

RWTH AACHEN
UNIVERSITY

Machine Learning – Lecture 22

Repetition

29.01.2018

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

Machine Learning Winter '17

RWTH AACHEN
UNIVERSITY

Announcements

- Exams
 - Special oral exams (for exchange students):
 - We're in the process of sending out the exam slots
 - You'll receive an email with details tonight
 - Format: 30 minutes, 4 questions, 3 answers
 - Regular exams:
 - We will send out an email with the assignment to lecture halls
 - Format: 120min, closed-book exam

B. Leibe 2

RWTH AACHEN
UNIVERSITY

Announcements (2)

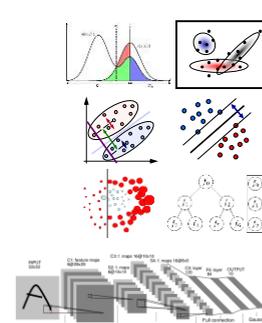
- Today, I'll summarize the most important points from the lecture.
 - It is an opportunity for you to ask questions...
 - ...or get additional explanations about certain topics.
 - So, please do ask.
- Today's slides are intended as an index for the lecture.
 - But they are not complete, won't be sufficient as only tool.
 - Also look at the exercises – they often explain algorithms in detail.

B. Leibe 3

RWTH AACHEN
UNIVERSITY

Course Outline

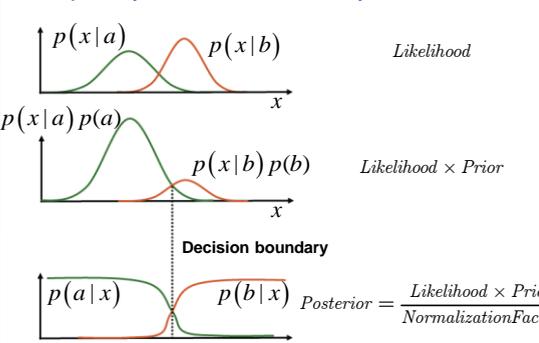
- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
 - Mixture Models and EM
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks



B. Leibe 4

RWTH AACHEN
UNIVERSITY

Recap: Bayes Decision Theory



5

Slide credit: Bernt Schiele B. Leibe Image source: G.M. Bishop, 2006

RWTH AACHEN
UNIVERSITY

Recap: Bayes Decision Theory

- Optimal decision rule
 - Decide for C_1 if

$$p(C_1|x) > p(C_2|x)$$
 - This is equivalent to

$$p(x|C_1)p(C_1) > p(x|C_2)p(C_2)$$
 - Which is again equivalent to (Likelihood-Ratio test)

$$\frac{p(x|C_1)}{p(x|C_2)} > \underbrace{\frac{p(C_2)}{p(C_1)}}_{\text{Decision threshold } \theta}$$

6

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Bayes Decision Theory

- Decision regions: $\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}_3, \dots$

$R1$ $R2$ $R3$

Machine Learning Winter '17 B. Leibe 7

Slide credit: Bernt Schiele

RWTH AACHEN UNIVERSITY

Recap: Classifying with Loss Functions

- In general, we can formalize this by introducing a loss matrix L_{kj}

$$L_{kj} = \text{loss for decision } \mathcal{C}_j \text{ if truth is } \mathcal{C}_k.$$

- Example: cancer diagnosis

		Decision	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

$$L_{\text{cancer diagnosis}} = \begin{matrix} & \text{cancer} & \text{normal} \\ \text{cancer} & 0 & 1000 \\ \text{normal} & 1 & 0 \end{matrix}$$

Machine Learning Winter '17 B. Leibe 8

RWTH AACHEN UNIVERSITY

Recap: Minimizing the Expected Loss

- Optimal solution minimizes the loss.
 - But: loss function depends on the true class, which is unknown.
- Solution: Minimize the expected loss

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) dx$$

- This can be done by choosing the regions \mathcal{R}_j such that

$$\mathbb{E}[L] = \sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

which is easy to do once we know the posterior class probabilities $p(\mathcal{C}_k | \mathbf{x})$

see Exercise 1.3

Machine Learning Winter '17 B. Leibe 9

RWTH AACHEN UNIVERSITY

Recap: The Reject Option

- Classification errors arise from regions where the largest posterior probability $p(\mathcal{C}_k | \mathbf{x})$ is significantly less than 1.
 - These are the regions where we are relatively uncertain about class membership.
 - For some applications, it may be better to reject the automatic decision entirely in such a case and e.g. consult a human expert.

Machine Learning Winter '17 B. Leibe 10
Image source: C.M. Bishop, 2006

RWTH AACHEN UNIVERSITY

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
 - Mixture Models and EM
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

Machine Learning Winter '17 B. Leibe 11

RWTH AACHEN UNIVERSITY

Recap: Gaussian (or Normal) Distribution

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

Machine Learning Winter '17 B. Leibe 12
Image source: C.M. Bishop, 2006

RWTH AACHEN UNIVERSITY

Recap: Maximum Likelihood Approach

- Computation of the likelihood
 - Single data point: $p(x_n|\theta)$
 - Assumption: all data points $X = \{x_1, \dots, x_n\}$ independent

$$L(\theta) = p(X|\theta) = \prod_{n=1}^N p(x_n|\theta)$$
 - Log-likelihood

$$E(\theta) = -\ln L(\theta) = -\sum_{n=1}^N \ln p(x_n|\theta)$$
- Estimation of the parameters θ (Learning)
 - Maximize the likelihood (= minimize the negative log-likelihood)
 - Take the derivative and set it to zero.

$$\frac{\partial}{\partial \theta} E(\theta) = -\sum_{n=1}^N \frac{\partial}{\partial \theta} \ln p(x_n|\theta) \stackrel{!}{=} 0$$

see Exercise 1.4

Machine Learning Winter '17

13

Slide credit: Bernt Schiele, B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Bayesian Learning Approach

- Bayesian view:
 - Consider the parameter vector θ as a random variable.
 - When estimating the parameters, what we compute is

$$p(x|X) = \int p(x, \theta|X) d\theta$$

Assumption: given θ , this doesn't depend on X anymore

$$p(x, \theta|X) = p(x|\theta, X)p(\theta|X)$$
 - $$p(x|X) = \int \underbrace{p(x|\theta)p(\theta|X)}_{\text{This is entirely determined by the parameter } \theta \text{ (i.e. by the parametric form of the pdf).}} d\theta$$

Machine Learning Winter '17

14

Slide adapted from Bernt Schiele, B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Bayesian Learning Approach

- Discussion
 - Likelihood of the parametric form θ given the data set X .
 - Prior for the parameters θ
 - Estimate for x based on parametric form θ
$$p(x|X) = \int \frac{p(x|\theta)L(\theta)p(\theta)}{\int L(\theta)p(\theta)d\theta} d\theta$$

Normalization: integrate over all possible values of θ
- The more uncertain we are about θ , the more we average over all possible parameter values.

Machine Learning Winter '17

15

B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Histograms

- Basic idea:
 - Partition the data space into distinct bins with widths Δ_i and count the number of observations, n_i , in each bin.

$$p_i = \frac{n_i}{N\Delta_i}$$
 - Often, the same width is used for all bins, $\Delta_i = \Delta$.
 - This can be done, in principle, for any dimensionality D ...

Machine Learning Winter '17

16

Image source: G.M. Bishop, 2006

RWTH AACHEN UNIVERSITY

Recap: Kernel Density Estimation

- Approximation formula:

$$p(\mathbf{x}) \approx \frac{K}{NV}$$
 - fixed V determine K → Kernel Methods
 - fixed K determine V → K-Nearest Neighbor
- Kernel methods
 - Place a kernel window k at location \mathbf{x} and count how many data points fall inside it.
- K-Nearest Neighbor
 - Increase the volume V until the K next data points are found.

see Exercise 1.5

Machine Learning Winter '17

17

Slide adapted from Bernt Schiele, B. Leibe

RWTH AACHEN UNIVERSITY

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
 - Mixture Models and EM
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

Machine Learning Winter '17

18

B. Leibe

Machine Learning Winter '17

Recap: Mixture of Gaussians (MoG)

RWTH AACHEN UNIVERSITY

- “Generative model”

$p(j) = \pi_j$ “Weight” of mixture component
 $p(x|\theta_j)$ Mixture component
 $p(x|\theta) = \sum_{j=1}^M p(x|\theta_j)p(j)$ Mixture density

Slide credit: Bernt Schiele B. Leibe 19

Machine Learning Winter '17

Recap: MoG – Iterative Strategy

RWTH AACHEN UNIVERSITY

- Assuming we knew the values of the hidden variable...

assumed known \rightarrow 1 111 22 2 2 j
 $h(j=1|x_n) =$ 1 111 00 0 0
 $h(j=2|x_n) =$ 0 000 11 1 1
 $\mu_1 = \frac{\sum_{n=1}^N h(j=1|x_n)x_n}{\sum_{i=1}^N h(j=1|x_n)}$ $\mu_2 = \frac{\sum_{n=1}^N h(j=2|x_n)x_n}{\sum_{i=1}^N h(j=2|x_n)}$

Slide credit: Bernt Schiele B. Leibe 20

Machine Learning Winter '17

Recap: MoG – Iterative Strategy

RWTH AACHEN UNIVERSITY

- Assuming we knew the mixture components...

