = \sum_{k=1}^K \mu_k = 1$$
  - Expectation
 
$$\mathbb{E}[\mathbf{x}|\boldsymbol{\mu}] = \sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) \mathbf{x} = (\mu_1, \dots, \mu_K)^T = \boldsymbol{\mu}$$
  - Likelihood given a data set  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ :
 
$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})} = \prod_{k=1}^K \mu_k^{m_k}$$

where  $m_k$  is the number of cases for which  $\mathbf{x}_n$  has output  $k$ .

Advanced Machine Learning Winter'12 31  
Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## ML Parameter Estimation

- Maximum Likelihood solution for  $\mu$ 
  - Need to maximize
 
$$\log p(\mathcal{D}|\mu) = \log \prod_{k=1}^K \mu_k^{m_k} = \sum_{k=1}^K m_k \log \mu_k$$
  - Under the constraint  $\sum_k \mu_k = 1$
- Solution with Lagrange multiplier
 
$$\arg \max_{\mu} \sum_{k=1}^K m_k \log \mu_k + \lambda \left( \sum_{k=1}^K \mu_k - 1 \right)$$
  - Setting the derivative to zero yields
 
$$\mu_k = -m_k/\lambda \quad \mu_k^{\text{ML}} = \frac{m_k}{N}$$

Advanced Machine Learning Winter'12 | Slide adapted from C. Bishop | B. Leibe | 32

RWTH AACHEN UNIVERSITY

## The Multinomial Distribution

- Multinomial Distribution
  - Joint distribution over  $m_1, \dots, m_K$  conditioned on  $\mu$  and  $N$ 

$$\text{Mult}(m_1, m_2, \dots, m_K | \mu, N) = \binom{N}{m_1 m_2 \dots m_K} \prod_{k=1}^K \mu_k^{m_k}$$
  - with the normalization coefficient
 
$$\binom{N}{m_1 m_2 \dots m_K} = \frac{N!}{m_1! m_2! \dots m_K!}$$
  - Properties
 
$$\begin{aligned} \mathbb{E}[m_k] &= N \mu_k \\ \text{var}[m_k] &= N \mu_k (1 - \mu_k) \\ \text{cov}[m_j, m_k] &= -N \mu_j \mu_k \end{aligned}$$

Advanced Machine Learning Winter'12 | Slide adapted from C. Bishop | B. Leibe | 33

RWTH AACHEN UNIVERSITY

## Bayesian Multinomial

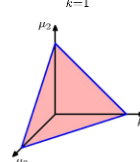
- Conjugate prior for the Multinomial
  - Introduce a family of prior distributions for the parameters  $\{\mu_k\}$  of the Multinomial.
  - The conjugate prior is given by
 
$$p(\mu|\alpha) \propto \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$
  - with the constraints
 
$$\forall k : 0 \leq \mu_k \leq 1 \quad \text{and} \quad \sum_{k=1}^K \mu_k = 1$$

Advanced Machine Learning Winter'12 | B. Leibe | 34

RWTH AACHEN UNIVERSITY

## The Dirichlet Distribution

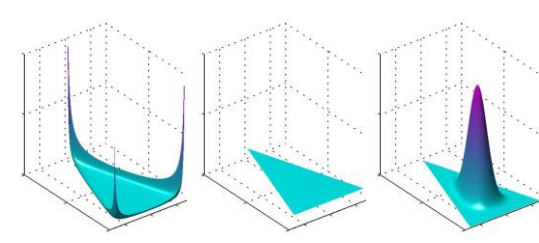
- Dirichlet Distribution
  - Multivariate generalization of the Beta distribution
 
$$\text{Dir}(\mu|\alpha) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1} \quad \text{with} \quad \alpha_0 = \sum_{k=1}^K \alpha_k$$
  - Properties
    - The Dirichlet distribution over  $K$  variables is confined to a  $K-1$  dimensional simplex.
    - Expectations:
 
$$\begin{aligned} \mathbb{E}[\mu_k] &= \frac{\alpha_k}{\alpha_0} \\ \text{var}[\mu_k] &= \frac{\alpha_k (\alpha_0 - \alpha_k)}{\alpha_0^2 (\alpha_0 + 1)} \\ \text{cov}[\mu_j, \mu_k] &= -\frac{\alpha_j \alpha_k}{\alpha_0^3 (\alpha_0 + 1)} \end{aligned}$$



