

RWTH AACHEN
UNIVERSITY

Machine Learning – Lecture 19

Repetition

31.01.2019

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

Machine Learning Winter '18

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

Announcements (3)

- Seminar in the summer semester
 - Current topics in Computer Vision and Machine Learning
 - Quick poll: Who is interested?

B. Leibe

4

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
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

B. Leibe

5

RWTH AACHEN
UNIVERSITY

Recap: Bayes Decision Theory

B. Leibe

6

Slide credit: Bernt Schiele Image source: C.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)

$$\frac{p(x|C_1)}{p(x|C_2)} > \frac{p(C_2)}{p(C_1)}$$

Decision threshold θ

Machine Learning Winter '18 | Slide credit: Bernt Schiele | B. Leibe | 7

RWTH AACHEN UNIVERSITY

Recap: Bayes Decision Theory

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

Machine Learning Winter '18 | Slide credit: Bernt Schiele | B. Leibe | 8

RWTH AACHEN UNIVERSITY

Recap: Classifying with Loss Functions

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

$$L_{kj} = \text{loss for decision } C_j \text{ if truth is } C_k.$$
- Example: cancer diagnosis

		Decision	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

Machine Learning Winter '18 | B. Leibe | 9

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**

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, C_k) d\mathbf{x}$$
- This can be done by choosing the regions \mathcal{R}_j such that

$$\mathbb{E}[L] = \sum_k L_{kj} p(C_k | \mathbf{x})$$
 which is easy to do once we know the posterior class probabilities $p(C_k | \mathbf{x})$

see Exercise 1.3

Machine Learning Winter '18 | B. Leibe | 10

RWTH AACHEN UNIVERSITY

Recap: The Reject Option

- Classification errors arise from regions where the largest posterior probability $p(C_k|x)$ 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 '18 | B. Leibe | Image source: G.M. Bishop, 2006 | 11

RWTH AACHEN UNIVERSITY

Course Outline

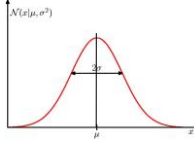
- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
 - Mixture Models and EM
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

Machine Learning Winter '18 | B. Leibe | 12

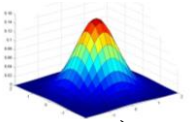
RWTH AACHEN UNIVERSITY

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

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


13
B. Leibe Image source: C.M. Bishop, 2006

RWTH AACHEN UNIVERSITY

Recap: Maximum Likelihood Approach

see Exercise 1.4

- 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.

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

14
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

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

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

15
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

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

Normalization: integrate over all possible values of θ

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

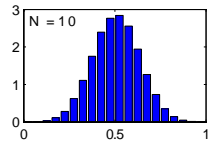
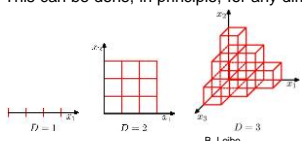
16
B. Leibe

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 ...

17
B. Leibe Image source: C.M. Bishop, 2006

RWTH AACHEN UNIVERSITY

Recap: Kernel Density Estimation

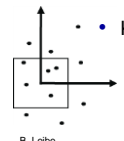
see Exercise 1.5

$$p(x) \approx \frac{K}{NV}$$

fixed V determine K fixed K determine V

Kernel Methods K-Nearest Neighbor

- Kernel methods
 - Place a kernel window k at location 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.



18
Slide adapted from Bernt Schiele B. Leibe

Machine Learning Winter '18

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 19

Machine Learning Winter '18

Recap: Mixture of Gaussians (MoG)

- "Generative model"

$$p(j) = \pi_j \quad \text{"Weight" of mixture component}$$

$$p(x|\theta) = \sum_{j=1}^M p(x|\theta_j) p(j)$$

B. Leibe 20

Machine Learning Winter '18

Recap: MoG – Iterative Strategy

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

$$h(j=1|x_n) = \begin{matrix} 1 & 1 & 1 \\ 0 & 0 & 0 \end{matrix} \quad h(j=2|x_n) = \begin{matrix} 0 & 0 & 0 \\ 1 & 1 & 1 \end{matrix}$$

$$\mu_1 = \frac{\sum_{n=1}^N h(j=1|x_n) x_n}{\sum_{i=1}^N h(j=1|x_n)} \quad \mu_2 = \frac{\sum_{n=1}^N h(j=2|x_n) x_n}{\sum_{i=1}^N h(j=2|x_n)}$$

B. Leibe 21

Machine Learning Winter '18

Recap: MoG – Iterative Strategy

- Assuming we knew the mixture components...

$$p(j=1|x) > p(j=2|x)$$

- Bayes decision rule: Decide $j=1$ if

B. Leibe 22

Machine Learning Winter '18

Recap: K-Means Clustering

- 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.