$p(j=1|x)$ $p(j=2|x)$
 1 111 22 2 2 j

- Bayes decision rule: Decide $j=1$ if $p(j=1|x_n) > p(j=2|x_n)$

Slide credit: Bernt Schiele B. Leibe 21

Machine Learning Winter '17

Recap: K-Means Clustering

RWTH AACHEN UNIVERSITY

- Iterative procedure
 - Initialization: pick K arbitrary centroids (cluster means)
 - Assign each sample to the closest centroid.
 - Adjust the centroids to be the means of the samples assigned to them.
 - Go to step 2 (until no change)
- Algorithm is guaranteed to converge after finite #iterations.
 - Local optimum
 - Final result depends on initialization.

Slide credit: Bernt Schiele B. Leibe 22

Machine Learning Winter '17

Recap: EM Algorithm

RWTH AACHEN UNIVERSITY

- Expectation-Maximization (EM) Algorithm
 - E-Step: softly assign samples to mixture components

$$\gamma_j(x_n) \leftarrow \frac{\pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)} \quad \forall j = 1, \dots, K, \quad n = 1, \dots, N$$
 - M-Step: re-estimate the parameters (separately for each mixture component) based on the soft assignments

$$\hat{N}_j \leftarrow \sum_{n=1}^N \gamma_j(x_n) = \text{soft number of samples labeled } j$$

$$\hat{\pi}_j^{\text{new}} \leftarrow \frac{\hat{N}_j}{N}$$

$$\hat{\mu}_j^{\text{new}} \leftarrow \frac{1}{\hat{N}_j} \sum_{n=1}^N \gamma_j(x_n) x_n$$

$$\hat{\Sigma}_j^{\text{new}} \leftarrow \frac{1}{\hat{N}_j} \sum_{n=1}^N \gamma_j(x_n) (x_n - \hat{\mu}_j^{\text{new}})(x_n - \hat{\mu}_j^{\text{new}})^T$$

see Exercise 1.6

Slide adapted from Bernt Schiele B. Leibe 23

Machine Learning Winter '17

Course Outline

RWTH AACHEN UNIVERSITY

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

Slide credit: Bernt Schiele B. Leibe 24

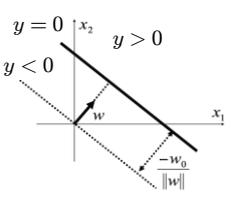
RWTH AACHEN UNIVERSITY

Recap: Linear Discriminant Functions

- Basic idea
 - Directly encode decision boundary
 - Minimize misclassification probability directly.
- Linear discriminant functions
 - $$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

\mathbf{w}
weight vector

w_0
"bias"
(= threshold)
 - \mathbf{w}, w_0 define a hyperplane in \mathbb{R}^D .
 - If a data set can be perfectly classified by a linear discriminant, then we call it **linearly separable**.



Machine Learning Winter '17 25

Slide adapted from Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Least-Squares Classification

- Simplest approach
 - Directly try to **minimize the sum-of-squares error**

$$E(\mathbf{w}) = \sum_{n=1}^N (y(\mathbf{x}_n; \mathbf{w}) - t_n)^2$$

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{\mathbf{X}} \tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}} \tilde{\mathbf{W}} - \mathbf{T}) \right\}$$

- Setting the derivative to zero yields
 - $$\tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T}$$
- We then obtain the discriminant function as
 - $$y(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T (\tilde{\mathbf{X}}^\dagger)^T \tilde{\mathbf{x}}$$
- \Rightarrow **Exact, closed-form solution** for the discriminant function parameters.

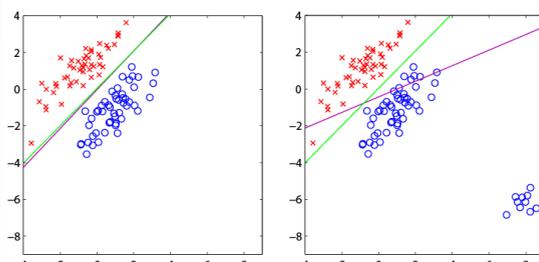


Machine Learning Winter '17 26

B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Problems with Least Squares



- Least-squares is very sensitive to outliers!
 - The error function penalizes predictions that are "too correct".

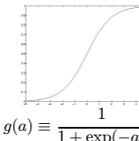
Machine Learning Winter '17 27

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

RWTH AACHEN UNIVERSITY

Recap: Generalized Linear Models

- Generalized linear model
 - $$y(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + w_0)$$
 - $g(\cdot)$ is called an **activation function** and may be nonlinear.
 - The decision surfaces correspond to
 - $$y(\mathbf{x}) = \text{const.} \Leftrightarrow \mathbf{w}^T \mathbf{x} + w_0 = \text{const.}$$
 - If g is monotonous (which is typically the case), the resulting decision boundaries are still linear functions of \mathbf{x} .
- Advantages of the non-linearity
 - Can be used to bound the influence of outliers and "too correct" data points.
 - When using a sigmoid for $g(\cdot)$, we can interpret the $y(\mathbf{x})$ as posterior probabilities.



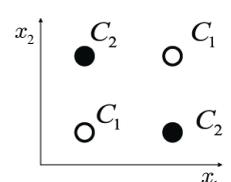
Machine Learning Winter '17 28

B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Linear Separability

- Up to now: restrictive assumption
 - Only consider linear decision boundaries
- Classical counterexample: XOR



Machine Learning Winter '17 29

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Extension to Nonlinear Basis Fcts.

- Generalization
 - Transform vector \mathbf{x} with M nonlinear basis functions $\phi_j(\mathbf{x})$:
 - $$y_k(\mathbf{x}) = \sum_{j=1}^M w_{kj} \phi_j(\mathbf{x}) + w_{k0}$$
- Advantages
 - Transformation allows non-linear decision boundaries.
 - By choosing the right ϕ_j , every continuous function can (in principle) be approximated with arbitrary accuracy.
- Disadvantage
 - The error function can in general no longer be minimized in closed form.
 - \Rightarrow Minimization with Gradient Descent

Machine Learning Winter '17 30

B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Probabilistic Discriminative Models

- Consider models of the form

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$
 with

$$p(\mathcal{C}_2|\phi) = 1 - p(\mathcal{C}_1|\phi)$$
- This model is called **logistic regression**.
- Properties
 - Probabilistic interpretation
 - But discriminative method: only focus on decision hyperplane
 - Advantageous for high-dimensional spaces, requires less parameters than explicitly modeling $p(\phi|\mathcal{C}_k)$ and $p(\mathcal{C}_k)$.

31

RWTH AACHEN UNIVERSITY

Recap: Logistic Regression

- Let's consider a data set $\{\phi_n, t_n\}$ with $n = 1, \dots, N$, where $\phi_n = \phi(\mathbf{x}_n)$ and $t_n \in \{0, 1\}$ $\mathbf{t} = (t_1, \dots, t_N)^T$
- With $y_n = p(\mathcal{C}_1|\phi_n)$, we can write the likelihood as

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$
- Define the error function as the negative log-likelihood

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w})$$

$$= -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$
 - This is the so-called **cross-entropy error function**.

