

# Advanced Machine Learning Summer 2019

## Part 20 – Repetition 11.07.2019

Prof. Dr. Bastian Leibe

RWTH Aachen University, Computer Vision Group  
<http://www.vision.rwth-aachen.de>



## Announcements

- Today, I'll summarize the most important points from the lecture.
  - It is an opportunity for you to ask questions...
  - ...or get additional explanations about certain topics.
  - So, please do ask.
- Today's slides are intended as an index for the lecture.
  - Summarizing the most important points from each class
  - But they are not complete, won't be sufficient as only tool.
  - Also look at the exercises – they often explain algorithms in detail.
- Exam procedure
  - Closed-book exam, the core exam time will be 2h.
  - We will send around an announcement with the exact starting times and places by email.

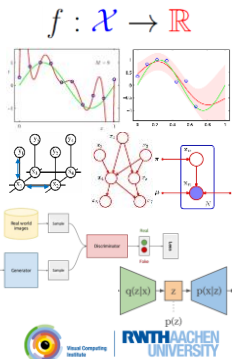
2

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders



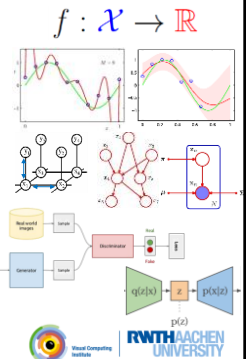
3

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders



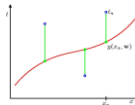
4

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Regression

- Learning to predict a continuous function value
  - Given: training set  $\mathbf{X} = \{x_1, \dots, x_N\}$  with target values  $\mathbf{T} = \{t_1, \dots, t_N\}$ .
  - Learn a continuous function  $y(x)$  to predict the function value for a new input  $x$ .
- Define an error function  $E(\mathbf{w})$  to optimize
  - E.g., sum-of-squares error
$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$
  - Procedure: Take the derivative and set it to zero
$$\frac{\partial E(\mathbf{w})}{\partial w_j} = \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\} \frac{\partial y(x_n, \mathbf{w})}{\partial w_j} \stackrel{!}{=} 0$$



5

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Least-Squares Regression

- Setup
  - Step 1: Define  $\tilde{\mathbf{x}}_i = \begin{pmatrix} \mathbf{x}_i \\ 1 \end{pmatrix}$ ,  $\tilde{\mathbf{w}} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$
  - Step 2: Rewrite  $\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} = t_i$ ,  $\forall i = 1, \dots, n$
  - Step 3: Matrix-vector notation  $\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} = \mathbf{t}$  with  $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n]$  and  $\mathbf{t} = [t_1, \dots, t_n]^T$
  - Step 4: Find least-squares solution  $\|\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{t}\|^2 \rightarrow \min$
  - Solution:  $\tilde{\mathbf{w}} = (\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T)^{-1} \tilde{\mathbf{X}} \mathbf{t}$

see Exercise 1.3

6

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernt Schiele



### Course Outline

- **Regression Techniques**
  - Linear Regression
  - **Regularization (Ridge, Lasso)**
  - Kernels (Kernel Ridge Regression)
- **Deep Reinforcement Learning**
- **Probabilistic Graphical Models**
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- **Deep Generative Models**
  - Generative Adversarial Networks
  - Variational Autoencoders

$f: \mathcal{X} \rightarrow \mathbb{R}$

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

### Recap: Regularization

- **Problem: Overfitting**
  - Many parameters & little data  $\Rightarrow$  tendency to overfit to the noise
  - Side effect: The coefficient values get very large.
- **Workaround: Regularization**
  - Penalize large coefficient values

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Here we've simply added a **quadratic regularizer**, which is simple to optimize

$$\|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + \dots + w_M^2$$

- The resulting form of the problem is called **Ridge Regression**.
- (Note:  $w_0$  is often omitted from the regularizer.)

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

### Recap: Probabilistic Regression

- **First assumption:**
  - Our target function values  $y$  are generated by adding noise to the function estimate:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon$$

Target function value  $\rightarrow$   $t$   $\leftarrow$  Noise

Regression function  $\rightarrow$   $y(\mathbf{x}, \mathbf{w})$   $\leftarrow$  Input value  $\leftarrow$  Weights or parameters

- **Second assumption:**
  - The noise is Gaussian distributed

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

Mean  $\rightarrow$   $y(\mathbf{x}, \mathbf{w})$   $\leftarrow$  Variance ( $\beta$  precision)

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernd Schiele

### Recap: Maximum Likelihood Regression

- **Given**
  - Training data points:  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$
  - Associated function values:  $\mathbf{t} = [t_1, \dots, t_n]^T$
- **Conditional likelihood (assuming i.i.d. data)**

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|y(\mathbf{x}_n, \mathbf{w}), \beta^{-1}) = \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

$\Rightarrow$  Maximize w.r.t.  $\mathbf{w}, \beta$

Generalized linear regression function

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernd Schiele

### Recap: Maximum Likelihood Regression

$$\nabla_{\mathbf{w}} \log p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = -\beta \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)$$

- **Setting the gradient to zero:**

$$0 = -\beta \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)$$

$$\Leftrightarrow \sum_{n=1}^N t_n \phi(\mathbf{x}_n) = \left[ \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right] \mathbf{w}$$

$$\Leftrightarrow \Phi \mathbf{t} = \Phi \Phi^T \mathbf{w} \quad \Phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]$$

$$\Leftrightarrow \mathbf{w}_{ML} = (\Phi \Phi^T)^{-1} \Phi \mathbf{t} \quad \leftarrow \text{Same as in least-squares regression!}$$

$\Rightarrow$  *Least-squares regression is equivalent to Maximum Likelihood under the assumption of Gaussian noise.*

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernd Schiele

### Recap: Role of the Precision Parameter

- Also use ML to determine the precision parameter  $\beta$ :

$$\log p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{N}{2} \log \beta - \frac{N}{2} \log(2\pi)$$

- **Gradient w.r.t.  $\beta$ :**

$$\nabla_{\beta} \log p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = -\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{N}{2} \frac{1}{\beta}$$

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$

$\Rightarrow$  *The inverse of the noise precision is given by the residual variance of the target values around the regression function.*

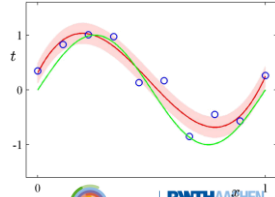
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

## Recap: Predictive Distribution

- Having determined the parameters  $\mathbf{w}$  and  $\beta$ , we can now make predictions for new values of  $\mathbf{x}$ .

$$p(t|\mathbf{X}, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1})$$

- This means
  - Rather than giving a point estimate, we can now also give an estimate of the estimation uncertainty.



## Recap: Maximum-A-Posteriori Estimation

- Introduce a prior distribution over the coefficients  $\mathbf{w}$ .
  - For simplicity, assume a zero-mean Gaussian distribution

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

- New hyperparameter  $\alpha$  controls the distribution of model parameters.
- Express the posterior distribution over  $\mathbf{w}$ .
  - Using Bayes' theorem:
 
$$p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \beta, \alpha) \propto p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)$$
  - We can now determine  $\mathbf{w}$  by maximizing the posterior.
  - This technique is called maximum-a-posteriori (MAP).

## Recap: MAP Solution

- Minimize the negative logarithm

$$-\log p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \beta, \alpha) \propto -\log p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) - \log p(\mathbf{w}|\alpha)$$

$$-\log p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 + \text{const}$$

$$-\log p(\mathbf{w}|\alpha) = \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const}$$

- The MAP solution is therefore the solution of

$$\frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

$\Rightarrow$  Maximizing the posterior distribution is equivalent to minimizing the regularized sum-of-squares error (with  $\lambda = \frac{\alpha}{\beta}$ ).

## Recap: MAP Solution (2)

$$\nabla_{\mathbf{w}} \log p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \beta, \alpha) = -\beta \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) + \alpha \mathbf{w}$$

- Setting the gradient to zero:

$$0 = -\beta \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) + \alpha \mathbf{w}$$

$$\Leftrightarrow \sum_{n=1}^N t_n \phi(\mathbf{x}_n) = \left[ \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right] \mathbf{w} + \frac{\alpha}{\beta} \mathbf{w}$$

$$\Leftrightarrow \Phi \mathbf{t} = \left( \Phi \Phi^T + \frac{\alpha}{\beta} \mathbf{I} \right) \mathbf{w} \quad \Phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)]$$

$$\Leftrightarrow \mathbf{w}_{\text{MAP}} = \left( \Phi \Phi^T + \frac{\alpha}{\beta} \mathbf{I} \right)^{-1} \Phi \mathbf{t}$$

Effect of regularization:  
Keeps the inverse well-conditioned

## Recap: Bayesian Curve Fitting

- Given
  - Training data points:  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{d \times N}$
  - Associated function values:  $\mathbf{t} = [t_1, \dots, t_N]^T$
  - Our goal is to predict the value of  $t$  for a new point  $\mathbf{x}$ .

- Evaluate the predictive distribution

$$p(t|x, \mathbf{X}, \mathbf{t}) = \int p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{X}, \mathbf{t}) d\mathbf{w}$$

What we just computed for MAP

- Noise distribution – again assume a Gaussian here

$$p(t|x, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

- Assume that parameters  $\alpha$  and  $\beta$  are fixed and known for now.

## Recap: Bayesian Curve Fitting

- Under those assumptions, the posterior distribution is a Gaussian and can be evaluated analytically:

$$p(t|x, \mathbf{X}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$

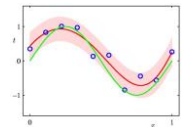
- where the mean and variance are given by

$$m(x) = \beta \phi(x)^T \mathbf{S} \sum_{n=1}^N \phi(\mathbf{x}_n) t_n$$

$$s(x)^2 = \beta^{-1} + \phi(x)^T \mathbf{S} \phi(x)$$

- and  $\mathbf{S}$  is the regularized covariance matrix

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T$$



### Recap: Loss Functions for Regression

Mean prediction

- Optimal prediction
  - Minimize the expected loss
 
$$\mathbb{E}[L] = \iint L(t, y(\mathbf{x}))p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$
  - Under squared loss, the optimal regression function is the mean  $\mathbb{E}[t|\mathbf{x}]$  of the posterior  $p(t|\mathbf{x})$  ("mean prediction").
  - For generalized linear regression function and squared loss:
 
$$y(\mathbf{x}) = \int t \mathcal{N}(t|\mathbf{w}^T \phi(\mathbf{x}), \beta^{-1}) \, dt = \mathbf{w}^T \phi(\mathbf{x})$$

19 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition  
Slide adapted from Stefan Roth. Image source: C.M. Bishop, 2006

### Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders

20 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition  
Image source: C.M. Bishop, 2006

### Recap: Loss Functions for Regression

- The squared loss is not the only possible choice
  - Poor choice when conditional distribution  $p(t|\mathbf{x})$  is multimodal.
- Simple generalization: Minkowski loss
 
$$L(t, y(\mathbf{x})) = |y(\mathbf{x}) - t|^q$$
  - Expectation
 
$$\mathbb{E}[L_q] = \iint |y(\mathbf{x}) - t|^q p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$
  - Minimum of  $\mathbb{E}[L_q]$  is given by
    - Conditional mean for  $q = 2$ ,
    - Conditional median for  $q = 1$ ,
    - Conditional mode for  $q = 0$ .

see Exercise 1.2

21 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition  
Image source: C.M. Bishop, 2006

### Recap: Linear Basis Function Models

- Generally, we consider models of the following form
 
$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$
  - where  $\phi_j(\mathbf{x})$  are known as *basis functions*.
  - In the simplest case, we use linear basis functions:  $\phi_d(\mathbf{x}) = x_d$ .
- Other popular basis functions
  - Polynomial
  - Gaussian
  - Sigmoid

22 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition  
Image source: C.M. Bishop, 2006

### Recap: Regularized Least-Squares

- Consider more general regularization functions
  - "L<sub>q</sub> norms":
 
$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$
- Effect: Sparsity for  $q \leq 1$ .
  - Minimization tends to set many coefficients to zero

23 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition  
Image source: C.M. Bishop, 2006

### Recap: Lasso as Bayes Estimation

- L<sub>1</sub> regularization ("The Lasso")
 
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \lambda \sum_{j=1}^M |w_j|$$
- Interpretation as Bayes Estimation
  - We can think of  $|w_j|^q$  as the log-prior density for  $w_j$ .
- Prior for Lasso ( $q = 1$ ): Laplacian distribution
 
$$p(\mathbf{w}) = \frac{1}{2\tau} \exp\{-|\mathbf{w}|/\tau\} \quad \text{with} \quad \tau = \frac{1}{\lambda}$$

24 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition  
Image source: Wikipedia

## Recap: The Lasso

- $L_1$  regularization ("The Lasso")

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \lambda \sum_{j=1}^M |w_j|$$

- The solution will be sparse (only few coefficients non-zero)
- The  $L_1$  penalty makes the problem non-linear.
- ⇒ There is no closed-form solution.