B. Leibe 23

Machine Learning Winter '18

Recap: EM Algorithm

- 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

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

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

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

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

see Exercise 1.6

B. Leibe 24

Machine Learning Winter '18

Course Outline

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

B. Leibe 25

Machine Learning Winter '18

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} , 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**.

Slide adapted from Bernd Schölkopf. B. Leibe 26

Machine Learning Winter '18

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

$$y(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T (\tilde{\mathbf{X}}^\dagger)^T \tilde{\mathbf{x}}$$
- ⇒ Exact, closed-form solution for the discriminant function parameters.

see Exercise 2.1

B. Leibe 27

Machine Learning Winter '18

Recap: Problems with Least Squares

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

Image source: G.M. Bishop, 2006. B. Leibe 28

Machine Learning Winter '18

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.

B. Leibe 29

Machine Learning Winter '18

Recap: Linear Separability

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

Slide credit: Bernd Schölkopf. B. Leibe 30

RWTH AACHEN UNIVERSITY

Recap: Extension to Nonlinear Basis Fcts.

- Generalization
 - Transform vector \mathbf{x} with M nonlinear basis functions $\phi_j(\mathbf{x})$:

$$y_k(\mathbf{x}) = \sum_{j=1}^M w_{ki} \phi_j(\mathbf{x}) + w_{k0}$$
- 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.
 - ➔ Minimization with Gradient Descent

31

RWTH AACHEN UNIVERSITY

Recap: Probabilistic Discriminative Models

- Consider models of the form

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$
 with

$$p(C_2|\phi) = 1 - p(C_1|\phi)$$
- 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 $p(\phi|C_k)$ and $p(C_k)$.

32

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 $y_n = p(C_1|\phi_n)$, 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**.

33

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 $\mathbf{H} = \nabla \nabla E(\mathbf{w})$; the Hessian matrix, i.e. the matrix of second derivatives.
 - Local quadratic approximation to the target function
 - Faster convergence

34

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)**

35

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.

36

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.

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

$$\nabla_{\mathbf{w}_k} E(\mathbf{w}) = - \sum_{n=1}^N \left[\mathbb{I}(t_n = k) \ln P(y_n = k | \mathbf{x}_n; \mathbf{w}) \right]$$

37

RWTH AACHEN UNIVERSITY

Course Outline

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

38

RWTH AACHEN UNIVERSITY

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.
 - \Rightarrow *Overfitting to the training set!*

39

RWTH AACHEN UNIVERSITY

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

40

RWTH AACHEN UNIVERSITY

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

$a_n \geq 0$ $t_n y(\mathbf{x}_n) - 1 \geq 0$ $a_n \{t_n y(\mathbf{x}_n) - 1\} = 0$	KKT: $\lambda \geq 0$ $f(\mathbf{x}) \geq 0$ $\lambda f(\mathbf{x}) = 0$
---	--

41

RWTH AACHEN UNIVERSITY

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
 - \Rightarrow Only the SVs actually influence the decision boundary!
 - Compute b by averaging over all support vectors:

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

42

Machine Learning Winter '18

RWTH AACHEN UNIVERSITY

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!

43

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

Machine Learning Winter '18

RWTH AACHEN UNIVERSITY

Recap: SVM – Dual Formulation

- Maximize

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

44

Slide adapted from Bernt Schiele B. Leibe

Machine Learning Winter '18

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} .

45

B. Leibe

Machine Learning Winter '18

RWTH AACHEN UNIVERSITY

Recap: SVM – New Dual Formulation

- New SVM Dual: Maximize

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

$$0 \leq a_n \leq C$$

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

This is all that changed!
- This is again a quadratic programming problem
⇒ Solve as before...