32

RWTH AACHEN UNIVERSITY

Recap: Iterative Methods for Estimation

- Gradient Descent (1st order)

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w})|_{\mathbf{w}^{(\tau)}}$$
 - Simple and general
 - Relatively slow to converge, has problems with some functions
- Newton-Raphson (2nd order)

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \mathbf{H}^{-1} \nabla E(\mathbf{w})|_{\mathbf{w}^{(\tau)}}$$
 where $\mathbf{H} = \nabla^2 E(\mathbf{w})$; the Hessian matrix, i.e. the matrix of second derivatives.
 - Local quadratic approximation to the target function
 - Faster convergence

33

RWTH AACHEN UNIVERSITY

Recap: Iteratively Reweighted Least Squares

- Update equations

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t})$$

$$= (\Phi^T \mathbf{R} \Phi)^{-1} \left\{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(\tau)} - \Phi^T (\mathbf{y} - \mathbf{t}) \right\}$$

$$= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}$$
 with $\mathbf{z} = \Phi \mathbf{w}^{(\tau)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})$
- Very similar form to pseudo-inverse (normal equations)
 - But now with non-constant weighing matrix \mathbf{R} (depends on \mathbf{w})
 - Need to apply normal equations iteratively.
 - ⇒ **Iteratively Reweighted Least-Squares (IRLS)**

34

RWTH AACHEN UNIVERSITY

Recap: Softmax Regression

- Multi-class generalization of logistic regression
 - In logistic regression, we assumed binary labels $t_n \in \{0, 1\}$
 - Softmax generalizes this to K values in 1-of- K notation.

$$\mathbf{y}(\mathbf{x}; \mathbf{w}) = \begin{bmatrix} P(y=1|\mathbf{x}; \mathbf{w}) \\ P(y=2|\mathbf{x}; \mathbf{w}) \\ \vdots \\ P(y=K|\mathbf{x}; \mathbf{w}) \end{bmatrix} = \frac{1}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \mathbf{x})} \begin{bmatrix} \exp(\mathbf{w}_1^T \mathbf{x}) \\ \exp(\mathbf{w}_2^T \mathbf{x}) \\ \vdots \\ \exp(\mathbf{w}_K^T \mathbf{x}) \end{bmatrix}$$

- This uses the **softmax** function

$$\frac{\exp(a_k)}{\sum_j \exp(a_j)}$$
- Note: the resulting distribution is normalized.

35

RWTH AACHEN UNIVERSITY

Recap: Softmax Regression Cost Function

- Logistic regression
 - Alternative way of writing the cost function

$$E(\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

$$= -\sum_{n=1}^N \sum_{k=0}^1 \{\mathbb{I}(t_n = k) \ln P(y_n = k|\mathbf{x}_n; \mathbf{w})\}$$
- Softmax regression
 - Generalization to K classes using indicator functions.

$$E(\mathbf{w}) = -\sum_{n=1}^N \sum_{k=1}^K \left\{ \mathbb{I}(t_n = k) \ln \frac{\exp(\mathbf{w}_k^T \mathbf{x})}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \mathbf{x})} \right\}$$
 - $$\nabla_{\mathbf{w}_k} E(\mathbf{w}) = -\sum_{n=1}^N \{\mathbb{I}(t_n = k) \ln P(y_n = k|\mathbf{x}_n; \mathbf{w})\}$$

36

Machine Learning Winter '17

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

B. Leibe 37

Machine Learning Winter '17

Recap: Generalization and Overfitting

- Goal: predict class labels of new observations
 - Train classification model on limited training set.
 - The further we optimize the model parameters, the more the **training error** will decrease.
 - However, at some point the **test error** will go up again.
 - ⇒ *Overfitting to the training set!*

B. Leibe 38
Image source: B. Schölkopf

Machine Learning Winter '17

Recap: Support Vector Machine (SVM)

- Basic idea
 - The SVM tries to find a classifier which maximizes the **margin** between pos. and neg. data points.
 - Up to now: consider linear classifiers

$$\mathbf{w}^T \mathbf{x} + b = 0$$
- Formulation as a **convex optimization problem**
 - Find the hyperplane satisfying

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$
 under the constraints

$$t_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \forall n$$
 based on training data points \mathbf{x}_n and target values $t_n \in \{-1, 1\}$

B. Leibe 39

Machine Learning Winter '17

Recap: SVM – Primal Formulation

- Lagrangian primal form

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^T \mathbf{x}_n + b) - 1\}$$

$$= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1\}$$
- The solution of L_p needs to fulfill the **KKT conditions**
 - Necessary and sufficient conditions

KKT:
$\lambda \geq 0$
$f(\mathbf{x}) \geq 0$
$\lambda f(\mathbf{x}) = 0$

B. Leibe 40

Machine Learning Winter '17

Recap: SVM – Solution

- Solution for the hyperplane
 - Computed as a linear combination of the training examples

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \mathbf{x}_n$$
 - Sparse solution: $a_n \neq 0$ only for some points, the support vectors
 - ⇒ Only the SVs actually influence the decision boundary!
 - Compute b by averaging over all support vectors:

$$b = \frac{1}{N_S} \sum_{n \in S} \left(t_n - \sum_{m \in S} a_m t_m \mathbf{x}_m^T \mathbf{x}_n \right)$$

B. Leibe 41

Machine Learning Winter '17

Recap: SVM – Support Vectors

- The training points for which $a_n > 0$ are called "**support vectors**".
- Graphical interpretation:
 - The support vectors are the points on the margin.
 - They *define* the margin and thus the hyperplane.

⇒ All other data points can be discarded!

Slide adapted from Bernt Schölkopf B. Leibe 42
Image source: C. Burges, 1998

RWTH AACHEN UNIVERSITY

Recap: SVM – Dual Formulation

- Maximize

$$L_d(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m (\mathbf{x}_m^T \mathbf{x}_n)$$

see Exercise 2.2

under the conditions

$$a_n \geq 0 \quad \forall n$$

$$\sum_{n=1}^N a_n t_n = 0$$

- Comparison
 - L_d is equivalent to the primal form L_p , but only depends on a_n .
 - L_p scales with $\mathcal{O}(D^3)$.
 - L_d scales with $\mathcal{O}(N^3)$ – in practice between $\mathcal{O}(N)$ and $\mathcal{O}(N^2)$.

Machine Learning Winter '17
B. Leibe
43
Slide adapted from Bernt Schiele

RWTH AACHEN UNIVERSITY

Recap: SVM for Non-Separable Data

- Slack variables
 - One slack variable $\xi_n \geq 0$ for each training data point.
- Interpretation
 - $\xi_n = 0$ for points that are on the correct side of the margin.
 - $\xi_n = |t_n - y(\mathbf{x}_n)|$ for all other points.

Point on decision boundary: $\xi_n = 1$

Misclassified point: $\xi_n > 1$

- We do not have to set the slack variables ourselves!
- They are jointly optimized together with \mathbf{w} .

Machine Learning Winter '17
B. Leibe
44

RWTH AACHEN UNIVERSITY

Recap: SVM – New Dual Formulation

- New SVM Dual: Maximize

$$L_d(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m (\mathbf{x}_m^T \mathbf{x}_n)$$

see Exercise 2.2

under the conditions

$$0 \cdot a_n \cdot C$$

This is all that changed!

$$\sum_{n=1}^N a_n t_n = 0$$

- This is again a quadratic programming problem
 - \Rightarrow Solve as before...

Machine Learning Winter '17
B. Leibe
45
Slide adapted from Bernt Schiele

RWTH AACHEN UNIVERSITY

Recap: Nonlinear SVMs

- General idea: The original input space can be mapped to some higher-dimensional feature space where the training set is separable:

Machine Learning Winter '17
Slide credit: Raymond Mooney
46

RWTH AACHEN UNIVERSITY

Recap: The Kernel Trick

- Important observation
 - $\phi(\mathbf{x})$ only appears in the form of dot products $\phi(\mathbf{x})^T \phi(\mathbf{y})$:

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

$$= \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) + b$$
 - Define a so-called **kernel function** $k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$.
 - Now, in place of the dot product, use the kernel instead:

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}_n, \mathbf{x}) + b$$
 - The kernel function *implicitly* maps the data to the higher-dimensional space (without having to compute $\phi(\mathbf{x})$ explicitly)!

Machine Learning Winter '17
B. Leibe
47

RWTH AACHEN UNIVERSITY

Recap: Kernels Fulfilling Mercer's Condition

- Polynomial kernel

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^p$$
- Radial Basis Function kernel

$$k(\mathbf{x}, \mathbf{y}) = \exp \left\{ -\frac{(\mathbf{x} - \mathbf{y})^2}{2\sigma^2} \right\}$$
 e.g. Gaussian
- Hyperbolic tangent kernel

$$k(\mathbf{x}, \mathbf{y}) = \tanh(\kappa \mathbf{x}^T \mathbf{y} + \delta)$$
 e.g. Sigmoid
- And many, many more, including kernels on graphs, strings, and symbolic data...