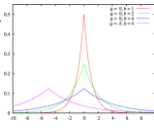
- Interpretation as Bayes Estimation

- We can think of  $|w_j|^{q/2}$  as the log-prior density for  $w_j$ .

- Prior for Lasso ( $q = 1$ ):

- Laplacian distribution

$$p(\mathbf{w}) = \frac{1}{2\tau} \exp\{-|\mathbf{w}|/\tau\} \quad \text{with} \quad \tau = \frac{1}{\lambda}$$



25

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Image source: Wikipedia

## Course Outline

- Regression Techniques

- Linear Regression
- Regularization (Ridge, Lasso)
- Kernels (Kernel Ridge Regression)

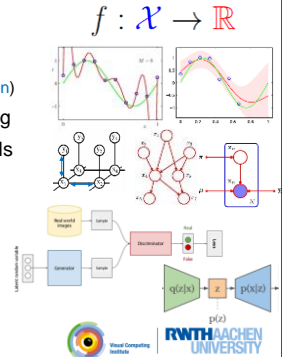
- Deep Reinforcement Learning

- Probabilistic Graphical Models

- Bayesian Networks
- Markov Random Fields
- Inference (exact & approximate)
- Latent Variable Models

- Deep Generative Models

- Generative Adversarial Networks
- Variational Autoencoders



26

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Kernel Ridge Regression

- Dual definition

- Instead of working with  $\mathbf{w}$ , substitute  $\mathbf{w} = \Phi^T \mathbf{a}$  into  $J(\mathbf{w})$  and write the result using the kernel matrix  $\mathbf{K} = \Phi \Phi^T$ :

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

- Solving for  $\mathbf{a}$ , we obtain

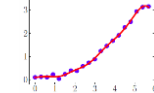
$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

- Prediction for a new input  $\mathbf{x}$ :

- Writing  $\mathbf{k}(\mathbf{x})$  for the vector with elements  $k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x})$

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

- ⇒ The dual formulation allows the solution to be entirely expressed in terms of the kernel function  $k(\mathbf{x}, \mathbf{x}')$ .



27

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Image source: Christoph Lampert

## Recap: Properties of Kernels

- Theorem

- Let  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a positive definite kernel function. Then there exists a Hilbert Space  $\mathcal{H}$  and a mapping  $\phi: \mathcal{X} \rightarrow \mathcal{H}$  such that

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$$

- where  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  is the inner product in  $\mathcal{H}$ .

- Translation

- Take any set  $\mathcal{X}$  and any function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ .
- If  $k$  is a positive definite kernel, then we can use  $k$  to learn a classifier for the elements in  $\mathcal{X}$ !

- Note

- $\mathcal{X}$  can be any set, e.g.  $\mathcal{X} = \text{"all videos on YouTube"}$  or  $\mathcal{X} = \text{"all permutations of } \{1, \dots, k\}$ , or  $\mathcal{X} = \text{"the internet"}$ .

28

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide credit: Christoph Lampert

## Recap: The "Kernel Trick"

Any algorithm that uses data only in the form of inner products can be kernelized.

- How to kernelize an algorithm

- Write the algorithm only in terms of inner products.
- Replace all inner products by kernel function evaluations.

⇒ The resulting algorithm will do the same as the linear version, but in the (hidden) feature space  $\mathcal{H}$ .

- Caveat: working in  $\mathcal{H}$  is not a guarantee for better performance. A good choice of  $k$  and model selection are important!

29

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide credit: Christoph Lampert

## Recap: How to Check if a Function is a Kernel

- Problem:

- Checking if a given  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  fulfills the conditions for a kernel is difficult:

- We need to prove or disprove

$$\sum_{i,j=1}^n t_i k(x_i, x_j) t_j \geq 0$$

for any set  $x_1, \dots, x_n \in \mathcal{X}$  and any  $\mathbf{t} \in \mathbb{R}^n$  for any  $n \in \mathbb{N}$ .

- Workaround:

- It is easy to construct functions  $k$  that are positive definite kernels.

30

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

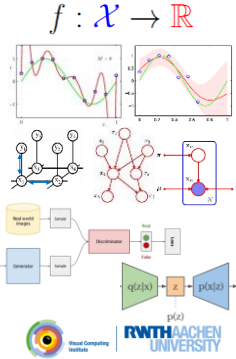


RWTH AACHEN  
UNIVERSITY

Slide credit: Christoph Lampert

## Course Outline

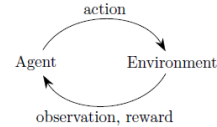
- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders



31 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

## Recap: Reinforcement Learning

- Motivation
  - General purpose framework for decision making.
  - Basis: Agent with the capability to interact with its environment
  - Each action influences the agent's future state.
  - Success is measured by a scalar reward signal.
  - Goal: select actions to maximize future rewards.

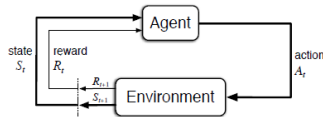


- Formalized as a partially observable Markov decision process (POMDP)

32 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

Slide adapted from: David Silver, Sergey Levine

## Recap: The Agent–Environment Interface



### Let's formalize this

- Agent and environment interact at discrete time steps  $t = 0, 1, 2, \dots$
- Agent observes state at time  $t$ :  $S_t \in \mathcal{S}$
- Produces an action at time  $t$ :  $A_t \in \mathcal{A}(S_t)$
- Gets a resulting reward  $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$
- And a resulting next state:  $S_{t+1}$

33 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

Slide adapted from: Sutton & Barto

## Recap: Reward vs. Return

- Objective of learning
  - We seek to maximize the expected return  $G_t$  as some function of the reward sequence  $R_{t+1}, R_{t+2}, R_{t+3}, \dots$
  - Standard choice: expected discounted return

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

where  $0 \leq \gamma \leq 1$  is called the discount rate.

- Difficulty
  - We don't know which past actions caused the reward.
  - ⇒ Temporal credit assignment problem

34 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

## Recap: Markov Decision Process (MDP)

- Markov Decision Processes
  - We consider decision processes that fulfill the Markov property.
  - I.e., where the environments response at time  $t$  depends only on the state and action representation at  $t$ .
- To define an MDP, we need to specify
  - State and action sets
  - One-step dynamics defined by state transition probabilities

$$p(s'|s, a) = \Pr\{S_{t+1} = s' | S_t = s, A_t = a\} = \sum_{r \in \mathcal{R}} p(s', r | s, a)$$

- Expected rewards for next state-action-next-state triplets

$$r(s, a, s') = \mathbb{E}[R_{t+1} | S_t = s, A_t = a, S_{t+1} = s'] = \frac{\sum_{r \in \mathcal{R}} r p(s', r | s, a)}{p(s'|s, a)}$$

35 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

## Recap: Policy

- Definition
  - A policy determines the agent's behavior
  - Map from state to action  $\pi: \mathcal{S} \rightarrow \mathcal{A}$
- Two types of policies
  - Deterministic policy:  $a = \pi(s)$
  - Stochastic policy:  $\pi(a|s) = \Pr\{A_t = a | S_t = s\}$
- Note
  - $\pi(a|s)$  denotes the probability of taking action  $a$  when in state  $s$ .

36 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

## Recap: Value Function

- **Idea**
  - Value function is a prediction of future reward
  - Used to evaluate the goodness/badness of states
  - And thus to select between actions
- **Definition**
  - The **value of a state**  $s$  under a policy  $\pi$ , denoted  $v_\pi(s)$ , is the expected return when starting in  $s$  and following  $\pi$  thereafter.

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s]$$

- The **value of taking action**  $a$  in state  $s$  under a policy  $\pi$ , denoted  $q_\pi(s, a)$ , is the expected return starting from  $s$ , taking action  $a$ , and following  $\pi$  thereafter.

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a]$$

37

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Optimal Value Functions

- **Bellman optimality equations**
  - For the **optimal state-value function**  $v_*$ :
 
$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a)$$

$$= \max_{a \in \mathcal{A}(s)} \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]$$
  - $v_*$  is the unique solution to this system of nonlinear equations.
  - For the **optimal action-value function**  $q_*$ :
 
$$q_*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]$$
  - $q_*$  is the unique solution to this system of nonlinear equations.
- ⇒ If the dynamics of the environment  $p(s', r | s, a)$  are known, then in principle one can solve those equation systems.

38

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Optimal Policies

- **Why optimal state-value functions are useful**
  - Any policy that is *greedy* w.r.t.  $v_*$  is an *optimal policy*.
  - ⇒ Given  $v_*$ , one-step-ahead search produces the long-term optimal results.
  - ⇒ Given  $q_*$ , we do not even have to do one-step-ahead search

$$\pi_*(s) = \operatorname{argmax}_{a \in \mathcal{A}(s)} q_*(s, a)$$

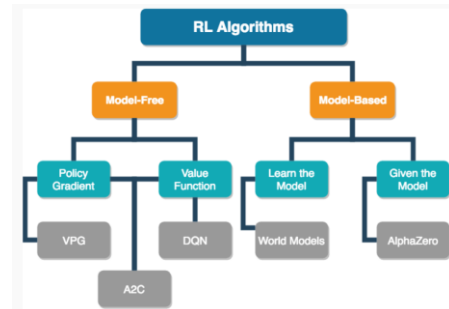
- **Challenge**
  - Many interesting problems have too many states for solving  $v_*$ .
  - Many Reinforcement Learning methods can be understood as approximately solving the Bellman optimality equations, using actually observed transitions instead of the ideal ones.

39

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Taxonomy of RL methods



40

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Tabular vs. Approximate methods

- **Tabular methods**
  - For problems with small discrete state and action spaces
  - Value function or Policy function can be expressed as a table of values.
- **Approximate methods**
  - If we cannot enumerate our states or actions we use function approximation.
  - E.g., Kernel methods, Deep Learning / Neural Networks
- **In practice, large problems with huge state spaces**
  - E.g. chess:  $10^{120}$  states.
  - Tabular methods don't scale well – they are a lookup table
    - Too many states to store in memory
    - Too slow to learn value function for every state/state-action.