Advanced Machine Learning Winter'12 | Slide adapted from C. Bishop | B. Leibe | Image source: C. Bishop, 2009 | 35

RWTH AACHEN UNIVERSITY

## Dirichlet Distribution: Visualization



$\alpha_k = 10^{-1}$        $\alpha_k = 10^0$        $\alpha_k = 10^1$

Advanced Machine Learning Winter'12 | Slide credit: C. Bishop | B. Leibe | Image source: C. Bishop, 2009 | 36

RWTH AACHEN UNIVERSITY

## Bayesian Multinomial

- Posterior distribution over the parameters  $\{\mu_k\}$ 

$$p(\mu|\mathcal{D}, \alpha) \propto p(\mathcal{D}|\mu) p(\mu|\alpha) \propto \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1}$$
- Comparison with the definition gives us the normalization factor
 
$$\begin{aligned} p(\mu|\mathcal{D}, \alpha) &= \text{Dir}(\mu|\alpha + m) \\ &= \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + m_1) \dots \Gamma(\alpha_K + m_K)} \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1} \end{aligned}$$

⇒ We can interpret the parameters  $\alpha_k$  of the Dirichlet prior as an effective number of observations of  $x_k = 1$ .

Advanced Machine Learning Winter'12 | Slide adapted from C. Bishop | B. Leibe | 37

RWTH AACHEN UNIVERSITY

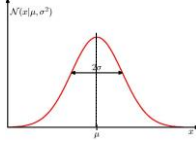
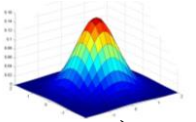
## Topics of This Lecture

- Probability Distributions
  - Bayesian Estimation Reloaded
- Binary Variables
  - Bernoulli distribution
  - Binomial distribution
  - Beta distribution
- Multinomial Variables
  - Multinomial distribution
  - Dirichlet distribution
- Continuous Variables
  - Gaussian distribution
  - Gamma distribution
  - Student's t distribution
  - Exponential Family

38

RWTH AACHEN UNIVERSITY

## The Gaussian Distribution

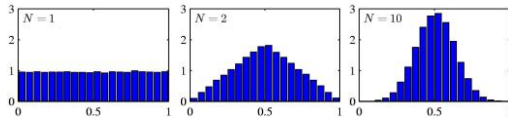
- One-dimensional case
  - Mean  $\mu$
  - Variance  $\sigma^2$
$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

- Multi-dimensional case
  - Mean  $\mu$
  - Covariance  $\Sigma$
$$\mathcal{N}(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)\right\}$$


39

RWTH AACHEN UNIVERSITY

## Gaussian Distribution - Properties

- Central Limit Theorem
  - "The distribution of the sum of  $N$  i.i.d. random variables becomes increasingly Gaussian as  $N$  grows."
  - In practice, the convergence to a Gaussian can be very rapid.
  - This makes the Gaussian interesting for many applications.
- Example:  $N$  uniform  $[0,1]$  random variables.



40

RWTH AACHEN UNIVERSITY

## Gaussian Distribution - Properties

- Properties
  - $\mathbb{E}[\mathbf{x}] = \mu$
  - $\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \mu\mu^T + \Sigma$
  - $\text{cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] = \Sigma$
- Limitations
  - Distribution is intrinsically unimodal, i.e. it is unable to provide a good approximation to multimodal distributions.
  - ⇒ We will see how to fix that with mixture distributions later...

41

RWTH AACHEN UNIVERSITY

## Bayes' Theorem for Gaussian Variables

- Marginal and Conditional Gaussians
  - Suppose we are given a Gaussian prior  $p(\mathbf{x})$  and a Gaussian conditional distribution  $p(\mathbf{y}|\mathbf{x})$  (a linear Gaussian model)
 
$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mu, \Lambda^{-1})$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$
  - From this, we can compute
 
$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mu + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\Lambda^{-1}\mathbf{A}^T)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\Sigma\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \Lambda\mu\}, \Sigma)$$

where

$$\Sigma = (\Lambda + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}$$

⇒ Closed-form solution for (Gaussian) marginal and posterior.