Machine Learning Winter '17
Slide credit: Bernt Schiele
B. Leibe
48

RWTH AACHEN UNIVERSITY

Recap: Kernels Fulfilling Mercer's Condition

- Polynomial kernel

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^p$$
- Radial Basis Function kernel

$$k(\mathbf{x}, \mathbf{y}) = \exp \left\{ -\frac{(\mathbf{x} - \mathbf{y})^2}{2\sigma^2} \right\}$$
 e.g. Gaussian
- Hyperbolic tangent kernel

$$k(\mathbf{x}, \mathbf{y}) = \tanh(\alpha \mathbf{x}^T \mathbf{y} + \delta)$$
 e.g. Sigmoid

Actually, that was wrong in the original SVM paper...

And many, many more, including kernels on graphs, strings, and symbolic data...

Machine Learning Winter '17 | Slide credit: Bernt Schiele | B. Leibe | 49

RWTH AACHEN UNIVERSITY

Recap: Nonlinear SVM – Dual Formulation

see Exercise 2.2

- SVM Dual: Maximize

$$L_d(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_m, \mathbf{x}_n)$$
- under the conditions

$$0 \leq a_n \leq C$$

$$\sum_{n=1}^N a_n t_n = 0$$
- Classify new data points using

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}_n, \mathbf{x}) + b$$

Machine Learning Winter '17 | B. Leibe | 50

RWTH AACHEN UNIVERSITY

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

Machine Learning Winter '17 | B. Leibe | 51

RWTH AACHEN UNIVERSITY

Recap: Classifier Combination

- We've seen already a variety of different classifiers
 - k-NN
 - Bayes classifiers
 - Fisher's Linear Discriminant
 - SVMs
- Each of them has their strengths and weaknesses...
 - Can we improve performance by combining them?

Machine Learning Winter '17 | B. Leibe | 52

RWTH AACHEN UNIVERSITY

Recap: Bayesian Model Averaging

- Model Averaging
 - Suppose we have H different models $h = 1, \dots, H$ with prior probabilities $p(h)$.
 - Construct the marginal distribution over the data set

$$p(\mathbf{X}) = \sum_{h=1}^H p(\mathbf{X}|h)p(h)$$
- Average error of committee

$$\mathbb{E}_{COM} = \frac{1}{M} \mathbb{E}_{AV}$$
 - This suggests that the average error of a model can be reduced by a factor of M simply by averaging M versions of the model!
 - Unfortunately, this assumes that the errors are all **uncorrelated**. In practice, they will typically be highly correlated.

Machine Learning Winter '17 | B. Leibe | 53

RWTH AACHEN UNIVERSITY

Recap: AdaBoost – “Adaptive Boosting”

- Main idea [Freund & Schapire, 1996]
 - Instead of resampling, reweight misclassified training examples.
 - Increase the chance of being selected in a sampled training set.
 - Or increase the misclassification cost when training on the full set.
- Components
 - $h_m(\mathbf{x})$: “weak” or base classifier
 - Condition: <50% training error over any distribution
 - $H(\mathbf{x})$: “strong” or final classifier
- AdaBoost:
 - Construct a strong classifier as a thresholded linear combination of the weighted weak classifiers:

$$H(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m h_m(\mathbf{x}) \right)$$

Machine Learning Winter '17 | B. Leibe | 54

Machine Learning Winter '17

Recap: AdaBoost – Intuition

Consider a 2D feature space with **positive** and **negative** examples.

Each weak classifier splits the training examples with at least 50% accuracy.

Examples misclassified by a previous weak learner are given more emphasis at future rounds.

Slide credit: Kristen Grauman B. Leibe Figure adapted from Freund & Schapire

Machine Learning Winter '17

Recap: AdaBoost – Intuition

Weights Increased

Weak Classifier 2

Slide credit: Kristen Grauman B. Leibe Figure adapted from Freund & Schapire

Machine Learning Winter '17

Recap: AdaBoost – Intuition

Weights Increased

Weak Classifier 2

Weak classifier 3

Final classifier is combination of the weak classifiers

Slide credit: Kristen Grauman B. Leibe Figure adapted from Freund & Schapire

Machine Learning Winter '17

Recap: AdaBoost – Algorithm

1. Initialization: Set $w_n^{(1)} = \frac{1}{N}$ for $n = 1, \dots, N$.
2. For $m = 1, \dots, M$ iterations
 - a) Train a new weak classifier $h_m(\mathbf{x})$ using the current weighting coefficients $\mathbf{W}^{(m)}$ by minimizing the weighted error function

$$J_m = \sum_{n=1}^N w_n^{(m)} I(h_m(\mathbf{x}) \neq t_n)$$

$$I(A) = \begin{cases} 1, & \text{if } A \text{ is true} \\ 0, & \text{else} \end{cases}$$
 - b) Estimate the weighted error of this classifier on \mathbf{X} :

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(h_m(\mathbf{x}) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}$$
 - c) Calculate a weighting coefficient for $h_m(\mathbf{x})$:

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}$$
 - d) Update the weighting coefficients:

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(h_m(\mathbf{x}_n) \neq t_n) \}$$

see Exercise 3.1

Slide credit: Kristen Grauman B. Leibe

Machine Learning Winter '17

Recap: Comparing Error Functions

- Ideal misclassification error function
- "Hinge error" used in SVMs
- Exponential error function
 - Continuous approximation to ideal misclassification function.
 - Sequential minimization leads to simple AdaBoost scheme.
 - Disadvantage: exponential penalty for large negative values!
 - ⇒ Less robust to outliers or misclassified data points!

Image source: Bishop, 2006

Machine Learning Winter '17

Recap: Comparing Error Functions

- Ideal misclassification error function
- "Hinge error" used in SVMs
- Exponential error function
- "Cross-entropy error"

$$E = - \sum \{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \}$$
 - Similar to exponential error for $z > 0$.
 - Only grows linearly with large negative values of z .
 - ⇒ Make AdaBoost more robust by switching ⇒ "GentleBoost"

Image source: Bishop, 2006

Machine Learning Winter '17

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

B. Leibe 61

Machine Learning Winter '17

Recap: Decision Trees

```

graph TD
    Outlook[Outlook] -- Sunny --> Humidity[Humidity]
    Outlook -- Overcast --> Yes[Yes]
    Outlook -- Rain --> Wind[Wind]
    Humidity -- High --> No[No]
    Humidity -- Normal --> Yes[Yes]
    Wind -- Strong --> No[No]
    Wind -- Weak --> Yes[Yes]
  
```

- Example:
 - “Classify Saturday mornings according to whether they’re suitable for playing tennis.”

B. Leibe 62
Image source: T. Mitchell, 1997

Machine Learning Winter '17

Recap: CART Framework

- Six general questions
 1. Binary or multi-valued problem?
 - I.e. how many splits should there be at each node?
 2. Which property should be tested at a node?
 - I.e. how to select the query attribute?
 3. When should a node be declared a leaf?
 - I.e. when to stop growing the tree?
 4. How can a grown tree be simplified or pruned?
 - Goal: reduce overfitting.
 5. How to deal with impure nodes?
 - I.e. when the data itself is ambiguous.
 6. How should missing attributes be handled?