41

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Model-based vs Model-free

- **Model-based**
  - Has a model of the environment dynamics and reward
  - Allows agent to plan: predict state and reward before taking action
  - Pro: Better sample efficiency
  - Con: Agent only as good as the environment - Model-bias
- **Model-free**
  - No explicit model of the environment dynamics and reward
  - Less structured. More popular and further developed and tested.
  - Pro: Can be easier to implement and tune
  - Con: Very sample inefficient

42

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Value-based RL vs Policy-based RL

- **Policy-based RL**
  - RL methods directly estimate a policy
  - A direct mapping of what action to take in each state.
 
$$\pi(a|s) = P(a|s, \theta)$$
- **Value-based RL**
  - RL methods estimate a value function and derive a policy from that
  - Either a **state-value function**

$$V(s; \theta) \approx V^\pi(s)$$
  - Or an **action-state value function** (Q function)
 
$$Q(s, a; \theta) \approx Q^\pi(s, a)$$
- Or both simultaneously: **Actor-Critic**
  - Actor-Critic methods learn both a policy (actor) and a value function (critic)

43

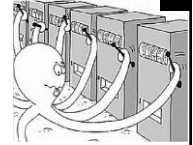
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Exploration-Exploitation Trade-off

- **Example: N-armed bandit problem**
  - Suppose we have the choice between  $N$  actions  $a_1, \dots, a_N$ .
  - If we knew their value functions  $q_*(s, a_i)$ , it would be trivial to choose the best.
  - However, we only have estimates based on our previous actions and their returns.
- We can now
  - **Exploit** our current knowledge
    - And choose the **greedy** action that has the highest value based on our current estimate.
  - **Explore** to gain additional knowledge
    - And choose a non-greedy action to improve our estimate of that action's value.



44

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Simple Action Selection Strategies

- **$\epsilon$ -greedy**
  - Select the greedy action with probability  $(1 - \epsilon)$  and a random one in the remaining cases.
  - ⇒ In the limit, every action will be sampled infinitely often.
  - ⇒ Probability of selecting the optimal action becomes  $> (1 - \epsilon)$ .
  - But: many bad actions are chosen along the way.
- **Softmax**
  - Choose action  $a_i$  at time  $t$  according to the softmax function
 
$$\frac{e^{q_t(a_i)/\tau}}{\sum_{j=1}^N e^{q_t(a_j)/\tau}}$$
  - where  $\tau$  is a temperature parameter (start high, then lower it).
  - Generalization: replace  $q_t$  by a preference function  $H_t$  that is learned by stochastic gradient ascent ("gradient bandit").

45

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: TD-Learning

- Policy evaluation (the prediction problem)
  - For a given policy  $\pi$ , compute the state-value function  $v_\pi$ .
- One option: **Monte-Carlo methods**
  - Play through a sequence of actions until a reward is reached, then backpropagate it to the states on the path.
$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

Target: the actual return after time  $t$
- **Temporal Difference Learning – TD( $\lambda$ )**
  - Directly perform an update using the estimate  $V(S_{t+\lambda+1})$ .
$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Target: an estimate of the return (here: TD(0))

46

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: SARSA – On-Policy TD Control

- **Idea**
  - Turn the TD idea into a control method by always updating the policy to be greedy w.r.t. the current estimate
- **Procedure**
  - Estimate  $q_\pi(s, a)$  for the current policy  $\pi$  and for all states  $s$  and actions  $a$ .
  - TD(0) update equation
 
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$
  - This rule is applied after every transition from a nonterminal state  $S_t$ .
  - It uses every element of the quintuple  $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$ .
  - ⇒ Hence, the name **SARSA**.

47

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Q-Learning – Off-Policy TD Control

- **Idea**
  - Directly approximate the optimal action-value function  $q_*$ , independent of the policy being followed.
- **Procedure**
  - TD(0) update equation
 
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$
  - Dramatically simplifies the analysis of the algorithm.
  - All that is required for correct convergence is that all pairs continue to be updated.

48

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY



### Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
RWTH AACHEN UNIVERSITY

### Recap: Deep Q-Learning

- Idea
  - Optimal Q-values should obey Bellman equation
 
$$Q_*(s, a) = \mathbb{E} \left[ r + \gamma \max_{a'} Q(s', a') \mid s, a \right]$$
  - Treat the right-hand side  $r + \gamma \max_{a'} Q(s', a', \mathbf{w})$  as a target
  - Minimize MSE loss by stochastic gradient descent
 
$$L(\mathbf{w}) = \left( r + \gamma \max_{a'} Q(s', a', \mathbf{w}) - Q(s, a, \mathbf{w}) \right)^2$$
  - This converges to  $Q_*$  using a lookup table representation.
  - Unfortunately, it **diverges** using neural networks due to
    - Correlations between samples
    - Non-stationary targets

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide adapted from David Silver  
RWTH AACHEN UNIVERSITY

### Recap: Deep Q-Networks (DQN)

- Adaptation: **Experience Replay**
  - To remove correlations, build a dataset from agent's own experience
 

$s_1, a_1, r_1, s'_1$
$s_2, a_2, r_2, s'_2$
$s_3, a_3, r_3, s'_3$
...
$s_t, a_t, r_t, s'_t$

 $\rightarrow s, a, r, s'$
- Perform minibatch updates to samples of experience drawn at random from the pool of stored samples
  - $(s, a, r, s') \sim U(D)$  where  $D = \{(s_t, a_t, r_{t+1}, s'_{t+1})\}$  is the dataset
- Advantages
  - Each experience sample is used in many updates (more efficient)
  - Avoids correlation effects when learning from consecutive samples
  - Avoids feedback loops from on-policy learning

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide adapted from David Silver  
RWTH AACHEN UNIVERSITY

### Recap: Deep Q-Networks (DQN)

- Adaptation: **Experience Replay**
  - To remove correlations, build a dataset from agent's own experience
 

$s_1, a_1, r_1, s'_1$
$s_2, a_2, r_2, s'_2$
$s_3, a_3, r_3, s'_3$
...
$s_t, a_t, r_t, s'_t$

 $\rightarrow s, a, r, s'$
- Sample from the dataset and apply an update
 
$$L(\mathbf{w}) = \left( r + \gamma \max_{a'} Q(s', a', \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right)^2$$
- To deal with non-stationary parameters,  $\mathbf{w}^-$  are held fixed.
  - Only update the target network parameters every  $C$  steps.
  - I.e., clone the network  $Q$  to generate a target network  $\hat{Q}$ . $\Rightarrow$  Again, this reduces oscillations to make learning more stable.

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide adapted from David Silver  
RWTH AACHEN UNIVERSITY

### Recap: Policy Gradients

- How to make high-value actions more likely
  - The gradient of a stochastic policy  $\pi(s, \mathbf{u})$  is given by
 
$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial}{\partial \mathbf{u}} \mathbb{E}_{\pi} [r_1 + \gamma r_2 + \gamma^2 r_3 + \dots \mid \pi(\cdot, \mathbf{u})]$$

$$= \mathbb{E}_{\pi} \left[ \frac{\partial \log \pi(a|s, \mathbf{u})}{\partial \mathbf{u}} Q_{\pi}(s, a) \right]$$
  - The gradient of a deterministic policy  $a = \pi(s)$  is given by
 
$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \mathbb{E}_{\pi} \left[ \frac{\partial Q_{\pi}(s, a)}{\partial a} \frac{\partial a}{\partial \mathbf{u}} \right]$$
- if  $a$  is continuous and  $Q$  is differentiable.

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide adapted from David Silver  
RWTH AACHEN UNIVERSITY

### Recap: Monte-Carlo Policy Gradient

- Execute policy to obtain sample episodes
- Update parameters by stochastic gradient ascent using the policy gradient theorem
- REINFORCE** algorithm
  - Initialize parameters arbitrarily
  - Repeat
    - Sample episode  $(s_0, a_0, r_1, s_1, a_1, \dots, r_T, s_T)$  using current policy
    - For each  $t \in \{0, \dots, T-1\}$ 
      - Update policy
 
$$\theta \leftarrow \theta + \eta \nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t) \left( \sum_{k=0}^{T-t-1} r_{t+k+1} \right)$$

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: David Silver  
RWTH AACHEN UNIVERSITY

## Recap: Deep Policy Gradients (DPG)

- DPG is the continuous analogue of DQN
  - **Experience replay**: build data-set from agent's experience
  - **Critic** estimates value of current policy by DQN

$$L_w(\mathbf{w}) = (r + \gamma Q(s', \pi(s', \mathbf{u}^-), \mathbf{w}^-) - Q(s, a, \mathbf{w}))^2$$

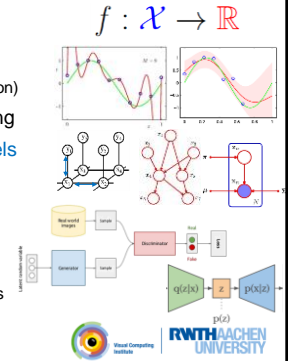
- To deal with non-stationarity, targets  $\mathbf{u}^-, \mathbf{w}^-$  are held fixed
- Actor updates policy in direction that improves Q

$$\frac{\partial L_w(\mathbf{u})}{\partial \mathbf{u}} = \frac{\partial Q(s, a, \mathbf{w})}{\partial a} \frac{\partial a}{\partial \mathbf{u}}$$

- In other words critic provides loss function for actor.

## Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders

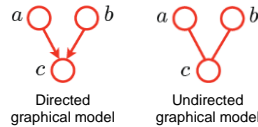


## Recap: Graphical Models

- Two basic kinds of graphical models
  - Directed graphical models or Bayesian Networks
  - Undirected graphical models or Markov Random Fields

### Key components

- **Nodes**
  - Random variables
- **Edges**
  - Directed or undirected

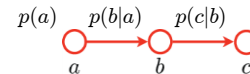


- The value of a random variable may be **known** or **unknown**.



## Recap: Directed Graphical Models

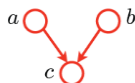
### Chains of nodes:



- Knowledge about a is expressed by the **prior probability**:  $p(a)$
- Dependencies are expressed through **conditional probabilities**:  $p(b|a), p(c|b)$
- **Joint distribution** of all three variables:  $p(a, b, c) = p(c|b)p(b|a)p(a) = p(c|b)p(b|a)p(a)$

## Recap: Directed Graphical Models

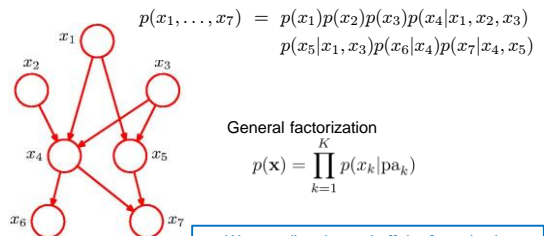
### Convergent connections:



- Here the value of  $c$  depends on both variables  $a$  and  $b$ .
- This is modeled with the conditional probability:  $p(c|a, b)$
- Therefore, the joint probability of all three variables is given as:  $p(a, b, c) = p(c|a, b)p(a)p(b) = p(c|a, b)p(a)p(b)$

## Recap: Factorization of the Joint Probability

### Exercise: Computing the joint probability



We can directly read off the factorization of the joint from the network structure!

