

RWTH AACHEN  
UNIVERSITY

# Advanced Machine Learning Lecture 6

## Approximate Inference

10.11.2016

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
  - Approximate Inference
- Deep Learning
  - Neural Networks
  - CNNs, RNNs, ResNets, 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:
 
$$\bar{\mathbf{f}}_* = K(X_*, X) K(X, X)^{-1} \mathbf{f}$$

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

Slide adapted from Bernt Schiele B. Leibe 3

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  $X$ .

Slide adapted from Bernt Schiele B. Leibe 4

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.

Slide credit: Bernt Schiele B. Leibe 5

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**.
- $$p(\mathbf{t} | X, \theta, \mathcal{H}_i) = \int p(\mathbf{t} | X, \mathbf{w}, \theta, \mathcal{H}_i) p(\mathbf{w} | \theta, \mathcal{H}_i) d\mathbf{w}$$

Slide credit: Bernt Schiele B. Leibe 6

RWTH AACHEN UNIVERSITY

## Recap: Model Selection at Mid Level

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

$$= \frac{p(\mathbf{t}|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i)}{p(\mathbf{t}|X, \mathcal{H}_i)}$$
- where
  - The marginal likelihood of the previous level  $p(\mathbf{t}|X, \theta, \mathcal{H}_i)$  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:
 
$$p(\mathbf{t}|X, \mathcal{H}_i) = \int p(\mathbf{t}|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i)d\theta$$

Slide credit: Bernt Schiele B. Leibe 7

RWTH AACHEN UNIVERSITY

## Recap: Model Selection at Top Level

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

Slide credit: Bernt Schiele B. Leibe 8

RWTH AACHEN UNIVERSITY

## Recap: Bayesian Model Selection

- Discussion
  - Marginal likelihood is main difference to non-Bayesian methods
 
$$p(\mathbf{t}|X, \mathcal{H}_i) = \int p(\mathbf{t}|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i)d\theta$$
  - 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 **marginal likelihood**.
    - A complex model can account for a large range of possible sets of target values - therefore, it can never attain a very high **marginal likelihood**.

Slide credit: Bernt Schiele B. Leibe image source: Rasmussen & Williams, 2006 9

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Approximate Inference
  - Variational methods
  - Sampling approaches
- Sampling approaches
  - Sampling from a distribution
  - Ancestral Sampling
  - Rejection Sampling
  - Importance Sampling
- Markov Chain Monte Carlo
  - Markov Chains
  - Metropolis Algorithm
  - Metropolis-Hastings Algorithm
  - Gibbs Sampling

Slide credit: Bernt Schiele B. Leibe 10

RWTH AACHEN UNIVERSITY

## Approximate Inference

- Exact Bayesian inference is often intractable.
  - Often infeasible to evaluate the posterior distribution or to compute expectations w.r.t. the distribution.
    - E.g. because the dimensionality of the latent space is too high.
    - Or because the posterior distribution has a too complex form.
  - Problems with continuous variables
    - Required integrations may not have closed-form solutions.
  - Problems with discrete variables
    - Marginalization involves summing over all possible configurations of the hidden variables.
    - There may be exponentially many such states.

⇒ We need to resort to approximation schemes.

Slide credit: Bernt Schiele B. Leibe 11

RWTH AACHEN UNIVERSITY

## Two Classes of Approximation Schemes

- Deterministic approximations (Variational methods)
  - Based on analytical approximations to the posterior distribution
    - E.g. by assuming that it factorizes in a certain form
    - Or that it has a certain parametric form (e.g. a Gaussian).
  - ⇒ Can never generate exact results, but are often scalable to large applications.
- Stochastic approximations (Sampling methods)
  - Given infinite computationally resources, they can generate exact results.
  - Approximation arises from the use of a finite amount of processor time.
  - ⇒ Enable the use of Bayesian techniques across many domains.
  - ⇒ But: computationally demanding, often limited to small-scale problems.