42

RWTH AACHEN UNIVERSITY

## Maximum Likelihood for the Gaussian

- Maximum Likelihood
  - Given i.i.d. data  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ , the log likelihood function is given by
 
$$\log p(\mathbf{X}|\mu, \Sigma) = -\frac{ND}{2} \log(2\pi) - \frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu)$$
- Sufficient statistics
  - The likelihood depends on the data set only through
 
$$\sum_{n=1}^N \mathbf{x}_n \quad \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$
  - Those are the **sufficient statistics** for the Gaussian distribution.

43

RWTH AACHEN UNIVERSITY

## ML for the Gaussian

- Setting the derivative to zero
 
$$\frac{\partial}{\partial \mu} \ln p(\mathbf{X}|\mu, \Sigma) = \sum_{n=1}^N \Sigma^{-1}(\mathbf{x}_n - \mu) = 0$$
- Solve to obtain
 
$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$
- And similarly, but a bit more involved
 
$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{\text{ML}})(\mathbf{x}_n - \mu_{\text{ML}})^T.$$

Slide credit: C. Bishop. B. Leibe. 44

RWTH AACHEN UNIVERSITY

## ML for the Gaussian

- Comparison with true results
  - Under the true distribution, we obtain
 
$$\begin{aligned} \mathbb{E}[\mu_{\text{ML}}] &= \mu \\ \mathbb{E}[\Sigma_{\text{ML}}] &= \frac{N-1}{N} \Sigma. \end{aligned}$$
  - ⇒ The ML estimate for the covariance is biased and underestimates the true covariance!
  - Therefore define the following unbiased estimator
 
$$\tilde{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \mu_{\text{ML}})(\mathbf{x}_n - \mu_{\text{ML}})^T.$$

Slide adapted from C. Bishop. B. Leibe. 45

RWTH AACHEN UNIVERSITY

## Bayesian Inference for the Gaussian

- Let's begin with a simple example
  - Consider a single Gaussian random variable  $x$ .
  - Assume  $\sigma^2$  is known and the task is to infer the mean  $\mu$ .
  - Given i.i.d. data  $\mathbf{X} = (x_1, \dots, x_N)^T$ , the likelihood function for  $\mu$  is given by
 
$$p(\mathbf{X}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2\right\}.$$
  - The likelihood function has a Gaussian shape as a function of  $\mu$ .
  - ⇒ The conjugate prior for this case is again a Gaussian.
 
$$p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2).$$

Slide adapted from C. Bishop. B. Leibe. 46

RWTH AACHEN UNIVERSITY

## Bayesian Inference for the Gaussian

- Combined with a Gaussian prior over  $\mu$ 

$$p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2).$$
- This results in the posterior
 
$$p(\mu|\mathbf{x}) \propto p(\mathbf{x}|\mu)p(\mu).$$
- Completing the square over  $\mu$ , we can derive that
 
$$p(\mu|\mathbf{x}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$$
- where
 
$$\begin{aligned} \mu_N &= \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{\text{ML}}, & \mu_{\text{ML}} &= \frac{1}{N} \sum_{n=1}^N x_n \\ \frac{1}{\sigma_N^2} &= \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}. \end{aligned}$$

Slide adapted from C. Bishop. B. Leibe. 47

RWTH AACHEN UNIVERSITY

## Visualization of the Results

- Bayes estimate:
 
$$\mu_N = \frac{\sigma^2 \mu_0 + N\sigma_0^2 \mu_{\text{ML}}}{\sigma^2 + N\sigma_0^2}$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$
- Behavior for large  $N$ 

	$N = 0$	$N \rightarrow \infty$
$\mu_N$	$\mu_0$	$\mu_{\text{ML}}$
$\sigma_N^2$	$\sigma_0^2$	0

Slide adapted from Bernt Schiele. B. Leibe. Image source: C. M. Bishop, 2006. 48

RWTH AACHEN UNIVERSITY

## Bayesian Inference for the Gaussian

- More complex case
  - Now assume  $\mu$  is known and the precision  $\lambda$  shall be inferred.
  - The likelihood function for  $\lambda = 1/\sigma^2$  is given by
 
$$p(\mathbf{X}|\lambda) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \lambda^{-1}) \propto \lambda^{N/2} \exp\left\{-\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}.$$
  - This has the shape of a Gamma function of  $\lambda$ .