## Recap: Factorized Representation

- **Reduction of complexity**
  - Joint probability of  $n$  binary variables requires us to represent values by brute force

$$\mathcal{O}(2^n) \text{ terms}$$

- The factorized form obtained from the graphical model only requires

$$\mathcal{O}(n \cdot 2^k) \text{ terms}$$

- $k$ : maximum number of parents of a node.

→ It's the edges that are missing in the graph that are important! They encode the simplifying assumptions we make.

61

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide credit: Bernd Schiele, Stefan Roth

## Recap: Conditional Independence

- $X$  is **conditionally independent** of  $Y$  given  $V$

– Definition:  $X \perp\!\!\!\perp Y | V \Leftrightarrow p(X|Y, V) = p(X|V)$

– Also:  $X \perp\!\!\!\perp Y | V \Leftrightarrow p(X, Y | V) = p(X|V)p(Y|V)$

- Special case: **Marginal Independence**

$$X \perp\!\!\!\perp Y \Leftrightarrow X \perp\!\!\!\perp Y | \emptyset \Leftrightarrow p(X, Y) = p(X)p(Y)$$

- Often, we are interested in conditional independence between sets of variables:

$$\mathcal{X} \perp\!\!\!\perp \mathcal{Y} | \mathcal{V} \Leftrightarrow \{X \perp\!\!\!\perp Y | \mathcal{V}, \forall X \in \mathcal{X} \text{ and } \forall Y \in \mathcal{Y}\}$$

see  
Exercise 3.2

62

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

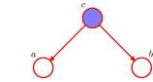


RWTH AACHEN  
UNIVERSITY

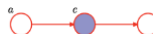
## Recap: Conditional Independence

- Three cases

- **Divergent** ("Tail-to-Tail")
  - Conditional independence when  $c$  is observed.



- **Chain** ("Head-to-Tail")
  - Conditional independence when  $c$  is observed.



- **Convergent** ("Head-to-Head")
  - Conditional independence when **neither  $c$ , nor any of its descendants** are observed.



63

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Image source: C. Bishop, 2006

## Recap: D-Separation

- **Definition**

- Let  $A$ ,  $B$ , and  $C$  be non-intersecting subsets of nodes in a directed graph.

- A path from  $A$  to  $B$  is **blocked** if it contains a node such that either

- The arrows on the path meet either **head-to-tail** or **tail-to-tail** at the node, and the **node is in the set  $C$** , or
- The arrows meet **head-to-head** at the node, and **neither the node, nor any of its descendants, are in the set  $C$** .



- If all paths from  $A$  to  $B$  are blocked,  $A$  is said to be **d-separated** from  $B$  by  $C$ .

- If  $A$  is d-separated from  $B$  by  $C$ , the joint distribution over all variables in the graph satisfies  $A \perp\!\!\!\perp B | C$ .

- Read: " $A$  is conditionally independent of  $B$  given  $C$ ."

see  
Exercise 3.2

64

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Chris Bishop

## Recap: "Bayes Ball" Algorithm

- Graph algorithm to compute d-separation

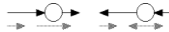
- Goal: Get a ball from  $X$  to  $Y$  without being blocked by  $V$ .

- Depending on its direction and the previous node, the ball can

- **Pass through** (from parent to all children, from child to all parents)
- **Bounce back** (from any parent/child to all parents/children)
- **Be blocked**

- Game rules

- An **unobserved** node ( $W \notin \mathcal{V}$ ) **passes through** balls from parents, but **also bounces back** balls from children.



- An **observed** node ( $W \in \mathcal{V}$ ) **bounces back** balls from parents, but **blocks** balls from children.



65

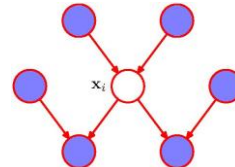
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Zoubin Ghahramani

## Recap: The Markov Blanket



- **Markov blanket** of a node  $x_i$

- Minimal set of nodes that isolates  $x_i$  from the rest of the graph.

- This comprises the set of

- Parents,
- Children, and
- Co-parents of  $x_i$ .

← This is what we have to watch out for!

66

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Image source: C. Bishop, 2006

## Recap: D-Separation

### • Definition

- Let  $A$ ,  $B$ , and  $C$  be non-intersecting subsets of nodes in a directed graph.
- A path from  $A$  to  $B$  is **blocked** if it contains a node such that either
  - The arrows on the path meet either **head-to-tail** or **tail-to-tail** at the node, and the **node is in the set  $C$** ; or
  - The arrows meet **head-to-head** at the node, and **neither the node, nor any of its descendants, are in the set  $C$** .
- If all paths from  $A$  to  $B$  are blocked,  $A$  is said to be **d-separated** from  $B$  by  $C$ .
- If  $A$  is d-separated from  $B$  by  $C$ , the joint distribution over all variables in the graph satisfies  $A \perp\!\!\!\perp B \mid C$ .
  - Read: “ $A$  is conditionally independent of  $B$  given  $C$ .”



67

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

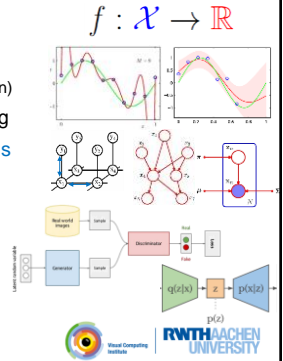


RWTH AACHEN  
UNIVERSITY

Slide adapted from Chris Bishop.

## Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders



68

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

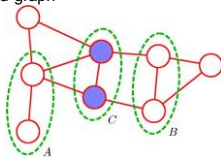


RWTH AACHEN  
UNIVERSITY

## Recap: Undirected Graphical Models

### • Undirected graphical models (“Markov Random Fields”)

- Given by undirected graph



### • Conditional independence for undirected graphs

- If every path from any node in set  $A$  to set  $B$  passes through at least one node in set  $C$ , then  $A \perp\!\!\!\perp B \mid C$ .
- Simple Markov blanket:



69

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Image source: C. Bishop, 2006

## Recap: Factorization in MRFs

### • Joint distribution

- Written as product of **potential functions** over **maximal cliques** in the graph:

$$p(\mathbf{x}) = \frac{1}{Z} \prod_C \psi_C(\mathbf{x}_C)$$

- The normalization constant  $Z$  is called the **partition function**.

$$Z = \sum_{\mathbf{x}} \prod_C \psi_C(\mathbf{x}_C)$$

### • Remarks

- BNs are automatically normalized. But for MRFs, we have to explicitly perform the normalization.
- Presence of normalization constant is major limitation!
  - Evaluation of  $Z$  involves summing over  $\mathcal{O}(K^M)$  terms for  $M$  nodes!

70

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Factorization in MRFs

### • Role of the potential functions

- General interpretation
  - No restriction to potential functions that have a specific probabilistic interpretation as marginals or conditional distributions.
- Convenient to express them as exponential functions (“**Boltzmann distribution**”)

$$\psi_C(\mathbf{x}_C) = \exp\{-E(\mathbf{x}_C)\}$$

- with an **energy function**  $E$ .
- Why is this convenient?
  - Joint distribution is the product of potentials  $\Rightarrow$  sum of energies.
  - We can take the log and simply work with the sums...

71

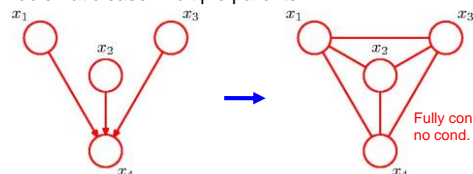
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Converting Directed to Undirected Graphs

### • Problematic case: multiple parents



$$p(\mathbf{x}) = p(x_1)p(x_2)p(x_3)p(x_4|x_1, x_2, x_3)$$

Need a clique of  $x_1, \dots, x_4$  to represent this factor!

- Need to introduce additional links (“**marry the parents**”).
- $\Rightarrow$  This process is called **moralization**. It results in the **moral graph**.

72

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Chris Bishop

Image source: C. Bishop, 2006

## Recap: Conversion Algorithm

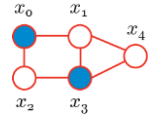
- General procedure to convert directed  $\rightarrow$  undirected
  1. Add undirected links to **marry the parents** of each node.
  2. **Drop the arrows** on the original links  $\Rightarrow$
  3. **Find maximal cliques** for each node and initialize all clique potentials to 1.
  4. Take each conditional distribution factor of the original directed graph and multiply it into one clique potential.
- **Restriction**
  - Conditional independence properties are often lost!
  - Moralization results in additional connections and larger cliques.

73 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide adapted from Chris Bishop.



## Recap: Computing Marginals

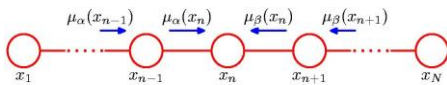
- How do we apply graphical models?
  - Given some observed variables, we want to **compute distributions of the unobserved variables**.
  - In particular, we want to **compute marginal distributions**, for example  $p(x_i)$ .
- How can we compute marginals?
  - Classical technique: **sum-product algorithm** by Judea Pearl.
  - In the context of (loopy) undirected models, this is also called (loopy) **belief propagation** [Weiss, 1997].
  - Basic idea: **message-passing**.



74 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernt Schiele, Stefan Roth.



## Recap: Message Passing on a Chain



- Idea: **Pass messages** from the two ends towards the query node  $x_n$ .
- Define the messages recursively:

$$\mu_\alpha(x_n) = \sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_n) \mu_\alpha(x_{n-1})$$

$$\mu_\beta(x_n) = \sum_{x_{n+1}} \psi_{n,n+1}(x_n, x_{n+1}) \mu_\beta(x_{n+1})$$

- Compute the normalization constant  $Z$  at any node  $x_m$ .

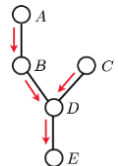
$$Z = \sum_{x_n} \mu_\alpha(x_n) \mu_\beta(x_n)$$

75 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide adapted from Chris Bishop.



## Summary: Message Passing on Trees

- **General procedure** for all tree graphs.
  - Root the tree at the variable that we want to compute the marginal of.
  - Start computing messages at the leaves.
  - Compute the messages for all nodes for which all incoming messages have already been computed.
  - Repeat until we reach the root.
- If we want to compute the marginals for all possible nodes (roots), we can reuse some of the messages.
  - Computational expense linear in the number of nodes.
- We already motivated message passing for inference.
  - How can we formalize this into a general algorithm?

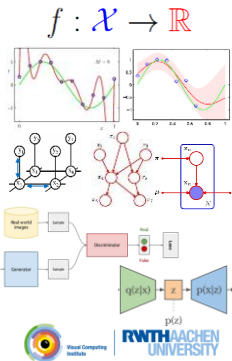


76 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernt Schiele, Stefan Roth.



## Course Outline

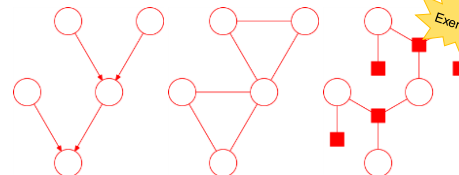
- **Regression Techniques**
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- **Deep Reinforcement Learning**
- **Probabilistic Graphical Models**
  - Bayesian Networks
  - Markov Random Fields
  - **Inference (exact & approximate)**
  - Latent Variable Models
- **Deep Generative Models**
  - Generative Adversarial Networks
  - Variational Autoencoders



77 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Recap: Factor Graphs



- **Joint probability**
  - Can be expressed as **product of factors**:  $p(\mathbf{x}) = \frac{1}{Z} \prod_s f_s(\mathbf{x}_s)$
  - Factor graphs make this explicit through separate factor nodes.
- **Converting a directed polytree**
  - Conversion to undirected tree creates loops due to moralization!
  - **Conversion to a factor graph again results in a tree!**

78 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



Image source: C. Bishop, 2006

## Recap: Sum-Product Algorithm

- Objectives
  - Efficient, **exact inference** algorithm for finding marginals.
- Procedure:
  - Pick an **arbitrary node** as root.
  - Compute and propagate messages **from the leaf nodes to the root**, storing received messages at every node.
  - Compute and propagate messages **from the root to the leaf nodes**, storing received messages at every node.
  - Compute the **product of received messages at each node** for which the marginal is required, and normalize if necessary.

$$p(x) \propto \prod_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x)$$

- Computational effort
  - Total number of messages = 2 · number of graph edges.