Slide credit: Bernt Schiele B. Leibe 12

Advanced Machine Learning Winter'12

## Topics of This Lecture

- Approximate Inference
  - Variational methods
  - Sampling approaches
- Sampling approaches
  - Sampling from a distribution
  - Ancestral Sampling
  - Rejection Sampling
  - Importance Sampling
- Markov Chain Monte Carlo
  - Markov Chains
  - Metropolis Algorithm
  - Metropolis-Hastings Algorithm
  - Gibbs Sampling

B. Leibe 13

Advanced Machine Learning Winter'12

## Sampling Idea

- Objective:
  - Evaluate expectation of a function  $f(\mathbf{z})$  w.r.t. a probability distribution  $p(\mathbf{z})$ .
$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z}$$
- Sampling idea
  - Draw  $L$  independent samples  $\mathbf{z}^{(l)}$  with  $l = 1, \dots, L$  from  $p(\mathbf{z})$ .
  - This allows the expectation to be approximated by a finite sum
$$\hat{f} = \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}^{(l)})$$
  - As long as the samples  $\mathbf{z}^{(l)}$  are drawn independently from  $p(\mathbf{z})$ , then
 
$$\|\mathbb{E}[\hat{f}] - \mathbb{E}[f]\|$$

⇒ Unbiased estimate, independent of the dimension of  $\mathbf{z}$ !

Slide adapted from Bernd Schiele B. Leibe Image source: C.M. Bishop, 2004 14

Advanced Machine Learning Winter'12

## Sampling - Challenges

- Problem 1: Samples might not be independent
  - ⇒ Effective sample size might be much smaller than apparent sample size.

- Problem 2:
  - If  $f(\mathbf{z})$  is small in regions where  $p(\mathbf{z})$  is large and vice versa, the expectation may be dominated by regions of small probability.
  - ⇒ Large sample sizes necessary to achieve sufficient accuracy.

B. Leibe Image source: C.M. Bishop, 2004 15

Advanced Machine Learning Winter'12

## Parametric Density Model

- Example:
  - A simple multivariate (d-dimensional) Gaussian model
$$p(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\}$$
  - This is a “generative” model in the sense that we can generate samples  $\mathbf{x}$  according to the distribution.

Slide adapted from Bernd Schiele B. Leibe 16

Advanced Machine Learning Winter'12

## Sampling from a Gaussian

- Given: 1-dim. Gaussian pdf (probability density function)  $p(x|\mu, \sigma^2)$  and the corresponding cumulative distribution:
 
$$F_{\mu, \sigma^2}(x) = \int_{-\infty}^x p(x|\mu, \sigma^2)dx$$
- To draw samples from a Gaussian, we can invert the cumulative distribution function:
 
$$u \sim \text{Uniform}(0, 1) \Rightarrow F_{\mu, \sigma^2}^{-1}(u) \sim p(x|\mu, \sigma^2)$$

Slide credit: Bernd Schiele B. Leibe 17

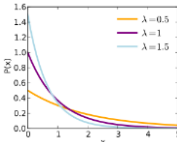
Advanced Machine Learning Winter'12

## Sampling from a pdf (Transformation method)

- In general, assume we are given the pdf  $p(\mathbf{x})$  and the corresponding cumulative distribution:
 
$$F(x) = \int_{-\infty}^x p(z)dz$$
- To draw samples from this pdf, we can invert the cumulative distribution function:
 
$$u \sim \text{Uniform}(0, 1) \Rightarrow F^{-1}(u) \sim p(x)$$

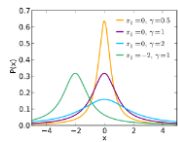
Slide credit: Bernd Schiele B. Leibe Image source: C.M. Bishop, 2004 18

**Example 1: Sampling from Exponential Distrib.**

- Exponential Distribution
 
$$p(y) = \lambda \exp(-\lambda y)$$
 where  $0 \leq y < \infty$ .
 
- Transformation sampling
  - Indefinite Integral  $h(y) = 1 - \exp(-\lambda y)$
  - Inverse function
 
$$y = h(y)^{-1} = -\lambda^{-1} \ln(1 - z)$$
 for a uniformly distributed input variable  $z$ .

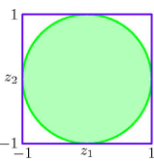
19  
B. Leibe  
Image source: Wikipedia

**Example 2: Sampling from Cauchy Distrib.**

- Cauchy Distribution
 
$$p(y) = \frac{1}{\pi} \frac{1}{1 + y^2}$$

- Transformation sampling
  - Inverse of integral can be expressed as a tan function.
 
$$y = h(y)^{-1} = \tan(z)$$
 for a uniformly distributed input variable  $z$ .