B. Leibe 63

Machine Learning Winter '17

Recap: Picking a Good Splitting Feature

see Exercise 3.2

- Goal
 - Select the query (=split) that decreases impurity the most
$$\Delta i(N) = i(N) - P_L i(N_L) - (1 - P_L) i(N_R)$$
- Impurity measures
 - Entropy impurity (information gain):

$$i(N) = - \sum_j p(C_j|N) \log_2 p(C_j|N)$$
 - Gini impurity:

$$i(N) = \sum_{i \neq j} p(C_i|N) p(C_j|N) = \frac{1}{2} \left[1 - \sum_j p^2(C_j|N) \right]$$

B. Leibe 64
Image source: R.O. Duda, P.F. Hart, D.G. Stork, 2001

Machine Learning Winter '17

Recap: Computational Complexity

- Given
 - Data points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
 - Dimensionality D
- Complexity
 - Storage: $O(N)$
 - Test runtime: $O(\log N)$
 - Training runtime: $O(DN^2 \log N)$
 - Most expensive part.
 - Critical step: selecting the optimal splitting point.
 - Need to check D dimensions, for each need to sort N data points.
$$O(DN \log N)$$

B. Leibe 65

Machine Learning Winter '17

Recap: Randomized Decision Trees

- Decision trees: main effort on finding good split
 - Training runtime: $O(DN^2 \log N)$
 - This is what takes most effort in practice.
 - Especially cumbersome with many attributes (large D).
- Idea: randomize attribute selection
 - No longer look for globally optimal split.
 - Instead randomly use subset of K attributes on which to base the split.
 - Choose best splitting attribute e.g. by maximizing the information gain (= reducing entropy):
$$\Delta E = \sum_{k=1}^K \frac{|S_k|}{|S|} \sum_{j=1}^N p_j \log_2(p_j)$$

B. Leibe 66

RWTH AACHEN UNIVERSITY

Recap: Ensemble Combination

- Ensemble combination
 - Tree leaves (l, η) store posterior probabilities of the target classes.

$$p_{l, \eta}(\mathcal{C} | \mathbf{x})$$
 - Combine the output of several trees by averaging their posteriors (Bayesian model combination)

$$p(\mathcal{C} | \mathbf{x}) = \frac{1}{L} \sum_{l=1}^L p_{l, \eta}(\mathcal{C} | \mathbf{x})$$

67

RWTH AACHEN UNIVERSITY

Recap: Random Forests (Breiman 2001)

- General ensemble method
 - Idea: Create ensemble of many (50 - 1,000) trees.
- Empirically very good results
 - Often as good as SVMs (and sometimes better)!
 - Often as good as Boosting (and sometimes better)!
- Injecting randomness
 - Bootstrap sampling process
 - On average only 63% of training examples used for building the tree
 - Remaining 37% out-of-bag samples used for validation.
 - Random attribute selection
 - Randomly choose subset of K attributes to select from at each node.
 - Faster training procedure.
- Simple majority vote for tree combination

68

RWTH AACHEN UNIVERSITY

Recap: A Graphical Interpretation

Different trees induce different partitions on the data.

By combining them, we obtain a finer subdivision of the feature space...

69

RWTH AACHEN UNIVERSITY

Recap: A Graphical Interpretation

Different trees induce different partitions on the data.

By combining them, we obtain a finer subdivision of the feature space... ..which at the same time also better reflects the uncertainty due to the bootstrapped sampling.

70

RWTH AACHEN UNIVERSITY

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

71

RWTH AACHEN UNIVERSITY

Recap: Perceptrons

- One output node per class

Output layer

Weights

Input layer
- Outputs
 - Linear outputs
$$y_k(\mathbf{x}) = \sum_{i=0}^d W_{ki} x_i$$

With output nonlinearity

$$y_k(\mathbf{x}) = g \left(\sum_{i=0}^d W_{ki} x_i \right)$$

⇒ Can be used to do multidimensional linear regression or multiclass classification.

72

RWTH AACHEN UNIVERSITY

Recap: Non-Linear Basis Functions

- Straightforward generalization

Output layer

Weights

Feature layer

Mapping (fixed)

Input layer

Output layer

Weights

Feature layer

Mapping (fixed)

Input layer

- Outputs
 - Linear outputs
$$y_k(\mathbf{x}) = \sum_{i=0}^d W_{ki} \phi(x_i)$$
- with output nonlinearity

$$y_k(\mathbf{x}) = g\left(\sum_{i=0}^d W_{ki} \phi(x_i)\right)$$

Machine Learning Winter '17 73

RWTH AACHEN UNIVERSITY

Recap: Non-Linear Basis Functions

- Straightforward generalization

Output layer

Weights

Feature layer

Mapping (fixed)

Input layer

Output layer

Weights

Feature layer

Mapping (fixed)

Input layer

- Remarks
 - Perceptrons are generalized linear discriminants!
 - Everything we know about the latter can also be applied here.
 - Note: feature functions $\phi(\mathbf{x})$ are kept fixed, not learned!

Machine Learning Winter '17 74

RWTH AACHEN UNIVERSITY

Recap: Perceptron Learning

- Process the training cases in some permutation
 - If the output unit is correct, leave the weights alone.
 - If the output unit incorrectly outputs a zero, add the input vector to the weight vector.
 - If the output unit incorrectly outputs a one, subtract the input vector from the weight vector.
- Translation

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta (y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn}) \phi_j(\mathbf{x}_n)$$
 - This is the **Delta rule** a.k.a. LMS rule!
 - ⇒ Perceptron Learning corresponds to 1st-order (stochastic) Gradient Descent of a quadratic error function!

Machine Learning Winter '17 75

Slide adapted from Geoff Hinton B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Loss Functions

- We can now also apply other loss functions
 - L_2 loss ⇒ Least-squares regression

$$L(t, y(\mathbf{x})) = \sum_n (y(\mathbf{x}_n) - t_n)^2$$
 - L_1 loss: ⇒ Median regression

$$L(t, y(\mathbf{x})) = \sum_n |y(\mathbf{x}_n) - t_n|$$
 - Cross-entropy loss ⇒ Logistic regression

$$L(t, y(\mathbf{x})) = -\sum_n \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$
 - Hinge loss ⇒ SVM classification

$$L(t, y(\mathbf{x})) = \sum_n [1 - t_n y(\mathbf{x}_n)]_+$$
 - Softmax loss ⇒ Multi-class probabilistic classification

$$L(t, y(\mathbf{x})) = -\sum_n \sum_k \{ \mathbb{I}(t_n = k) \ln \frac{\exp(y_k(\mathbf{x}_n))}{\sum_j \exp(y_j(\mathbf{x}_n))} \}$$

Machine Learning Winter '17 76

B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Multi-Layer Perceptrons

- Adding more layers

Output layer

Hidden layer

Input layer

Output layer

Hidden layer

Input layer

- Output

$$y_k(\mathbf{x}) = g^{(2)}\left(\sum_{i=0}^h W_{ki}^{(2)} g^{(1)}\left(\sum_{j=0}^d W_{ij}^{(1)} x_j\right)\right)$$

Machine Learning Winter '17 77

Slide adapted from Stefan Roth B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Learning with Hidden Units

- How can we train multi-layer networks efficiently?
 - Need an efficient way of adapting **all** weights, not just the last layer.
- Idea: Gradient Descent
 - Set up an error function

$$E(\mathbf{W}) = \sum_n L(t_n, y(\mathbf{x}_n; \mathbf{W})) + \lambda \Omega(\mathbf{W})$$

with a loss $L(\cdot)$ and a regularizer $\Omega(\cdot)$.
 - E.g., $L(t, y(\mathbf{x}; \mathbf{W})) = \sum_n (y(\mathbf{x}_n; \mathbf{W}) - t_n)^2$ L_2 loss
 - $\Omega(\mathbf{W}) = \|\mathbf{W}\|_F^2$ L_2 regularizer ("weight decay")
 - ⇒ Update each weight $W_{ij}^{(k)}$ in the direction of the gradient $\frac{\partial L(\mathbf{W})}{\partial W_{ij}^{(k)}}$

Machine Learning Winter '17 78

B. Leibe

Machine Learning Winter '17

Recap: Gradient Descent

- Two main steps
 - Computing the gradients for each weight
 - Adjusting the weights in the direction of the gradient
- We consider those two steps separately
 - Computing the gradients: [Backpropagation](#)
 - Adjusting the weights: [Optimization techniques](#)

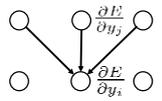
B. Leibe 79

Machine Learning Winter '17

Recap: Backpropagation Algorithm

- Core steps
 - Convert the discrepancy between each output and its target value into an error derivative.

$$E = \frac{1}{2} \sum_{j \in \text{output}} (t_j - y_j)^2$$

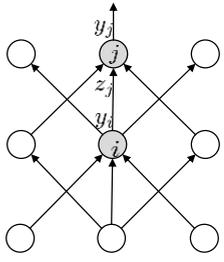
$$\frac{\partial E}{\partial y_j} = -(t_j - y_j)$$
 - Compute error derivatives in each hidden layer from error derivatives in the layer above.
 