79

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Chris Bishop.

## Recap: Sum-Product Algorithm

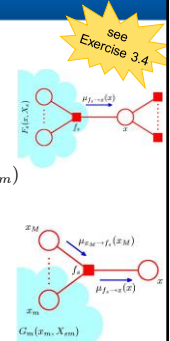
- Two kinds of messages
  - Message from factor node to variable nodes:
    - Sum of factor contributions

$$\begin{aligned} \mu_{f_s \rightarrow x}(x) &\equiv \sum_{X_s} F_s(x, X_s) \\ &= \sum_{X_s} f_s(\mathbf{x}_s) \prod_{m \in \text{ne}(f_s) \setminus x} \mu_{x_m \rightarrow f_s}(x_m) \end{aligned}$$

- Message from variable node to factor node:
  - Product of incoming messages

$$\mu_{x_m \rightarrow f_s}(x_m) \equiv \prod_{l \in \text{ne}(x_m) \setminus f_s} \mu_{f_l \rightarrow x_m}(x_m)$$

⇒ Simple propagation scheme.



80

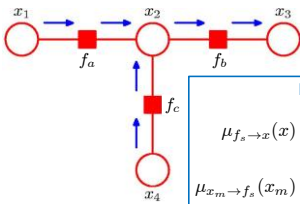
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Chris Bishop.

## Recap: Sum-Product from Leaves to Root



Message definitions:

$$\begin{aligned} \mu_{f_s \rightarrow x}(x) &\equiv \sum_{X_s} f_s(\mathbf{x}_s) \prod_{m \in \text{ne}(f_s) \setminus x} \mu_{x_m \rightarrow f_s}(x_m) \\ \mu_{x_m \rightarrow f_s}(x_m) &\equiv \prod_{l \in \text{ne}(x_m) \setminus f_s} \mu_{f_l \rightarrow x_m}(x_m) \end{aligned}$$

$$\begin{aligned} \mu_{x \rightarrow f}(x) &= 1 & \mu_{f \rightarrow x}(x) &= f(x) \end{aligned}$$

81

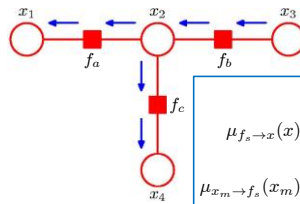
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition



RWTH AACHEN  
UNIVERSITY

Image source: C. Bishop, 2006

## Recap: Sum-Product from Root to Leaves



Message definitions:

$$\begin{aligned} \mu_{f_s \rightarrow x}(x) &\equiv \sum_{X_s} f_s(\mathbf{x}_s) \prod_{m \in \text{ne}(f_s) \setminus x} \mu_{x_m \rightarrow f_s}(x_m) \\ \mu_{x_m \rightarrow f_s}(x_m) &\equiv \prod_{l \in \text{ne}(x_m) \setminus f_s} \mu_{f_l \rightarrow x_m}(x_m) \end{aligned}$$

$$\begin{aligned} \mu_{x \rightarrow f}(x) &= 1 & \mu_{f \rightarrow x}(x) &= f(x) \end{aligned}$$

82

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition



RWTH AACHEN  
UNIVERSITY

Image source: C. Bishop, 2006

## Recap: Max-Sum Algorithm

- Objective: an efficient algorithm for finding
  - Value  $\mathbf{x}^{\max}$  that maximises  $p(\mathbf{x})$ ;
  - Value of  $p(\mathbf{x}^{\max})$ .
 ⇒ Application of dynamic programming in graphical models.
- Key ideas
  - We are interested in the **maximum value of the joint distribution**

$$p(\mathbf{x}^{\max}) = \max_{\mathbf{x}} p(\mathbf{x})$$
 ⇒ Maximize the product  $p(\mathbf{x})$ .
  - For numerical reasons, use the logarithm.
 
$$\ln \left( \max_{\mathbf{x}} p(\mathbf{x}) \right) = \max_{\mathbf{x}} \ln p(\mathbf{x}).$$
 ⇒ Maximize the sum (of log-probabilities).

83

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Chris Bishop.

## Recap: Max-Sum Algorithm

- Initialization (leaf nodes)

$$\mu_{f \rightarrow x}(x) = 0 \quad \mu_{f \rightarrow x}(x) = \ln f(x)$$

- Recursion

$$\mu_{f \rightarrow x}(x) = \max_{x_1, \dots, x_M} \left[ \ln f(x, x_1, \dots, x_M) + \sum_{m \in \text{ne}(f_s) \setminus x} \mu_{x_m \rightarrow f}(x_m) \right]$$

$$\mu_{x \rightarrow f}(x) = \sum_{l \in \text{ne}(x) \setminus f} \mu_{f_l \rightarrow x}(x)$$

- For each node, keep a record of which values of the variables gave rise to the maximum state:

$$\phi(x) = \arg \max_{x_1, \dots, x_M} \left[ \ln f(x, x_1, \dots, x_M) + \sum_{m \in \text{ne}(f_s) \setminus x} \mu_{x_m \rightarrow f}(x_m) \right]$$

84

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Chris Bishop.

## Recap: Max-Sum Algorithm

- Termination (root node)

- Score of maximal configuration

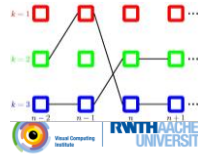
$$p^{\max} = \max_x \left[ \sum_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x) \right]$$

- Value of root node variable giving rise to that maximum

$$x^{\max} = \arg \max_x \left[ \sum_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x) \right]$$

- Back-track to get the remaining variable values

$$x_{n-1}^{\max} = \phi(x_n^{\max})$$



## Recap: Junction Tree Algorithm

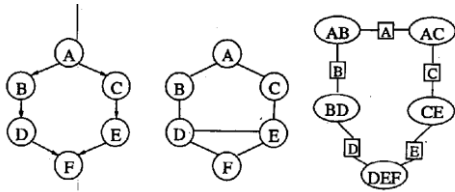
- Motivation

- Exact inference on general graphs.
- Works by turning the initial graph into a **junction tree** and then running a sum-product-like algorithm.
- Intractable on graphs with large cliques.

- Main steps

1. If starting from directed graph, first convert it to an undirected graph by **moralization**.
2. Introduce additional links by **triangulation** in order to reduce the size of cycles.
3. Find **cliques** of the moralized, triangulated graph.
4. Construct a new graph from the **maximal cliques**.
5. Remove minimal links to **break cycles** and get a **junction tree**.  
⇒ Apply regular **message passing** to perform inference.

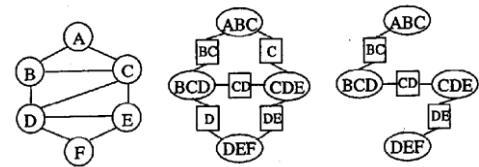
## Recap: Junction Tree Example



- Without triangulation step

- The final graph will contain cycles that we cannot break without losing the running intersection property!

## Recap: Junction Tree Example

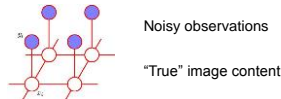


- When applying the triangulation

- Only small cycles remain that are easy to break.
- Running intersection property is maintained.

## Recap: MRF Structure for Images

- Basic structure



- Two components

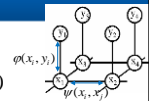
- Observation model
  - How likely is it that node  $x_i$  has label  $L_i$  given observation  $y_i$ ?
  - This relationship is usually learned from training data.
- Neighborhood relations
  - Simplest case: 4-neighborhood
  - Serve as smoothing terms.
  - ⇒ Discourage neighboring pixels to have different labels.
  - This can either be learned or be set to fixed "penalties".



## Recap: Energy Formulation

- Energy function

$$E(x, y) = \sum_i \underbrace{\varphi(x_i, y_i)}_{\text{Single-node potentials}} + \sum_{i,j} \underbrace{\psi(x_i, x_j)}_{\text{Pairwise potentials}}$$



- Single-node (unary) potentials  $\varphi$

- Encode local information about the given pixel/patch.
- How likely is a pixel/patch to belong to a certain class (e.g. foreground/background)?

- Pairwise potentials  $\psi$

- Encode neighborhood information.
- How different is a pixel/patch's label from that of its neighbor? (e.g. based on intensity/color/texture difference, edges)

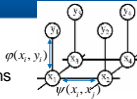
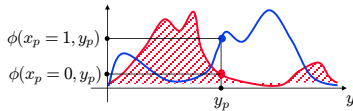
## Recap: How to Set the Potentials?

### Unary potentials

- E.g. color model, modeled with a Mixture of Gaussians

$$\phi(x_i, y_i; \theta_\phi) = \log \sum_k \theta_\phi(x_i, k) p(k|x_i) \mathcal{N}(y_i; \bar{y}_k, \Sigma_k)$$

⇒ Learn color distributions for each label



91

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: How to Set the Potentials?

### Pairwise potentials

- Potts Model

$$\psi(x_i, x_j; \theta_\psi) = \theta_\psi \delta(x_i \neq x_j)$$

- Simplest discontinuity preserving model.
- Discontinuities between any pair of labels are penalized equally.
- Useful when labels are unordered or number of labels is small.

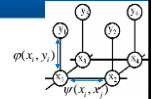
- Extension: “contrast sensitive Potts model”

$$\psi(x_i, x_j, g_{ij}(y); \theta_\psi) = \theta_\psi g_{ij}(y) \delta(x_i \neq x_j)$$

where

$$g_{ij}(y) = e^{-\beta \|y_i - y_j\|^2} \quad \beta = 2 \cdot \text{avg}(\|y_i - y_j\|^2)$$

- Discourages label changes except in places where there is also a large change in the observations.



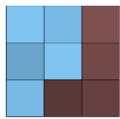
92

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



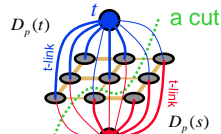
RWTH AACHEN  
UNIVERSITY

## Recap: Graph Cuts for Binary Problems



“expected” intensities of  
object and background  
 $I^s$  and  $I^t$   
can be re-estimated

EM-style optimization



$$D_p(s) \propto \exp(-\|I_p - I^s\|^2 / 2\sigma^2)$$

$$D_p(t) \propto \exp(-\|I_p - I^t\|^2 / 2\sigma^2)$$

93

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

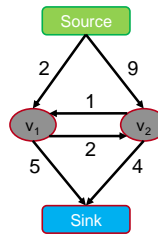
Slide credit: Yuri Boykov

[Boykov & Jolly, ICCV'01]

## Recap: Maxflow Algorithms

Flow = 0

Augmenting Path Based Algorithms



1. Find path from source to sink with positive capacity
2. Push maximum possible flow through this path
3. Adjust the capacity of the used edges and record “residual flows”
4. Repeat until no path can be found

Algorithms assume non-negative capacity

94

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide credit: Pushmeet Kohli

## Recap: When Can s-t Graph Cuts Be Applied?

$$E(L) = \sum_p E_p(L_p) + \sum_{p,q \in N} E(L_p, L_q)$$

t-links                      n-links                       $L_p \in \{s, t\}$

- s-t graph cuts can only globally minimize binary energies that are submodular. [Boros & Hummer, 2002, Kolmogorov & Zabih, 2004]

$$E(L) \text{ can be minimized by s-t graph cuts} \iff E(s,s) + E(t,t) \leq E(s,t) + E(t,s)$$

Submodularity (“convexity”)

- Submodularity is the discrete equivalent to convexity.
  - Implies that every local energy minimum is a global minimum.
  - ⇒ Solution will be globally optimal.