20  
B. Leibe  
Image source: Wikipedia

**Note: Efficient Sampling from a Gaussian**

- Problem with transformation method
  - Integral over Gaussian cannot be expressed in analytical form.
  - Standard transformation approach is very inefficient.
- More efficient: Box-Muller Algorithm
  - Generate pairs of uniformly distributed random numbers  $z_1, z_2 \in (-1, 1)$ .
  - Discard each pair unless it satisfies  $r^2 = z_1^2 + z_2^2 \leq 1$ .
  - This leads to a uniform distribution of points inside the unit circle with  $p(z_1, z_2) = 1/\pi$ .

21  
B. Leibe  
Image source: C. M. Bishop, 2006

**Box-Muller Algorithm (cont'd)**

- Box-Muller Algorithm (cont'd)
  - For each pair  $z_1, z_2$  evaluate
 
$$y_1 = z_1 \left( \frac{-2 \ln r^2}{r^2} \right)^{1/2} \quad y_2 = z_2 \left( \frac{-2 \ln r^2}{r^2} \right)^{1/2}$$
  - Then the joint distribution of  $y_1$  and  $y_2$  is given by
 
$$p(y_1, y_2) = p(z_1, z_2) \left| \frac{\partial(z_1, z_2)}{\partial(y_1, y_2)} \right|$$

$$= \left[ \frac{1}{\sqrt{2\pi}} \exp(-y_1^2/2) \right] \left[ \frac{1}{\sqrt{2\pi}} \exp(-y_2^2/2) \right]$$
  - $\Rightarrow y_1$  and  $y_2$  are independent and each has a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ .
  - If  $y \sim \mathcal{N}(0, 1)$ , then  $\sigma y + \mu \sim \mathcal{N}(\mu, \sigma^2)$ .


22  
B. Leibe

**Box-Muller Algorithm (cont'd)**

- Multivariate extension
  - If  $z$  is a vector valued random variable whose components are independent and Gaussian distributed with  $\mathcal{N}(0, 1)$ ,
  - Then  $y = \mu + Lz$  will have mean  $\mu$  and covariance  $\Sigma$ .
  - Where  $\Sigma = LL^T$  is the Cholesky decomposition of  $\Sigma$ .

23  
B. Leibe

**Ancestral Sampling**

- Generalization of this idea to directed graphical models.
  - Joint probability factorizes into conditional probabilities:
 
$$p(x) = \prod_{k=1}^K p(x_k | pa_k)$$

  - Ancestral sampling
    - Assume the variables are ordered such that there are no links from any node to a lower-numbered node.
    - Start with lowest-numbered node and draw a sample from its distribution.
 
$$\hat{x}_1 \sim p(x_1)$$
    - Cycle through each of the nodes in order and draw samples from the conditional distribution (where the parent variable is set to its sampled value).
 
$$\hat{x}_n \sim p(x_n | pa_n)$$

24  
B. Leibe  
Image source: C. M. Bishop, 2006

Advanced Machine Learning Winter'12

## Logic Sampling

RWTH AACHEN UNIVERSITY

- Extension of Ancestral sampling
  - Directed graph where some nodes are instantiated with observed values.
- Use ancestral sampling, except
  - When sample is obtained for an observed variable, if they agree then sample value is retained and proceed to next variable.
  - If they don't agree, whole sample is discarded.
- Result
  - Approach samples correctly from the posterior distribution.
  - However, probability of accepting a sample decreases rapidly as the number of observed variables increases.

⇒ Approach is rarely used in practice.

B. Leibe 25

Advanced Machine Learning Winter'12

## Discussion

RWTH AACHEN UNIVERSITY

- Transformation method
  - Limited applicability, as we need to invert the indefinite integral of the required distribution  $p(z)$ .
  - This will only be feasible for a limited number of simple distributions.
- More general
  - Rejection Sampling
  - Importance Sampling

Slide adapted from Bernd Schiele B. Leibe 26

Advanced Machine Learning Winter'12

## Rejection Sampling

RWTH AACHEN UNIVERSITY

- Assumptions
  - Sampling directly from  $p(z)$  is difficult.
  - But we can easily evaluate  $p(z)$  (up to some normalization factor  $Z_p$ ):
 
$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$
- Idea
  - We need some simpler distribution  $q(z)$  (called **proposal distribution**) from which we can draw samples.
  - Choose a constant  $k$  such that:  $\forall z : kq(z) \geq \tilde{p}(z)$