46

Slide adapted from Bernt Schiele B. Leibe

Machine Learning Winter '18

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:

$\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$

47

Slide credit: Raymond Moorey

Machine Learning Winter '18

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!)

48

B. Leibe

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...

49

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

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

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

50

RWTH AACHEN UNIVERSITY

Recap: Nonlinear SVM – Dual Formulation

- SVM Dual: Maximize

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

see Exercise 2.2

 under the conditions

$$0 \leq a_n \leq C$$

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

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

51

RWTH AACHEN UNIVERSITY

Course Outline

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

52

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?

53

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.

54

RWTH AACHEN UNIVERSITY

Recap: AdaBoost – “Adaptive Boosting”

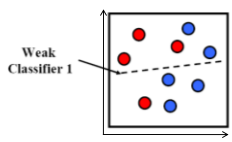
- Main idea [Freund & Schapire, 1996]
 - Instead of resampling, reweight misclassified training examples.
 - Increase the chance of being selected in a sampled training set.
 - Or increase the misclassification cost when training on the full set.
- Components
 - $h_m(\mathbf{x})$: “weak” or base classifier
 - 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:

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

55

RWTH AACHEN UNIVERSITY

Recap: AdaBoost – Intuition



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

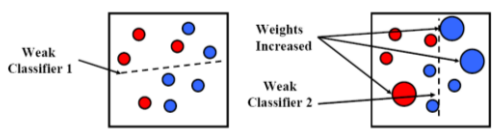
Each weak classifier splits the training examples with at least 50% accuracy.

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

56

RWTH AACHEN UNIVERSITY

Recap: AdaBoost – Intuition



Weak Classifier 1

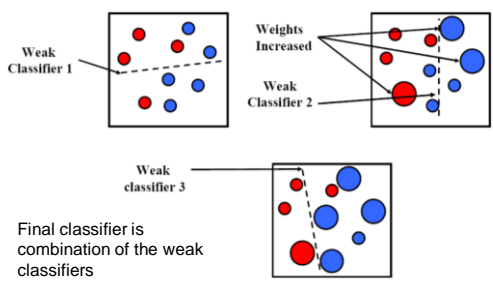
Weights Increased

Weak Classifier 2

57

RWTH AACHEN UNIVERSITY

Recap: AdaBoost – Intuition



Weak Classifier 1

Weights Increased

Weak Classifier 2

Weak classifier 3

Final classifier is combination of the weak classifiers

58

RWTH AACHEN UNIVERSITY

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

$$J_m = \sum_{n=1}^N w_n^{(m)} I(h_m(\mathbf{x}) \neq t_n) \quad I(A) = \begin{cases} 1, & \text{if } A \text{ is true} \\ 0, & \text{else} \end{cases}$$
 - 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 $h_m(\mathbf{x})$:

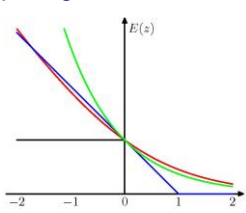
$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}$$
 - d) Update the weighting coefficients:

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

59

RWTH AACHEN UNIVERSITY

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!

60

Machine Learning Winter '18

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"

61
Image source: Bishop, 2006

Machine Learning Winter '18

Course Outline

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

B. Leibe

62

Machine Learning Winter '18

Recap: Perceptrons

- One output node per class

Output layer
Weights
Input layer

- Outputs
 - Linear outputs

With output nonlinearity

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

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

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

Slide adapted from Stefan Roth

B. Leibe

63

Machine Learning Winter '18

Recap: Non-Linear Basis Functions

- Straightforward generalization

Output layer
Weights
Feature layer
Mapping (fixed)
Input layer

- Outputs
 - Linear outputs

with output nonlinearity

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

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

Slide adapted from Geoff Hinton

B. Leibe

64

Machine Learning Winter '18

Recap: Non-Linear Basis Functions

- Straightforward generalization

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!

B. Leibe

65

Machine Learning Winter '18

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!