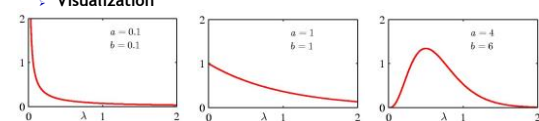
Slide adapted from C. Bishop. B. Leibe. 49



RWTH AACHEN UNIVERSITY

## The Gamma Distribution

- Gamma distribution
  - Product of a power of  $\lambda$  and the exponential of a linear function of  $\lambda$ .
 
$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$
- Properties
  - Finite integral if  $a > 0$  and the distribution itself is finite if  $a \geq 1$ .
  - Moments  $\mathbb{E}[\lambda] = \frac{a}{b}$   $\text{var}[\lambda] = \frac{a}{b^2}$
  - Visualization



Slide adapted from C. Bishop B. Leibe Image source: C.M. Bishop, 2004

RWTH AACHEN UNIVERSITY

## Bayesian Inference for the Gaussian

- Bayesian estimation
  - Combine a Gamma prior  $\text{Gam}(\lambda|a_0, b_0)$  with the likelihood function to obtain
 
$$p(\lambda|\mathbf{X}) \propto \lambda^{a_0-1} \lambda^{N/2} \exp\left\{-b_0\lambda - \frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$
  - We recognize this again as a Gamma function  $\text{Gam}(\lambda|a_N, b_N)$  with
 
$$a_N = a_0 + \frac{N}{2}$$

$$b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{\text{ML}}^2$$

Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## Bayesian Inference for the Gaussian

- Even more complex case
  - Assume that both  $\mu$  and  $\lambda$  are unknown
  - The joint likelihood function is given by
 
$$p(\mathbf{X}|\mu, \lambda) = \prod_{n=1}^N \left(\frac{\lambda}{2\pi}\right)^{1/2} \exp\left\{-\frac{\lambda}{2}(x_n - \mu)^2\right\}$$

$$\propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^N \exp\left\{\lambda\mu \sum_{n=1}^N x_n - \frac{\lambda}{2} \sum_{n=1}^N x_n^2\right\}$$

⇒ Need a prior with the same functional dependence on  $\mu$  and  $\lambda$ .

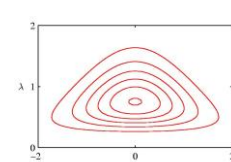
Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## The Gaussian-Gamma Distribution

- Gaussian-Gamma distribution
 
$$p(\mu, \lambda) = \mathcal{N}(\mu|\mu_0, (\beta\lambda)^{-1}) \text{Gam}(\lambda|a, b)$$

$$\propto \underbrace{\exp\left\{-\frac{\beta\lambda}{2}(\mu - \mu_0)^2\right\}}_{\text{Quadratic in } \mu} \underbrace{\lambda^{a-1} \exp\{-b\lambda\}}_{\text{Linear in } \lambda}$$
- Visualization



Slide adapted from C. Bishop B. Leibe Image source: C.M. Bishop, 2004

RWTH AACHEN UNIVERSITY

## Bayesian Inference for the Gaussian

- Multivariate conjugate priors
  - $\mu$  unknown,  $\Lambda$  known:  $p(\mu)$  Gaussian.
  - $\Lambda$  unknown,  $\mu$  known:  $p(\Lambda)$  Wishart,
 
$$\mathcal{W}(\Lambda|\mathbf{W}, \nu) = B|\Lambda|^{(\nu-D-1)/2} \exp\left(-\frac{1}{2}\text{Tr}(\mathbf{W}^{-1}\Lambda)\right)$$
  - $\Lambda$  and  $\mu$  unknown:  $p(\mu, \Lambda)$  Gaussian-Wishart,
 
$$p(\mu, \Lambda|\mu_0, \beta, \mathbf{W}, \nu) = \mathcal{N}(\mu|\mu_0, (\beta\Lambda)^{-1}) \mathcal{W}(\Lambda|\mathbf{W}, \nu)$$

Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## Student's t-Distribution

- Gaussian estimation
  - The conjugate prior for the precision of a Gaussian is a Gamma distribution.
  - Suppose we have a univariate Gaussian  $\mathcal{N}(x|\mu, \tau^{-1})$  together with a Gamma prior  $\text{Gam}(\tau|a, b)$ .
  - By integrating out the precision, obtain the marginal distribution
 
$$p(x|\mu, a, b) = \int_0^\infty \mathcal{N}(x|\mu, \tau^{-1}) \text{Gam}(\tau|a, b) d\tau$$

$$= \int_0^\infty \mathcal{N}(x|\mu, (\eta\lambda)^{-1}) \text{Gam}(\eta|\nu/2, \nu/2) d\eta$$
  - This corresponds to an infinite mixture of Gaussians having the same mean, but different precision.

Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## Student's t-Distribution

- Student's t-Distribution
  - We reparametrize the infinite mixture of Gaussians to get
 
$$St(x|\mu, \lambda, \nu) = \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left(\frac{\lambda}{\pi\nu}\right)^{1/2} \left[1 + \frac{\lambda(x - \mu)^2}{\nu}\right]^{-\nu/2 - 1/2}$$
- Parameters
  - "Precision"  $\lambda = a/b$
  - "Degrees of freedom"  $\nu = 2a$ .

Advanced Machine Learning Winter'12

56

Slide adapted from C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## Student's t-Distribution: Visualization

- Behavior
 

	$\nu = 1$	$\nu \rightarrow \infty$
$St(x \mu, \lambda, \nu)$	Cauchy	$\mathcal{N}(x \mu, \lambda^{-1})$

Advanced Machine Learning Winter'12

57

Slide adapted from C. Bishop B. Leibe Image source: C.M. Bishop, 2006

RWTH AACHEN UNIVERSITY

## Student's t-Distribution

- Robustness to outliers: **Gaussian vs t-distribution.**

⇒ The t-distribution is much less sensitive to outliers, can be used for robust regression.

⇒ Downside: ML solution for t-distribution requires EM algorithm.

Advanced Machine Learning Winter'12

58

Slide adapted from C. Bishop B. Leibe Image source: C. M. Bishop, 2006

RWTH AACHEN UNIVERSITY

## Student's t-Distribution: Multivariate Case

- Multivariate case in  $D$  dimensions
 
$$St(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) = \int_0^\infty \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, (\eta\boldsymbol{\Lambda})^{-1}) \text{Gam}(\eta|\nu/2, \nu/2) d\eta$$

$$= \frac{\Gamma(D/2 + \nu/2)}{\Gamma(\nu/2)} \frac{|\boldsymbol{\Lambda}|^{1/2}}{(\pi\nu)^{D/2}} \left[1 + \frac{\Delta^2}{\nu}\right]^{-D/2 - \nu/2}$$

where  $\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})$  is the Mahalanobis distance.
- Properties
  - $\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$ , if  $\nu > 1$
  - $\text{cov}[\mathbf{x}] = \frac{\nu}{(\nu - 2)} \boldsymbol{\Lambda}^{-1}$ , if  $\nu > 2$
  - $\text{mode}[\mathbf{x}] = \boldsymbol{\mu}$

Advanced Machine Learning Winter'12

59

Slide credit: C. Bishop B. Leibe

RWTH AACHEN UNIVERSITY

## References and Further Reading

- Probability distributions and their properties are described in Chapter 2 of Bishop's book.

Christopher M. Bishop  
Pattern Recognition and Machine Learning  
Springer, 2006

Advanced Machine Learning Winter'15

74

B. Leibe