Slide credit: Bernd Schiele B. Leibe Image source: C.M. Bishop, 2006 27

Advanced Machine Learning Winter'12

## Rejection Sampling

RWTH AACHEN UNIVERSITY

- Sampling procedure
  - Generate a number  $z_0$  from  $q(z)$ .
  - Generate a number  $u_0$  from the uniform distribution over  $[0, kq(z_0)]$ .
  - If  $u_0 > \tilde{p}(z_0)$  reject sample, otherwise accept.
    - Sample is rejected if it lies in the grey shaded area.
    - The remaining pairs  $(u_0, z_0)$  have uniform distribution under the curve  $\tilde{p}(z)$ .
- Discussion
  - Original values of  $z$  are generated from the distribution  $q(z)$ .
  - Samples are accepted with probability  $\tilde{p}(z)/kq(z)$ 

$$p(\text{accept}) = \int \frac{\tilde{p}(z)}{kq(z)} q(z) dz = \frac{1}{k} \int \tilde{p}(z) dz$$

⇒  $k$  should be as small as possible!

Slide credit: Bernd Schiele B. Leibe Image source: C.M. Bishop, 2006 28

Advanced Machine Learning Winter'12

## Rejection Sampling - Discussion

RWTH AACHEN UNIVERSITY

- Limitation: high-dimensional spaces
  - For rejection sampling to be of practical value, we require that  $kq(z)$  be close to the required distribution, so that the rate of rejection is minimal.
- Artificial example
  - Assume that  $p(z)$  is Gaussian with covariance matrix  $\sigma_p^2 I$
  - Assume that  $q(z)$  is Gaussian with covariance matrix  $\sigma_q^2 I$
  - Obviously:  $\sigma_q^2 \geq \sigma_p^2$
  - In  $D$  dimensions:  $k = (\sigma_q/\sigma_p)^D$ .
    - Assume  $\sigma_q$  is just 1% larger than  $\sigma_p$ ,
    - $D = 1000 \Rightarrow k = 1.01^{1000} \geq 20,000$
    - And  $p(\text{accept}) \cdot \frac{1}{20000}$

⇒ Often impractical to find good proposal distributions for high dimensions!

Slide credit: Bernd Schiele B. Leibe Image source: C.M. Bishop, 2006 29

Advanced Machine Learning Winter'12

## Example: Sampling from a Gamma Distrib.

RWTH AACHEN UNIVERSITY

- Gamma distribution
 
$$\text{Gam}(z|a, b) = \frac{1}{\Gamma(a)} b^a z^{a-1} \exp(-bz) \quad a > 1$$
- Rejection sampling approach
  - For  $a > 1$ , Gamma distribution has a bell-shaped form.
  - Suitable proposal distribution is Cauchy (for which we can use the transformation method).
  - Generalize Cauchy slightly to ensure it is nowhere smaller than Gamma:  $y = b \tan y + c$  for uniform  $y$ .
  - This gives random numbers distributed according to
 
$$q(z) = \frac{k}{1 + (z - c)^2/b^2} \quad \text{with optimal rejection rate for } \begin{matrix} c = a - 1 \\ b^2 = 2a - 1 \end{matrix}$$

B. Leibe Image source: C.M. Bishop, 2006 30

RWTH AACHEN UNIVERSITY

## Importance Sampling

- Approach
  - Approximate expectations directly (but does not enable to draw samples from  $p(\mathbf{z})$  directly).
  - Goal: 
$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z}$$
- Simplistic strategy: Grid sampling
  - Discretize  $\mathbf{z}$ -space into a uniform grid.
  - Evaluate the integrand as a sum of the form
 
$$\mathbb{E}[f] \simeq \sum_{l=1}^L f(\mathbf{z}^{(l)})p(\mathbf{z}^{(l)})d\mathbf{z}$$
  - But: number of terms grows exponentially with number of dimensions!

Advanced Machine Learning Winter'12 31

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

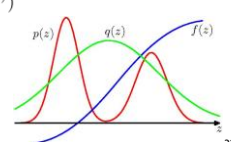
## Importance Sampling

- Idea
  - Use a proposal distribution  $q(\mathbf{z})$  from which it is easy to draw samples.
  - Express expectations in the form of a finite sum over samples  $\{\mathbf{z}^{(l)}\}$  drawn from  $q(\mathbf{z})$ .