 - Use error derivatives w.r.t. activities to get error derivatives w.r.t. the incoming weights

$$\frac{\partial E}{\partial y_j} \rightarrow \frac{\partial E}{\partial w_{ik}}$$

Slide adapted from Geoff Hinton B. Leibe 80

Machine Learning Winter '17

Recap: Backpropagation Algorithm



$$\frac{\partial E}{\partial z_j} = \frac{\partial y_j}{\partial z_j} \frac{\partial E}{\partial y_j} = y_j(1 - y_j) \frac{\partial E}{\partial y_j}$$

$$\frac{\partial E}{\partial y_i} = \sum_j \frac{\partial z_j}{\partial y_i} \frac{\partial E}{\partial z_j} = \sum_j w_{ij} \frac{\partial E}{\partial z_j}$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial z_j}{\partial w_{ij}} \frac{\partial E}{\partial z_j} = y_i \frac{\partial E}{\partial z_j}$$

- Efficient propagation scheme
 - y_i is already known from forward pass! (Dynamic Programming)
 - ⇒ Propagate back the gradient from layer j and multiply with y_i .

Slide adapted from Geoff Hinton B. Leibe 81

Machine Learning Winter '17

Recap: MLP Backpropagation Algorithm

- Forward Pass


```

y(0) = x
for k = 1, ..., l do
  z(k) = W(k) y(k-1)
  y(k) = gk(z(k))
endfor
y = y(l)
E = L(t, y) + λΩ(W)
      
```
- Backward Pass

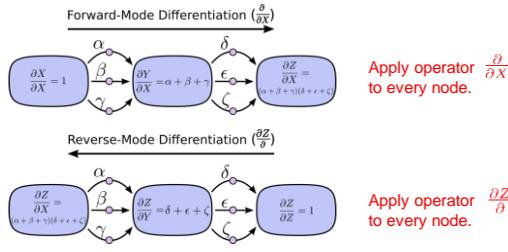

```

h ← ∂E/∂y = ∂E/∂y L(t, y) + λ ∂Ω/∂y
for k = l, l-1, ..., 1 do
  h ← ∂E/∂z(k) = h ∘ g'(k)(y(k))
  ∂E/∂W(k) = h y(k-1)T + λ ∂Ω/∂W(k)
  h ← ∂E/∂y(k-1) = W(k)T h
endfor
      
```
- Notes
 - For efficiency, an entire batch of data \mathbf{X} is processed at once.
 - \odot denotes the element-wise product

B. Leibe 82

Machine Learning Winter '17

Recap: Computational Graphs



Forward-Mode Differentiation ($\frac{\partial Z}{\partial X}$)

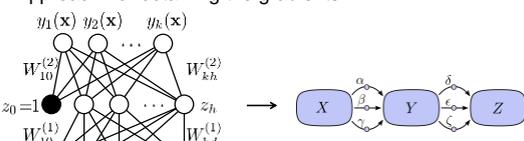
Reverse-Mode Differentiation ($\frac{\partial Z}{\partial X}$)

- Forward differentiation needs one pass per node. Reverse-mode differentiation can compute all derivatives in one single pass.
- ⇒ Speed-up in $\mathcal{O}(\#\text{inputs})$ compared to forward differentiation!

Slide inspired by Christopher Olah B. Leibe Image source: Christopher Olah, colah.github.io 83

Machine Learning Winter '17

Recap: Automatic Differentiation

- Approach for obtaining the gradients
 
 - Convert the network into a computational graph.
 - Each new layer/module just needs to specify how it affects the forward and backward passes.
 - Apply reverse-mode differentiation.
 - ⇒ Very general algorithm, used in today's Deep Learning packages

B. Leibe Image source: Christopher Olah, colah.github.io 84

Recap: Choosing the Right Learning Rate

- Convergence of Gradient Descent
 - Simple 1D example

$$W^{(\tau-1)} = W^{(\tau)} - \eta \frac{dE(W)}{dW}$$
 - What is the optimal learning rate η_{opt} ?
 - If E is quadratic, the optimal learning rate is given by the inverse of the Hessian

$$\eta_{opt} = \left(\frac{d^2 E(W^{(\tau)})}{dW^2} \right)^{-1}$$
 - Advanced optimization techniques try to approximate the Hessian by a simplified form.
 - If we exceed the optimal learning rate, bad things happen!

85
B. Leibe Image source: Yann LeCun et al., Efficient BackProp (1998)

Recap: Advanced Optimization Techniques

- Momentum
 - Instead of using the gradient to change the position of the weight "particle", use it to change the velocity.
 - Effect: dampen oscillations in directions of high curvature
 - Nesterov-Momentum: Small variation in the implementation
- RMS-Prop
 - Separate learning rate for each weight: Divide the gradient by a running average of its recent magnitude.
- AdaGrad
- AdaDelta
- Adam

Some more recent techniques, work better for some problems. Try them.

86
B. Leibe Image source: Geoff Hinton

Recap: Patience

- Saddle points dominate in high-dimensional spaces!

⇒ Learning often doesn't get stuck, you just may have to wait...

87
B. Leibe Image source: Yoshua Bengio

Recap: Reducing the Learning Rate

- Final improvement step after convergence is reached
 - Reduce learning rate by a factor of 10.
 - Continue training for a few epochs.
 - Do this 1-3 times, then stop training.
- Effect
 - Turning down the learning rate will reduce the random fluctuations in the error due to different gradients on different minibatches.
- Be careful: Do not turn down the learning rate too soon!
 - Further progress will be much slower after that.

88
Slide adapted from Geoff Hinton B. Leibe

Recap: Data Augmentation

- Effect
 - Much larger training set
 - Robustness against expected variations
- During testing
 - When cropping was used during training, need to again apply crops to get same image size.
 - Beneficial to also apply flipping during test.
 - Applying several ColorPCA variations can bring another ~1% improvement, but at a significantly increased runtime.

Augmented training data (from one original image)

89
B. Leibe Image source: Lucas Beyer

Recap: Normalizing the Inputs

- Convergence is fastest if
 - The mean of each input variable over the training set is zero.
 - The inputs are scaled such that all have the same covariance.
 - Input variables are uncorrelated if possible.
- Diagram illustrating the process: Mean Cancellation (shifting data to zero mean), Covariance Equalization (scaling to unit variance), and KL-Expansion (decorrelating inputs).
- Advisable normalization steps (for MLPs only, not for CNNs)
 - Normalize all inputs that an input unit sees to zero-mean, unit covariance.
 - If possible, try to decorrelate them using PCA (also known as Karhunen-Loeve expansion).

90
B. Leibe Image source: Yann LeCun et al., Efficient BackProp (1998)

Machine Learning Winter '17

Recap: Another Note on Error Functions

$E(z_n)$

Ideal misclassification error
Squared error
Squared error on tanh

Zero gradient!

No penalty for "too correct" data points!

$t_n \in \{-1, 1\}$

$z_n = t_n y(x_n)$

- Squared error on sigmoid/tanh output function
 - Avoids penalizing "too correct" data points.
 - But: zero gradient for confidently incorrect classifications!
 - ⇒ Do not use L_2 loss with sigmoid outputs (instead: cross-entropy)!

91
Image source: Bishop, 2006

Machine Learning Winter '17

Recap: Commonly Used Nonlinearities

- Sigmoid

$$g(a) = \sigma(a) = \frac{1}{1 + \exp\{-a\}}$$
- Hyperbolic tangent

$$g(a) = \tanh(a) = 2\sigma(2a) - 1$$
- Softmax

$$g(\mathbf{a}) = \frac{\exp\{-a_i\}}{\sum_j \exp\{-a_j\}}$$

92
B. Leibe

Machine Learning Winter '17

Recap: Commonly Used Nonlinearities (2)

- Rectified linear unit (ReLU)

$$g(a) = \max\{0, a\}$$
- Leaky ReLU

$$g(a) = \max\{\beta a, a\} \quad \beta \in [0.01, 0.3]$$
 - Avoids stuck-at-zero units
 - Weaker offset bias
- ELU

$$g(a) = \begin{cases} a, & a \geq 0 \\ e^a - 1, & a < 0 \end{cases}$$
 - No offset bias anymore
 - BUT: need to store activations

93
B. Leibe

Machine Learning Winter '17

Recap: Glorot Initialization

[Glorot & Bengio, '10]

- Variance of neuron activations
 - Suppose we have an input X with n components and a linear neuron with random weights W that spits out a number Y .
 - We want the variance of the input and output of a unit to be the same, therefore $n \text{Var}(W_i)$ should be 1. This means

$$\text{Var}(W_i) = \frac{1}{n} = \frac{1}{n_{\text{in}}}$$
 - Or for the backpropagated gradient

$$\text{Var}(W_i) = \frac{1}{n_{\text{out}}}$$
 - As a compromise, Glorot & Bengio propose to use

$$\text{Var}(W) = \frac{2}{n_{\text{in}} + n_{\text{out}}}$$
 - ⇒ Randomly sample the weights with this variance. That's it.