95

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

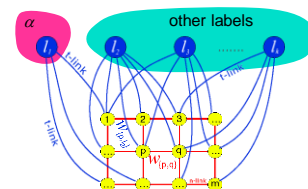


RWTH AACHEN  
UNIVERSITY

## Recap: $\alpha$ -Expansion

### Basic idea:

- Break multi-way cut computation into a sequence of binary s-t cuts.



96

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

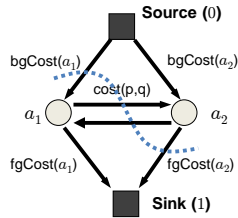
Slide credit: Yuri Boykov



## Recap: Converting an MRF to an s-t Graph

```

Graph *g;
For all pixels p
    /* Add a node to the graph */
    nodeID(p) = g->add_node();
    /* Set cost of terminal edges */
    set_weights(nodeID(p), fgCost(p), bgCost(p));
end
for all adjacent pixels p,q
    add_weights(nodeID(p), nodeID(q), cost);
end
g->compute_maxflow();
label_p = g->is_connected_to_source(nodeID(p));
// is the label of pixel p (0 or 1)
    
```



$$a_1 = \text{bg} \quad a_2 = \text{fg}$$

97

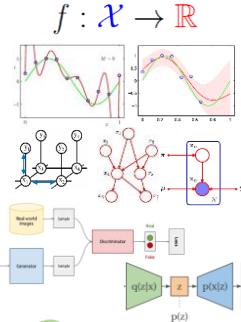
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition

Slide credit: Pushmeet Kohli



## Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders



98

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition

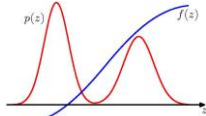


## Recap: Sampling Idea

### 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  $\mathbf{z}$ !

99

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition

Slide adapted from Bernt Schiele



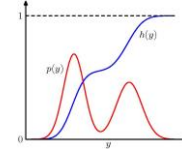
## Recap: Transformation Method

- 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:

$$u \sim \text{Uniform}(0, 1) \Rightarrow F^{-1}(u) \sim p(x)$$



100

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition

Slide credit: Bernt Schiele



## Recap: Rejection Sampling

### Assumptions

- Sampling directly from  $p(\mathbf{z})$  is difficult.
- But we can easily evaluate  $p(\mathbf{z})$  (up to some norm. factor  $Z_p$ ):

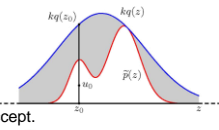
$$p(\mathbf{z}) = \frac{1}{Z_p} \tilde{p}(\mathbf{z})$$

### Idea

- We need some simpler distribution  $q(\mathbf{z})$  (called proposal distribution) from which we can draw samples.
- Choose a constant  $k$  such that:  $\forall \mathbf{z} : kq(\mathbf{z}) \geq \tilde{p}(\mathbf{z})$

### 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.



101

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition

Slide adapted from Bernt Schiele



## Recap: Importance Sampling

### Approach

- Approximate expectations directly  $\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z}$  (but does not enable to draw samples from  $p(\mathbf{z})$  directly).

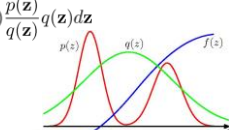
### Idea

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

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

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

Importance weights



102

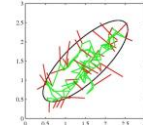
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 - Repetition

Slide adapted from Bernt Schiele



## Recap: MCMC – Markov Chain Monte Carlo

- **Overview**
  - Allows to sample from a large class of distributions.
  - Scales well with the dimensionality of the sample space.
- **Idea**
  - We maintain a record of the current state  $\mathbf{z}^{(r)}$
  - The proposal distribution depends on the current state:  $q(\mathbf{z}|\mathbf{z}^{(r)})$
  - The sequence of samples forms a Markov chain  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots$
- **Approach**
  - At each time step, we generate a candidate sample from the proposal distribution and accept the sample according to a criterion.
  - Different variants of MCMC for different criteria.



103 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide adapted from Bernt Schiele



RWTH AACHEN UNIVERSITY  
Image source: C.M. Bishop, 2006

## Recap: Markov Chains – Properties

- **Invariant distribution**
  - A distribution is said to be **invariant** (or **stationary**) w.r.t. a Markov chain if each step in the chain leaves that distribution invariant.
  - Transition probabilities:
$$T(\mathbf{z}^{(m)}, \mathbf{z}^{(m+1)}) = p(\mathbf{z}^{(m+1)}|\mathbf{z}^{(m)})$$
  - For homogeneous Markov chain, distribution  $p^*(\mathbf{z})$  is invariant if:
$$p^*(\mathbf{z}) = \sum_{\mathbf{z}'} T(\mathbf{z}', \mathbf{z}) p^*(\mathbf{z}')$$
- **Detailed balance**
  - Sufficient (but not necessary) condition to ensure that a distribution is invariant:
$$p^*(\mathbf{z})T(\mathbf{z}, \mathbf{z}') = p^*(\mathbf{z}')T(\mathbf{z}', \mathbf{z})$$
  - A Markov chain which respects *detailed balance* is **reversible**.

104 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernt Schiele



RWTH AACHEN UNIVERSITY

## Recap: Detailed Balance

- **Detailed balance** means
  - If we pick a state from the target distribution  $p(\mathbf{z})$  and make a transition under  $T$  to another state, it is just as likely that we will pick  $\mathbf{z}_A$  and go from  $\mathbf{z}_A$  to  $\mathbf{z}_B$  than that we will pick  $\mathbf{z}_B$  and go from  $\mathbf{z}_B$  to  $\mathbf{z}_A$ .
  - It can easily be seen that a transition probability that satisfies detailed balance w.r.t. a particular distribution will leave that distribution invariant, because

$$\begin{aligned} \sum_{\mathbf{z}'} p^*(\mathbf{z}')T(\mathbf{z}', \mathbf{z}) &= \sum_{\mathbf{z}'} p^*(\mathbf{z})T(\mathbf{z}, \mathbf{z}') \\ &= p^*(\mathbf{z}) \sum_{\mathbf{z}'} p(\mathbf{z}'|\mathbf{z}) = p^*(\mathbf{z}) \end{aligned}$$

105 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN UNIVERSITY

## Recap: MCMC – Metropolis Algorithm

see Exercise 4.4

- **Metropolis algorithm** [Metropolis et al., 1953]
  - Proposal distribution is symmetric:  $q(\mathbf{z}_A|\mathbf{z}_B) = q(\mathbf{z}_B|\mathbf{z}_A)$
  - The new candidate sample  $\mathbf{z}^*$  is accepted with probability
$$A(\mathbf{z}^*, \mathbf{z}^{(r)}) = \min\left(1, \frac{\tilde{p}(\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(r)})}\right)$$
  - ⇒ New candidate samples always accepted if  $\tilde{p}(\mathbf{z}^*) \geq \tilde{p}(\mathbf{z}^{(r)})$
  - The algorithm sometimes accepts a state with lower probability.
- **Metropolis-Hastings algorithm**
  - Generalization: Proposal distribution not necessarily symmetric.
  - The new candidate sample  $\mathbf{z}^*$  is accepted with probability
$$A(\mathbf{z}^*, \mathbf{z}^{(r)}) = \min\left(1, \frac{\tilde{p}(\mathbf{z}^*)q_k(\mathbf{z}^{(r)}|\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(r)})q_k(\mathbf{z}^*|\mathbf{z}^{(r)})}\right)$$
  - where  $k$  labels the members of the set of considered transitions.

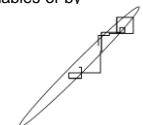
106 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide adapted from Bernt Schiele



RWTH AACHEN UNIVERSITY

## Recap: Gibbs Sampling

- **Approach**
  - MCMC-algorithm that is simple and widely applicable.
  - May be seen as a special case of Metropolis-Hastings.
- **Idea**
  - Sample variable-wise: replace  $z_i$  by a value drawn from the distribution  $p(z_i|\mathbf{z}_{-i})$ .
    - This means we update one coordinate at a time.
  - Repeat procedure either by cycling through all variables or by choosing the next variable.
- **Properties**
  - The algorithm **always** accepts!
  - Completely parameter free.
  - Can also be applied to subsets of variables.



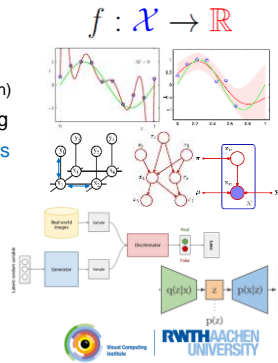
107 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernt Schiele



RWTH AACHEN UNIVERSITY

## Course Outline

- **Regression Techniques**
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- **Deep Reinforcement Learning**
- **Probabilistic Graphical Models**
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- **Deep Generative Models**
  - Generative Adversarial Networks
  - Variational Autoencoders



108 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN UNIVERSITY

### Recap: Mixtures of Gaussians

- "Generative model"
 
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$p(j) = \pi_j$   
 $p(\mathbf{x} | \theta_j) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$   
 $p(\mathbf{x}) = \sum_{j=1}^3 \pi_j p(\mathbf{x} | \theta_j)$

109 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Bernt Schiele  
Image source: C.M. Bishop, 2006

### Recap: GMMs as Latent Variable Models

- Write GMMs in terms of latent variables  $\mathbf{z}$ 
  - Marginal distribution of  $\mathbf{x}$ 

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z}} p(\mathbf{z}) p(\mathbf{x} | \mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
- Advantage of this formulation
  - We have represented the marginal distribution in terms of **latent variables  $\mathbf{z}$** .
  - Since  $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$ , there is a corresponding latent variable  $\mathbf{z}_n$  for each data point  $\mathbf{x}_n$ .
  - We are now able to work with the joint distribution  $p(\mathbf{x}, \mathbf{z})$  instead of the marginal distribution  $p(\mathbf{x})$ .

⇒ This will lead to significant simplifications...

110 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Image source: C.M. Bishop, 2006

### Recap: Sampling from a Gaussian Mixture

- MoG Sampling
  - We can use **ancestral sampling** to generate random samples from a Gaussian mixture model.
    1. Generate a value  $\tilde{\mathbf{z}}$  from the marginal distribution  $p(\mathbf{z})$ .
    2. Generate a value  $\tilde{\mathbf{x}}$  from the conditional distribution  $p(\mathbf{x} | \tilde{\mathbf{z}})$ .

111 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Image source: C.M. Bishop, 2006

### Recap: Gaussian Mixtures Revisited

- Applying the latent variable view of EM
  - Goal is to maximize the log-likelihood using the observed data  $\mathbf{X}$ 

$$\log p(\mathbf{X} | \boldsymbol{\theta}) = \log \left\{ \sum_{\mathbf{z}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) \right\}$$
  - Corresponding graphical model:
  - Suppose we are additionally given the values of the latent variables  $\mathbf{Z}$ .
  - The corresponding graphical model for the complete data now looks like this:
- ⇒ Straightforward to marginalize...