Slide adapted from Geoff Hinton

B. Leibe

66

RWTH AACHEN UNIVERSITY

Recap: Loss Functions

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

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

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

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

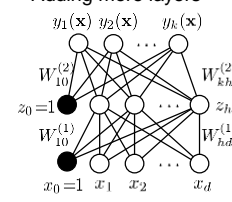
$$L(t, y(\mathbf{x})) = \sum_n [1 - t_n y(\mathbf{x}_n)]_+$$
 - Softmax loss \Rightarrow Multi-class probabilistic classification

$$L(t, y(\mathbf{x})) = -\sum_n \sum_k \left\{ \mathbb{1}(t_n = k) \ln \frac{\exp(y_k(\mathbf{x}))}{\sum_j \exp(y_j(\mathbf{x}))} \right\}$$

67

RWTH AACHEN UNIVERSITY

Recap: Multi-Layer Perceptrons

- Adding more layers
 
- Output

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

68

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")
 - \Rightarrow Update each weight $W_{ij}^{(k)}$ in the direction of the gradient $\frac{\partial E(\mathbf{W})}{\partial W_{ij}^{(k)}}$

69

RWTH AACHEN UNIVERSITY

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**

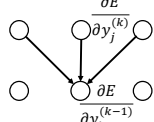
70

RWTH AACHEN UNIVERSITY

Recap: Backpropagation Algorithm

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

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

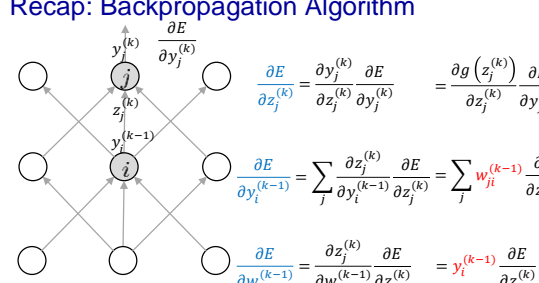
$$\frac{\partial E}{\partial y_j} = -(t_j - y_j)$$
 - Compute error derivatives in each hidden layer from error derivatives in the layer above.
 
 - Use error derivatives w.r.t. activities to get error derivatives w.r.t. the incoming weights

$$\frac{\partial E}{\partial y_j^{(k)}} \rightarrow \frac{\partial E}{\partial w_{ji}^{(k-1)}}$$

71

RWTH AACHEN UNIVERSITY

Recap: Backpropagation Algorithm



- Efficient propagation scheme
 - $y_i^{(k-1)}$ is already known from forward pass! (Dynamic Programming)
 \Rightarrow Propagate back the gradient from layer k and multiply with $y_i^{(k-1)}$.

72

RWTH AACHEN UNIVERSITY

Recap: MLP Backpropagation Algorithm

- Forward Pass**

$$\mathbf{y}^{(0)} = \mathbf{x}$$
for $k = 1, \dots, l$ **do**

$$\mathbf{z}^{(k)} = \mathbf{W}^{(k)} \mathbf{y}^{(k-1)}$$

$$\mathbf{y}^{(k)} = g_k(\mathbf{z}^{(k)})$$
endfor

$$\mathbf{y} = \mathbf{y}^{(l)}$$

$$E = L(\mathbf{t}, \mathbf{y}) + \lambda \Omega(\mathbf{W})$$
- Backward Pass**

$$\mathbf{h} \leftarrow \frac{\partial E}{\partial \mathbf{y}} = \frac{\partial}{\partial \mathbf{y}} L(\mathbf{t}, \mathbf{y}) + \lambda \frac{\partial}{\partial \mathbf{y}} \Omega$$
for $k = l, l-1, \dots, 1$ **do**

$$\mathbf{h} \leftarrow \frac{\partial E}{\partial \mathbf{z}^{(k)}} = \mathbf{h} \odot g'(\mathbf{y}^{(k)})$$

$$\frac{\partial E}{\partial \mathbf{W}^{(k)}} = \mathbf{h} \mathbf{y}^{(k-1)\top} + \lambda \frac{\partial \Omega}{\partial \mathbf{W}^{(k)}}$$

$$\mathbf{h} \leftarrow \frac{\partial E}{\partial \mathbf{y}^{(k-1)}} = \mathbf{W}^{(k)\top} \mathbf{h}$$
endfor

Machine Learning Winter '18

73

RWTH AACHEN UNIVERSITY

Recap: Computational Graphs

Forward-Mode Differentiation ($\frac{\partial}{\partial \mathbf{x}}$)

Reverse-Mode Differentiation ($\frac{\partial \mathbf{z}}{\partial \mathbf{y}}$)

Apply operator $\frac{\partial}{\partial \mathbf{x}}$ to every node.