94
B. Leibe

Machine Learning Winter '17

Recap: He Initialization

[He et al., '15]

- Extension of Glorot Initialization to ReLU units
 - Use Rectified Linear Units (ReLU)

$$g(a) = \max\{0, a\}$$
 - Effect: gradient is propagated with a constant factor

$$\frac{\partial g(a)}{\partial a} = \begin{cases} 1, & a > 0 \\ 0, & \text{else} \end{cases}$$
- Same basic idea: Output should have the input variance
 - However, the Glorot derivation was based on tanh units, linearity assumption around zero does not hold for ReLU.
 - He et al. made the derivations, proposed to use instead

$$\text{Var}(W) = \frac{2}{n_{\text{in}}}$$

95
B. Leibe

Machine Learning Winter '17

Recap: Batch Normalization

[Ioffe & Szegedy '14]

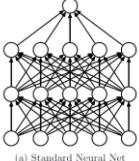
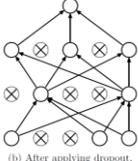
- Motivation
 - Optimization works best if all inputs of a layer are normalized.
- Idea
 - Introduce intermediate layer that centers the activations of the previous layer per minibatch.
 - I.e., perform transformations on all activations and undo those transformations when backpropagating gradients
- Effect
 - Much improved convergence

96
B. Leibe

RWTH AACHEN UNIVERSITY

Recap: Dropout

[Srivastava, Hinton '12]

(a) Standard Neural Net (b) After applying dropout.

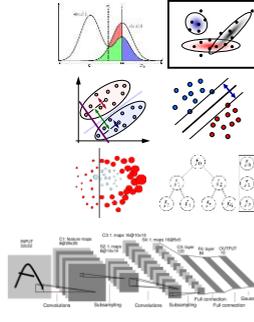
- Idea
 - Randomly switch off units during training.
 - Change network architecture for each data point, effectively training many different variants of the network.
 - When applying the trained network, multiply activations with the probability that the unit was set to zero.
 - ⇒ Improved performance

Machine Learning Winter '17 B. Leibe 97

RWTH AACHEN UNIVERSITY

Course Outline

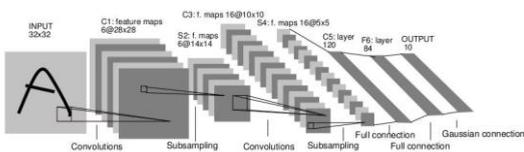
- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks



Machine Learning Winter '17 B. Leibe 98

RWTH AACHEN UNIVERSITY

Recap: Convolutional Neural Networks



- Neural network with specialized connectivity structure
 - Stack multiple stages of feature extractors
 - Higher stages compute more global, more invariant features
 - Classification layer at the end

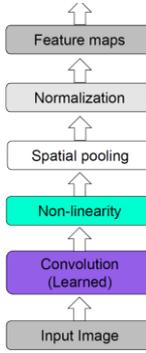
Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, [Gradient-based learning applied to document recognition](#), Proceedings of the IEEE 86(11): 2278–2324, 1998.

Machine Learning Winter '17 B. Leibe 99 Slide credit: Svetlana Lazebnik

RWTH AACHEN UNIVERSITY

Recap: CNN Structure

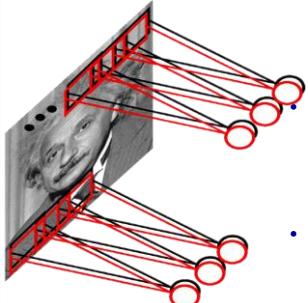
- Feed-forward feature extraction
 1. Convolve input with learned filters
 2. Non-linearity
 3. Spatial pooling
 4. (Normalization)
- Supervised training of convolutional filters by back-propagating classification error



Machine Learning Winter '17 B. Leibe 100 Slide credit: Svetlana Lazebnik

RWTH AACHEN UNIVERSITY

Recap: Intuition of CNNs

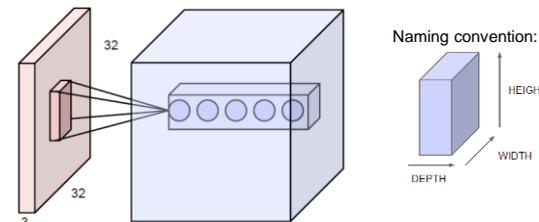


- Convolutional net
 - Share the same parameters across different locations
 - Convolutions with learned kernels
- Learn *multiple* filters
 - E.g. 1000×1000 image
 - 100 filters
 - 10×10 filter size
 - ⇒ only 10k parameters
- Result: Response map
 - size: 1000×1000×100
 - Only memory, not params!

Machine Learning Winter '17 B. Leibe 101 Slide adapted from Marc'Aurelio Ranzato Image source: Yann LeCun

RWTH AACHEN UNIVERSITY

Recap: Convolution Layers



- All Neural Net activations arranged in 3 dimensions
 - Multiple neurons all looking at the same input region, stacked in depth
 - Form a single [1×1×depth] depth column in output volume.

Machine Learning Winter '17 B. Leibe 102 Slide credit: FeiFei Li, Andrej Karpathy

Recap: Visualizing CNNs

Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Machine Learning Winter '17

Slide credit: Yann LeCun

B. Leibe

109

Recap: Residual Networks

AlexNet, 8 layers (ILSVRC 2012) VGG, 19 layers (ILSVRC 2014) ResNet, 152 layers (ILSVRC 2015)

- Core component
 - Skip connections bypassing each layer
 - Better propagation of gradients to the deeper layers
 - This makes it possible to train (much) deeper networks.

$$F(x) = \text{weight layer} \rightarrow \text{relu}$$

$$H(x) = F(x) + x \oplus \text{relu}$$

Machine Learning Winter '17

B. Leibe

110

Recap: Analysis of ResNets

- The effective paths in ResNets are relatively shallow
 - Effectively only 5-17 active modules
- This explains the resilience to deletion
 - Deleting any single layer only affects a subset of paths (and the shorter ones less than the longer ones).
- New interpretation of ResNets
 - ResNets work by creating an ensemble of relatively shallow paths
 - Making ResNets deeper increases the size of this ensemble
 - Excluding longer paths from training does not negatively affect the results.

Machine Learning Winter '17

Image source: Veit et al. 2016

111

Recap: R-CNN for Object Detection

Machine Learning Winter '17

Slide credit: Ross Girshick

B. Leibe

112

Recap: Faster R-CNN

- One network, four losses
 - Remove dependence on external region proposal algorithm.
 - Instead, infer region proposals from same CNN.
 - Feature sharing
 - Joint training

Object detection in a single pass becomes possible.

Machine Learning Winter '17

Slide credit: Ross Girshick

113

Recap: Fully Convolutional Networks

- CNN
- FCN
 - convolutionalization
 - tabby cat heatmap
- Intuition
 - Think of FCNs as performing a sliding-window classification, producing a heatmap of output scores for each class

Machine Learning Winter '17

Image source: Long, Shelhamer, Darrell

114

RWTH AACHEN UNIVERSITY

Recap: Semantic Image Segmentation

- Encoder-Decoder Architecture
 - Problem: FCN output has low resolution
 - Solution: perform upsampling to get back to desired resolution
 - Use skip connections to preserve higher-resolution information

115
Image source: Newell et al.

RWTH AACHEN UNIVERSITY

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

B. Leibe
116

RWTH AACHEN UNIVERSITY

Recap: Neural Probabilistic Language Model

- Core idea
 - Learn a shared distributed encoding (word embedding) for the words in the vocabulary.

117
Slide adapted from Geoff Hinton
B. Leibe
Image source: Geoff Hinton

RWTH AACHEN UNIVERSITY

Recap: word2vec

- Goal
 - Make it possible to learn high-quality word embeddings from huge data sets (billions of words in training set).
- Approach
 - Define two alternative learning tasks for learning the embedding:
 - "Continuous Bag of Words" (CBOW)
 - "Skip-gram"
 - Designed to require fewer parameters.