112 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Image source: C.M. Bishop, 2006

### Recap: Alternative View of EM

- In practice, however...
  - We are not given the complete data set  $\{\mathbf{X}, \mathbf{Z}\}$ , but only the incomplete data  $\mathbf{X}$ . All we can compute about  $\mathbf{Z}$  is the posterior distribution  $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})$ .
  - Since we cannot use the complete-data log-likelihood, we consider instead its **expected value under the posterior distribution of the latent variables**:
 
$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{z}} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$$
  - This corresponds to the **E-step** of the EM algorithm.
  - In the subsequent **M-step**, we then maximize the expectation to obtain the revised parameter set  $\boldsymbol{\theta}^{\text{new}}$ .
 
$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$$

113 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Image source: C.M. Bishop, 2006

### Recap: General EM Algorithm

- Algorithm
  1. Choose an initial setting for the parameters  $\boldsymbol{\theta}^{\text{old}}$
  2. **E-step**: Evaluate  $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$
  3. **M-step**: Evaluate  $\boldsymbol{\theta}^{\text{new}}$  given by
 
$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$$
 where
 
$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{z}} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$$
  4. While not converged, let  $\boldsymbol{\theta}^{\text{old}} \leftarrow \boldsymbol{\theta}^{\text{new}}$  and return to step 2.

114 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Image source: C.M. Bishop, 2006

## Recap: MAP-EM

- **Modification for MAP**
  - The EM algorithm can be adapted to find MAP solutions for models for which a prior  $p(\theta)$  is defined over the parameters.
  - Only changes needed:
    2. **E-step:** Evaluate  $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$
    3. **M-step:** Evaluate  $\theta^{\text{new}}$  given by
 
$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}) + \log p(\theta)$$
- ⇒ Suitable choices for the prior will remove the ML singularities!

## Recap: Monte Carlo EM

- **EM procedure**
  - **M-step:** Maximize expectation of complete-data log-likelihood
 
$$Q(\theta, \theta^{\text{old}}) = \int p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\theta) d\mathbf{Z}$$
  - For more complex models, we may not be able to compute this analytically anymore...
- **Idea**
  - Use sampling to approximate this integral by a finite sum over samples  $\{\mathbf{Z}^{(l)}\}$  drawn from the current estimate of the posterior
 
$$Q(\theta, \theta^{\text{old}}) \sim \frac{1}{L} \sum_{l=1}^L \log p(\mathbf{X}, \mathbf{Z}^{(l)}|\theta)$$
  - This procedure is called the **Monte Carlo EM algorithm**.

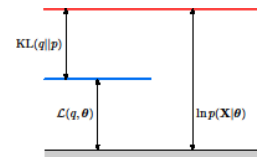
## Recap: EM as Variational Inference

- **Decomposition**
  - Introduce a distribution  $q(\mathbf{Z})$  over the latent variables. For any choice of  $q(\mathbf{Z})$ , the following decomposition holds
 
$$\log p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + KL(q \parallel p)$$
  - where
 
$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \right\}$$

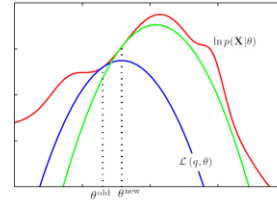
$$KL(q \parallel p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \theta)}{q(\mathbf{Z})} \right\}$$
  - $KL(q \parallel p)$  is the **Kullback-Leibler divergence** between the distribution  $q(\mathbf{Z})$  and the posterior distribution  $p(\mathbf{Z}|\mathbf{X}, \theta)$ .
  - $\mathcal{L}(q, \theta)$  is a **functional** of the distribution  $q(\mathbf{Z})$  and a function of the parameters  $\theta$ . Since  $KL \geq 0$ ,  $\mathcal{L}(q, \theta)$  is a **lower bound** on  $\log p(\mathbf{X}|\theta)$ .

## Recap: Analysis of EM

- **Decomposition**

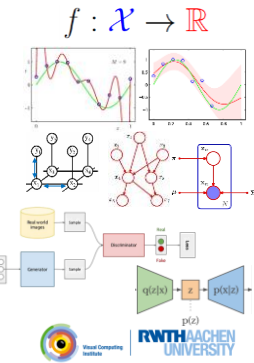
$$\log p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + KL(q \parallel p)$$

- **Interpretation**
  - $\mathcal{L}(q, \theta)$  is a **lower bound** on  $\log p(\mathbf{X}|\theta)$ .
  - The approximation comes from the fact that we use an approximative distribution  $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$  instead of the (unknown) real posterior.
  - The KL divergence measures the difference between the approximative distribution  $q(\mathbf{Z})$  and the real posterior  $p(\mathbf{Z}|\mathbf{X}, \theta)$ .
  - In every EM iteration, we try to make this difference smaller.

## Recap: Analysis of EM

- **Visualization in the space of parameters**

- **The EM algorithm alternately...**
  - Computes a lower bound on the log-likelihood for the current parameters values
  - And then maximizes this bound to obtain the new parameter values.

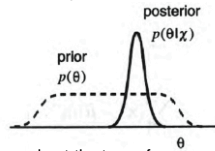
## Course Outline

- **Regression Techniques**
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- **Deep Reinforcement Learning**
- **Probabilistic Graphical Models**
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - **Bayesian Latent Variable Models**
- **Deep Generative Models**
  - Generative Adversarial Networks
  - Variational Autoencoders



## Recap: Bayesian Estimation

- **Conceptual shift**
  - **Maximum Likelihood** views the true parameter vector  $\theta$  to be unknown, but fixed.
  - In **Bayesian learning**, we consider  $\theta$  to be a random variable.
- This allows us to use knowledge about the parameters  $\theta$ 
  - i.e., to use a prior for  $\theta$
  - Training data then converts this prior distribution on  $\theta$  into a posterior probability density.



- The prior thus encodes knowledge we have about the type of distribution we expect to see for  $\theta$ .

## Recap: Bayesian Estimation

- **Discussion**

Likelihood of the parametric form  $\theta$  given the data set  $X$ .

Estimate for  $x$  based on parametric form  $\theta$       Prior for the parameters  $\theta$

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

Normalization: integrate over all possible values of  $\theta$

⇒ The parameter values  $\theta$  are not the goal, just a means to an end.

## Recap: 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.
- What to do when we cannot use the conjugate prior?
  - ⇒ Use approximate inference methods.

## Recap: The Dirichlet Distribution

- **Dirichlet Distribution**

- **Conjugate prior** for the Categorical and the Multinomial distrib.

$$\text{Dir}(\mu|\alpha) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1} \quad \text{with} \quad \alpha_0 = \sum_{k=1}^K \alpha_k$$

- **Symmetric version (with concentration parameter  $\alpha$ )**

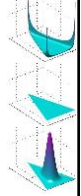
$$\text{Dir}(\mu|\alpha) = \frac{\Gamma(\alpha)^K}{\Gamma(\alpha/K)^K} \prod_{k=1}^K \mu_k^{\alpha/K - 1}$$

- **Properties** (symmetric version)

$$\mathbb{E}[\mu_k] = \frac{\alpha_k}{\alpha_0} = \frac{1}{K}$$

$$\text{var}[\mu_k] = \frac{\alpha_k(\alpha_0 - \alpha_k)}{\alpha_0^2(\alpha_0 + 1)} = \frac{K-1}{K^2(\alpha+1)}$$

$$\text{cov}[\mu_j, \mu_k] = -\frac{\alpha_j \alpha_k}{\alpha_0^2(\alpha_0 + 1)} = -\frac{1}{K^2(\alpha+1)}$$



## Recap: Bayesian Mixture Models

- Let's be Bayesian about mixture models

- Place priors over our parameters
- Again, introduce variable  $z_n$  as indicator which component data point  $x_n$  belongs to.

$$z_n | \pi \sim \text{Multinomial}(\pi)$$

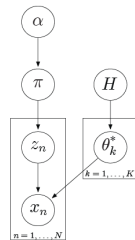
$$x_n | z_n = k, \mu, \Sigma \sim \mathcal{N}(\mu_k, \Sigma_k)$$

- Introduce **conjugate priors** over parameters

$$\pi \sim \text{Dirichlet}\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right)$$

$$\mu_k, \Sigma_k \sim H = \mathcal{N} - \text{IW}(0, s, d, \phi)$$

"Normal – Inverse Wishart"



## Recap: Bayesian Mixture Models

- **Full Bayesian Treatment**

- Given a dataset, we are interested in the cluster assignments

$$p(\mathbf{Z}|\mathbf{X}) = \frac{p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})}{\sum_{\mathbf{Z}} p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})}$$

where the likelihood is obtained by marginalizing over the parameters  $\theta$

$$p(\mathbf{X}|\mathbf{Z}) = \int p(\mathbf{X}|\mathbf{Z}, \theta)p(\theta)d\theta$$

$$= \int \prod_{n=1}^N \prod_{k=1}^K p(x_n | z_n, \theta_k)p(\theta_k|H)d\theta$$

- The posterior over assignments is intractable!

- Denominator requires summing over all possible partitions of the data into  $K$  groups!

⇒ Need efficient approximate inference methods to solve this...

## Recap: Mixture Models with Dirichlet Priors

- Integrating out the mixing proportions  $\pi$

$$p(\mathbf{z}|\alpha) = \int p(\mathbf{z}|\pi)p(\pi|\alpha)d\pi$$

$$= \frac{\Gamma(\alpha)}{\Gamma(N+\alpha)} \prod_{k=1}^K \frac{\Gamma(N_k + \alpha/K)}{\Gamma(\alpha/K)}$$

- Conditional probabilities

- Examine the conditional of  $\mathbf{z}_n$  given all other variables  $\mathbf{z}_{-n}$

$$p(z_{nk} = 1 | \mathbf{z}_{-n}, \alpha) = \frac{p(z_{nk} = 1, \mathbf{z}_{-n} | \alpha)}{p(\mathbf{z}_{-n} | \alpha)}$$

$$= \frac{N_{-n,k} + \alpha/K}{N - 1 + \alpha} \quad N_{-n,k} \stackrel{\text{def}}{=} \sum_{i=1, i \neq n}^N z_{ik}$$

⇒ The **more populous** a class is, the more likely it is to be joined!

127

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Zoubin Ghahramani

## Recap: Infinite Dirichlet Mixture Models

- Conditional probabilities: Finite  $K$

$$p(z_{nk} = 1 | \mathbf{z}_{-n}, \alpha) = \frac{N_{-n,k} + \alpha/K}{N - 1 + \alpha}, \quad N_{-n,k} \stackrel{\text{def}}{=} \sum_{i=1, i \neq n}^N z_{ik}$$

- Conditional probabilities: Infinite  $K$

- Taking the limit as  $K \rightarrow \infty$  yields the conditionals

$$p(z_{nk} = 1 | \mathbf{z}_{-n}, \alpha) = \begin{cases} \frac{N_{-n,k}}{N-1+\alpha} & \text{if } k \text{ represented} \\ \frac{\alpha}{N-1+\alpha} & \text{if all } k \text{ not represented} \end{cases}$$

- Left-over mass  $\alpha$  ⇒ countably infinite number of indicator settings

128

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Zoubin Ghahramani

## Recap: Gibbs Sampling for Finite Mixtures

- We need approximate inference here

- Gibbs Sampling: Conditionals are simple to compute

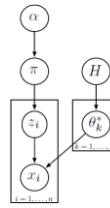
$$p(\mathbf{z}_n = k | \text{others}) \propto \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\pi | \mathbf{z} \sim \text{Dir}(N_1 + \alpha/K, \dots, N_K + \alpha/K)$$

$$\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k | \text{others} \sim \mathcal{N} - \mathcal{IW}(v', s', d', \phi')$$

- However, this will be rather inefficient...