$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int f(\mathbf{z})\frac{p(\mathbf{z})}{q(\mathbf{z})}q(\mathbf{z})d\mathbf{z}$$

$$\simeq \frac{1}{L} \sum_{l=1}^L \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})} f(\mathbf{z}^{(l)})$$

- with importance weights
 
$$r_l = \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})}$$



Advanced Machine Learning Winter'12 32

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

## Importance Sampling

- Typical setting:
  - $p(\mathbf{z})$  can only be evaluated up to an unknown normalization constant
 
$$p(\mathbf{z}) = \tilde{p}(\mathbf{z})/Z_p$$
  - $q(\mathbf{z})$  can also be treated in a similar fashion.
 
$$q(\mathbf{z}) = \tilde{q}(\mathbf{z})/Z_q$$
- Then
 
$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} = \frac{Z_q}{Z_p} \int f(\mathbf{z})\frac{\tilde{p}(\mathbf{z})}{\tilde{q}(\mathbf{z})}q(\mathbf{z})d\mathbf{z}$$

$$\simeq \frac{Z_q}{Z_p} \frac{1}{L} \sum_{l=1}^L \tilde{r}_l f(\mathbf{z}^{(l)})$$
- with:  $\tilde{r}_l = \frac{\tilde{p}(\mathbf{z}^{(l)})}{\tilde{q}(\mathbf{z}^{(l)})}$

Advanced Machine Learning Winter'12 33

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

## Importance Sampling

- Ratio of normalization constants can be evaluated
 
$$\frac{Z_p}{Z_q} = \frac{1}{Z_q} \int \tilde{p}(\mathbf{z})d\mathbf{z} = \int \frac{\tilde{p}(\mathbf{z}^{(l)})}{\tilde{q}(\mathbf{z}^{(l)})}q(\mathbf{z})d\mathbf{z} \simeq \frac{1}{L} \sum_{l=1}^L \tilde{r}_l$$
- and therefore
 
$$\mathbb{E}[f] \simeq \sum_{l=1}^L w_l f(\mathbf{z}^{(l)})$$
- with
 
$$w_l = \frac{\tilde{r}_l}{\sum_m \tilde{r}_m} = \frac{\frac{\tilde{p}(\mathbf{z}^{(l)})}{\tilde{q}(\mathbf{z}^{(l)})}}{\sum_m \frac{\tilde{p}(\mathbf{z}^{(m)})}{\tilde{q}(\mathbf{z}^{(m)})}}$$

Advanced Machine Learning Winter'12 34

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

## Importance Sampling - Discussion

- Observations
  - Success of importance sampling depends crucially on how well the sampling distribution  $q(\mathbf{z})$  matches the desired distribution  $p(\mathbf{z})$ .
  - Often,  $p(\mathbf{z})f(\mathbf{z})$  is strongly varying and has a significant proportion of its mass concentrated over small regions of  $\mathbf{z}$ -space.
  - ⇒ Weights  $r_l$  may be dominated by a few weights having large values.
  - Practical issue: if none of the samples falls in the regions where  $p(\mathbf{z})f(\mathbf{z})$  is large...
    - The results may be arbitrary in error.
    - And there will be no diagnostic indication (no large variance in  $r_l$ )!
  - Key requirement for sampling distribution  $q(\mathbf{z})$ :
    - Should not be small or zero in regions where  $p(\mathbf{z})$  is significant!

Advanced Machine Learning Winter'12 35

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Approximate Inference
  - Variational methods
  - Sampling approaches
- Sampling approaches
  - Sampling from a distribution
  - Ancestral Sampling
  - Rejection Sampling
  - Importance Sampling
- Markov Chain Monte Carlo
  - Markov Chains
  - Metropolis Algorithm
  - Metropolis-Hastings Algorithm
  - Gibbs Sampling

Advanced Machine Learning Winter'12 36

Slide credit: Bernt Schiele B. Leibe

## References and Further Reading

- Sampling methods for approximate inference are described in detail in Chapter 11 of Bishop's book.



David MacKay  
Information Theory, Inference, and Learning Algorithms  
Cambridge University Press, 2003



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

- Another good introduction to Monte Carlo methods can be found in Chapter 29 of MacKay's book (also available online: <http://www.inference.phy.cam.ac.uk/mackay/itprnn/book.html>)