B. Leibe
118
Image source: Mikolov et al., 2013

RWTH AACHEN UNIVERSITY

Recap: word2vec CBOW Model

- Continuous BOW Model
 - Remove the non-linearity from the hidden layer
 - Share the projection layer for all words (their vectors are averaged)

⇒ Bag-of-Words model (order of the words does not matter anymore)

B. Leibe
119
Image source: Xin Bing, 2015

RWTH AACHEN UNIVERSITY

Recap: word2vec Skip-Gram Model

- Continuous Skip-Gram Model
 - Similar structure to CBOW
 - Instead of predicting the current word, predict words within a certain range of the current word.
 - Give less weight to the more distant words

B. Leibe
120
Image source: Xin Bing, 2015

Machine Learning Winter '17

Recap: Problems with 100k-1M outputs

- Weight matrix gets huge!
 - Example: CBOW model
 - One-hot encoding for inputs
 - Input-hidden connections are just vector lookups.
 - This is not the case for the hidden-output connections!
 - State h is not one-hot, and vocabulary size is 1M.
 - $\Rightarrow W'_{N \times V}$ has $300 \times 1M$ entries
- Softmax gets expensive!
 - Need to compute normalization over 100k-1M outputs

B. Leibe 121
Image source: Yin Bonn, 2015

Machine Learning Winter '17

Recap: Hierarchical Softmax

- Idea
 - Organize words in binary search tree, words are at leaves
 - Factorize probability of word w_0 as a product of node probabilities along the path.
 - Learn a linear decision function $y = v_{n(w,j)} \cdot h$ at each node to decide whether to proceed with left or right child node.
 - \Rightarrow Decision based on output vector of hidden units directly.

B. Leibe 122
Image source: Yin Bonn, 2015

Machine Learning Winter '17

Recap: Recurrent Neural Networks

- Up to now
 - Simple neural network structure: 1-to-1 mapping of inputs to outputs
- Recurrent Neural Networks
 - Generalize this to arbitrary mappings

B. Leibe 123
Image source: Andrej Karpathy

Machine Learning Winter '17

Recap: Recurrent Neural Networks (RNNs)

- RNNs are regular NNs whose hidden units have additional connections over time.
 - You can unroll them to create a network that extends over time.
 - When you do this, keep in mind that the weights for the hidden are shared between temporal layers.
- RNNs are very powerful
 - With enough neurons and time, they can compute anything that can be computed by your computer.

B. Leibe 124
Image source: Andrej Karpathy

Machine Learning Winter '17

Recap: Backpropagation Through Time (BPTT)

- Configuration

$$h_t = \sigma(W_{xh}x_t + W_{hh}h_{t-1} + b)$$
- Backpropagated gradient

$$\hat{y}_t = \text{softmax}(W_{hy}h_t)$$
 - For weight w_{ij} :

$$\frac{\partial E_t}{\partial w_{ij}} = \sum_{1 \leq k \leq t} \left(\frac{\partial E_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial^+ h_k}{\partial w_{ij}} \right)$$

B. Leibe 125

Machine Learning Winter '17

Recap: Backpropagation Through Time (BPTT)

- Analyzing the terms
 - For weight w_{ij} :

$$\frac{\partial E_t}{\partial w_{ij}} = \sum_{1 \leq k \leq t} \left(\frac{\partial E_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial^+ h_k}{\partial w_{ij}} \right)$$
 - This is the "immediate" partial derivative (with h_{k-1} as constant)

B. Leibe 126

Machine Learning Winter '17

Recap: Backpropagation Through Time (BPTT)

- Analyzing the terms
 - For weight w_{ij} :
$$\frac{\partial E_t}{\partial w_{ij}} = \sum_{1 \leq k \leq t} \left(\frac{\partial E_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial^+ h_k}{\partial w_{ij}} \right)$$
 - Propagation term:
$$\frac{\partial h_t}{\partial h_k} = \prod_{t \geq i > k} \frac{\partial h_i}{\partial h_{i-1}}$$

127

Machine Learning Winter '17

Recap: Backpropagation Through Time (BPTT)

- Summary
 - Backpropagation equations

$$E = \sum_{1 \leq t \leq T} E_t$$

$$\frac{\partial E_t}{\partial w_{ij}} = \sum_{1 \leq k \leq t} \left(\frac{\partial E_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial^+ h_k}{\partial w_{ij}} \right)$$

$$\frac{\partial h_t}{\partial h_k} = \prod_{t \geq i > k} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{t \geq i > k} \mathbf{W}_{hh}^\top \text{diag}(\sigma'(\mathbf{h}_{i-1}))$$
 - Remaining issue: how to set the initial state \mathbf{h}_0 ?
 - ⇒ Learn this together with all the other parameters.

B. Leibe 128

Machine Learning Winter '17

Recap: Exploding / Vanishing Gradient Problem

- BPTT equations:

$$\frac{\partial E_t}{\partial w_{ij}} = \sum_{1 \leq k \leq t} \left(\frac{\partial E_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial^+ h_k}{\partial w_{ij}} \right)$$

$$\frac{\partial h_t}{\partial h_k} = \prod_{t \geq i > k} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{t \geq i > k} \mathbf{W}_{hh}^\top \text{diag}(\sigma'(\mathbf{h}_{i-1}))$$

$$= (\mathbf{W}_{hh}^\top)^l$$

(if t goes to infinity and $l = t - k$.)
- ⇒ We are effectively taking the weight matrix to a high power.
 - The result will depend on the eigenvalues of \mathbf{W}_{hh} .
 - Largest eigenvalue > 1 ⇒ Gradients *may* explode.
 - Largest eigenvalue < 1 ⇒ Gradients *will* vanish.
 - This is very bad...

B. Leibe 129

Machine Learning Winter '17

Recap: Gradient Clipping

- Trick to handle exploding gradients
 - If the gradient is larger than a threshold, clip it to that threshold.

Algorithm 1 Pseudo-code

```

 $\bar{\mathbf{g}} \leftarrow \frac{\partial E}{\partial \theta}$ 
if  $\|\bar{\mathbf{g}}\| \geq \text{threshold}$  then
   $\bar{\mathbf{g}} \leftarrow \frac{\text{threshold}}{\|\bar{\mathbf{g}}\|} \bar{\mathbf{g}}$ 
end if

```

- This makes a big difference in RNNs

Slide adapted from Richard Socher B. Leibe 130

Machine Learning Winter '17

Recap: Long Short-Term Memory

- LSTMs
 - Inspired by the design of memory cells
 - Each module has 4 layers, interacting in a special way.

Image source: Christopher Olah, <http://crdlob.github.io/colah/2015-08-11-understanding-lstm-1/>

131

Machine Learning Winter '17

Recap: Elements of LSTMs

- Forget gate layer
 - Look at \mathbf{h}_{t-1} and \mathbf{x}_t and output a number between 0 and 1 for each dimension in the cell state \mathbf{C}_{t-1} .
 - 0: completely delete this,
 - 1: completely keep this.
- Update gate layer
 - Decide what information to store in the cell state.
 - Sigmoid network (input gate layer) decides which values are updated.
 - tanh layer creates a vector of new candidate values that could be added to the state.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{\mathbf{C}}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

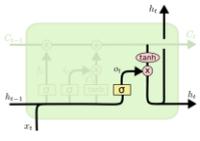
Source: Christopher Olah, <http://crdlob.github.io/colah/2015-08-11-understanding-lstm-1/>

132

RWTH AACHEN
UNIVERSITY

Recap: Elements of LSTMs

- **Output gate layer**
 - Output is a filtered version of our gate state.
 - First, apply sigmoid layer to decide what parts of the cell state to output.
 - Then, pass the cell state through a tanh (to push the values to be between -1 and 1) and multiply it with the output of the sigmoid gate.



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

133

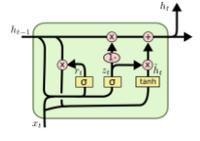
Source: Christopher Olah, <http://colah.github.io/posts/2015-08-Understandin-LSTMs/>

Machine Learning Winter '17

RWTH AACHEN
UNIVERSITY

Recap: Gated Recurrent Units (GRU)

- **Simpler model than LSTM**
 - Combines the forget and input gates into a single **update gate** z_t .
 - Similar definition for a **reset gate** r_t , but with different weights.
 - In both cases, merge the cell state and hidden state.
- **Empirical results**
 - Both LSTM and GRU can learn much longer-term dependencies than regular RNNs
 - GRU performance similar to LSTM (no clear winner yet), but fewer parameters.



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

134

Source: Christopher Olah, <http://colah.github.io/posts/2015-08-Understandin-LSTMs/>

Machine Learning Winter '17

RWTH AACHEN
UNIVERSITY

Any Questions?

So what can you do with all of this?

135

Machine Learning Winter '17

RWTH AACHEN
UNIVERSITY

Any More Questions?

Good luck for the exam!

136

Machine Learning Winter '17