- In each iteration, algorithm can only change the assignment for individual data points.
- There are often groups of data points that are associated with high probability to the same component. ⇒ Unlikely that group is moved.
- Better performance by **collapsed Gibbs sampling** which integrates out the parameters  $\pi, \mu, \Sigma$ .



129

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Yee Whye Teh

Image source: Yee Whye Teh

## Recap: Collapsed Finite Bayesian Mixture

- More efficient algorithm

- Conjugate priors allow analytic integration of some parameters
- Resulting sampler operates on reduced space of cluster assignments (implicitly considers all possible cluster shapes)

- Procedure

- The model implies the factorization

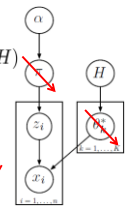
$$p(\mathbf{z}_n | \mathbf{z}_{-n}, \mathbf{x}, \alpha, H) \propto p(\mathbf{z}_n | \mathbf{z}_{-n}, \alpha) p(\mathbf{x}_n | \mathbf{z}_n, H)$$

- Derive

$$p(\mathbf{z} | \alpha) = \int p(\mathbf{z} | \pi) p(\pi | \alpha) d\pi \quad \checkmark$$

$$p(\mathbf{x}_n | \mathbf{z}_n, H) = \int \sum_{k=1}^K z_{nk} p(\mathbf{x}_n | \boldsymbol{\theta}_k) p(\boldsymbol{\theta}_k | H) d\boldsymbol{\theta} \quad \checkmark$$

⇒ Conjugate prior, Normal - Inverse Wishart



130

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Slide adapted from Erik Sudderth

Image source: Yee Whye Teh

## Course Outline

- Regression Techniques

- Linear Regression
- Regularization (Ridge, Lasso)
- Kernels (Kernel Ridge Regression)

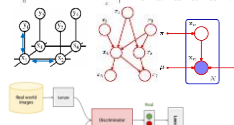
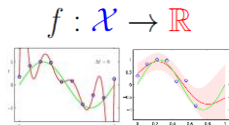
- Deep Reinforcement Learning

- Probabilistic Graphical Models

- Bayesian Networks
- Markov Random Fields
- Inference (exact & approximate)
- Latent Variable Models

- Deep Generative Models

- Generative Adversarial Networks
- Variational Autoencoders



131

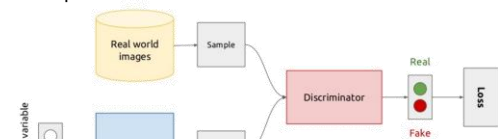
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Generative Adversarial Networks (GANs)

- Conceptual view



- Main idea

- Simultaneously train an image generator  $G$  and a discriminator  $D$ .
- Interpreted as a two-player game

132

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

Image credit: Kevin McGuinness

## Recap: GAN Loss Function

- This corresponds to a two-player minimax game:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z)))]$$

- Explanation

- Train  $D$  to maximize the probability of assigning the correct label to both training examples and samples from  $G$ .
- Simultaneously train  $G$  to minimize  $\log(1 - D(G(z)))$ .

- The Nash equilibrium of this game is achieved at

- $p_g(x) = p_{data}(x) \quad \forall x$
- $D(x) = \frac{1}{2} \quad \forall x$

133

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## GAN Algorithm

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

**for** number of training iterations **do**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Sample minibatch of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{data}(x)$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_D} \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log (1 - D(G(z^{(i)})))] \quad \text{Discriminator updates}$$

**end for**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_G} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))) \quad \text{Generator updates}$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

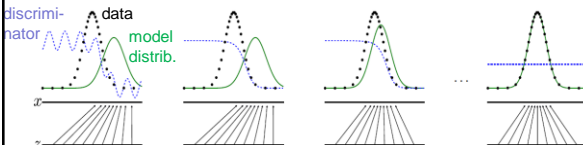
134

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Intuition behind GANs



- Behavior near convergence

- In the inner loop,  $D$  is trained to discriminate samples from data.
- Gradient of  $D$  guides  $G$  to flow to regions that are more likely to be classified as data.
- After several steps of training,  $G$  and  $D$  will reach a point at which they cannot further improve, because  $p_g = p_{data}$ .
- Now, the discriminator is unable to differentiate between the two distributions, i.e.,  $D(x) = 0.5$ .

135

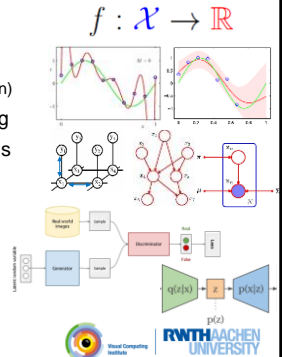
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders



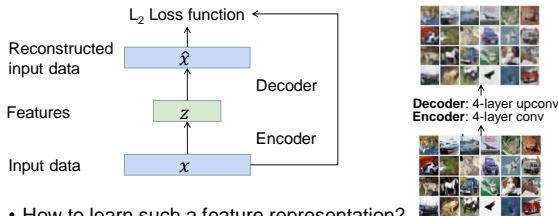
136

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



RWTH AACHEN  
UNIVERSITY

## Recap: Autoencoders



- How to learn such a feature representation?

- Unsupervised learning approach for learning a lower-dimensional feature representation  $z$  from unlabeled input data  $x$ .
- $z$  usually smaller than  $x$  (dimensionality reduction)
- Want to capture meaningful factors of variation in the data. Train such that features can be used to reconstruct original data.

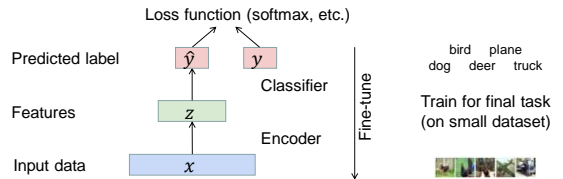
137

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Feifei Li



RWTH AACHEN  
UNIVERSITY

## Recap: Autoencoders



- After training

- Throw away the decoder part
- Encoder can be used to initialize a supervised model
- Fine-tune encoder jointly with supervised model
- Idea used in the 90s and early 2000s to pre-train deeper models

138

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
Slide credit: Feifei Li



RWTH AACHEN  
UNIVERSITY



### Recap: Variants of Autoencoders

- **Regularized Autoencoders**
  - Include a regularization term to the loss function:  $L(x, g(f(x))) + \Omega(z)$
  - E.g., enforce sparsity by an  $L_1$  regularizer  $\Omega(z) = \lambda \sum_i |z_i|$

139 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
RWTH AACHEN UNIVERSITY

### Recap: Variants of Autoencoders

- **Denoising Autoencoder (DAE)**
  - Rather than the reconstruction loss, minimize  $L(x, g(f(\tilde{x})))$  where  $\tilde{x}$  is a copy of  $x$  that has been corrupted by some noise.
  - Denoising forces  $f$  and  $g$  to implicitly learn the structure of  $p_{data}(x)$ .

140 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
RWTH AACHEN UNIVERSITY  
Image source: (Gooifellow 2016)

### Recap: Probabilistic Spin on Autoencoders

- **Idea: Sample the model to generate data**
  - We want to estimate the true parameters  $\theta^*$  of this generative model.
- **How should we represent the model?**
  - Choose prior  $p(z)$  to be simple, e.g., Gaussian
  - Conditional  $p(x|z)$  is complex (generates image)  $\Rightarrow$  Represent with neural network
  - Learn model parameters to maximize likelihood of training data

$$p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz \quad \text{Intractable!}$$

141 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
RWTH AACHEN UNIVERSITY  
Slide adapted from Feifei Li

### Recap: Variational Autoencoders

- Define additional encoder network  $q_\phi(z|x)$ 
  - Since we are modelling probabilistic generation of data, encoder and decoder networks are probabilistic

– Encoder and decoder networks are also called **recognition/inference** and **generation** networks

D. Kingma, M. Welling, *Auto-Encoding Variational Bayes*, ICLR 2014

142 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
RWTH AACHEN UNIVERSITY  
Slide credit: Feifei Li

### Recap: Variational Autoencoders

- We can now work out the log-likelihood

$$\begin{aligned} \log p_\theta(x^{(l)}) &= \mathbb{E}_{z \sim q_\phi(z|x^{(l)})} [\log p_\theta(x^{(l)})] && (p_\theta(x^{(l)}) \text{ does not depend on } z) \\ &= \mathbb{E}_z \left[ \log \frac{p_\theta(x^{(l)}|z)p_\theta(z)}{p_\theta(z|x^{(l)})} \right] && \text{(Bayes' Rule)} \\ &= \mathbb{E}_z \left[ \log \frac{p_\theta(x^{(l)}|z)p_\theta(z) q_\phi(z|x^{(l)})}{p_\theta(z|x^{(l)}) q_\phi(z|x^{(l)})} \right] && \text{(Multiply by constant)} \\ &= \mathbb{E}_z [\log p_\theta(x^{(l)}|z)] - \mathbb{E}_z \left[ \log \frac{q_\phi(z|x^{(l)})}{p_\theta(z)} \right] + \mathbb{E}_z \left[ \log \frac{q_\phi(z|x^{(l)})}{p_\theta(z|x^{(l)})} \right] \\ &= \mathbb{E}_z [\log p_\theta(x^{(l)}|z)] - D_{KL}(q_\phi(z|x^{(l)})||p_\theta(z)) + D_{KL}(q_\phi(z|x^{(l)})||p_\theta(z|x^{(l)})) \end{aligned}$$

$$\underbrace{\mathbb{E}_z [\log p_\theta(x^{(l)}|z)] - D_{KL}(q_\phi(z|x^{(l)})||p_\theta(z))}_{\mathcal{L}(x^{(l)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z|x^{(l)})||p_\theta(z|x^{(l)}))}_{\geq 0}$$

Tractable lower bound, which we can take gradient of and optimize

143 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
RWTH AACHEN UNIVERSITY  
Slide credit: Feifei Li

### Recap: Variational Autoencoders

- Variational Lower Bound (“ELBO”)
 
$$\begin{aligned} \log p_\theta(x^{(l)}) &\geq \mathcal{L}(x^{(l)}, \theta, \phi) \\ &= \mathbb{E}_z [\log p_\theta(x^{(l)}|z)] - D_{KL}(q_\phi(z|x^{(l)})||p_\theta(z)) \end{aligned}$$
  - “Reconstruct the input data”
  - “Make approximate posterior distribution close to prior”
- Training: Maximize lower bound
 
$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

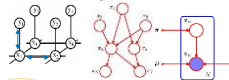
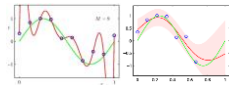
144 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition  
RWTH AACHEN UNIVERSITY  
Slide adapted from Feifei Li



## We're Done!

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- Deep Generative Models
  - Generative Adversarial Networks
  - Variational Autoencoders

$$f: \mathcal{X} \rightarrow \mathbb{R}$$



145

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition



## Any More Questions?

*Good luck for the exam!*

146

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Advanced Machine Learning  
Part 20 – Repetition