Apply operator $\frac{\partial \mathbf{z}}{\partial \mathbf{y}}$ to every node.

- 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!

Machine Learning Winter '18

74

RWTH AACHEN UNIVERSITY

Recap: Automatic Differentiation

- Approach for obtaining the gradients**

see Exercise 4.2

 - 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

Machine Learning Winter '18

75

RWTH AACHEN UNIVERSITY

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

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

Machine Learning Winter '18

76

RWTH AACHEN UNIVERSITY

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.

Machine Learning Winter '18

77

RWTH AACHEN UNIVERSITY

Recap: Patience

- Saddle points dominate in high-dimensional spaces!

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

Machine Learning Winter '18

78

Machine Learning Winter '18

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.

Slide adapted from Geoff Hinton

B. Leibe

79

Machine Learning Winter '18

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)

B. Leibe

Image source: Lucas Beyer

80

Machine Learning Winter '18

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.
- 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).

B. Leibe

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

81

Machine Learning Winter '18

Recap: Another Note on Error Functions

- 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)!

B. Leibe

Image source: Bishop, 2006

82

Machine Learning Winter '18

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

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

B. Leibe

83

Machine Learning Winter '18

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

B. Leibe

84

RWTH AACHEN UNIVERSITY

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

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

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

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

⇒ Randomly sample the weights with this variance. That's it.

85

RWTH AACHEN UNIVERSITY

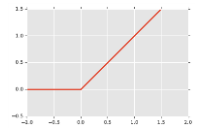
Recap: He Initialization [He et al., '15]

- 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

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



86

RWTH AACHEN UNIVERSITY

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
 - (Typically) much improved convergence

87

RWTH AACHEN UNIVERSITY

Recap: Dropout [Srivastava, Hinton '12]

- 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.

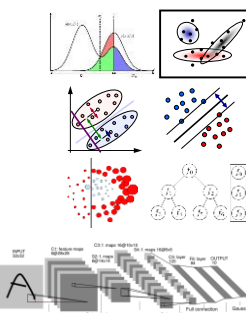
⇒ Improved performance

88

RWTH AACHEN UNIVERSITY

Course Outline

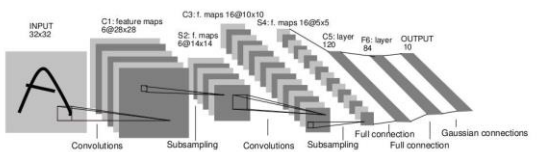
- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks



89

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.

90

Machine Learning Winter '18

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

Slide credit: Svetlana Lazebnik B. Leibe 91

Machine Learning Winter '18

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!

Slide adapted from Marc'Aurelio Ranzato B. Leibe Image source: Yann Lecun 92

Machine Learning Winter '18

Recap: Convolution Layers

Naming convention:

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

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe 93

Machine Learning Winter '18

Recap: Activation Maps

Activations: one filter = one depth slice (or activation map)

5×5 filters

Each activation map is a depth slice through the output volume.

Activation maps

Slide adapted from FeiFei Li, Andrei Karpathy B. Leibe 94

Machine Learning Winter '18

Recap: Pooling Layers

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters and stride 2

6	8
3	4

- Effect:
 - Make the representation smaller without losing too much information
 - Achieve robustness to translations

Slide adapted from FeiFei Li, Andrei Karpathy B. Leibe 95

Machine Learning Winter '18

Recap: AlexNet (2012)

- Similar framework as LeNet, but
 - Bigger model (7 hidden layers, 650k units, 60M parameters)
 - More data (10^6 images instead of 10^3)
 - GPU implementation
 - Better regularization and up-to-date tricks for training (Dropout)

A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012.

Image source: A. Krizhevsky, I. Sutskever and G.F. Hinton, NIPS 2012 96

RWTH AACHEN UNIVERSITY

Recap: R-CNN for Object Detection

Machine Learning Winter '18

Slide credit: Ross Girshick

B. Leibe

103

RWTH AACHEN UNIVERSITY

Recap: Faster R-CNN for Object Detection

- 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 '18

Slide credit: Ross Girshick

104

RWTH AACHEN UNIVERSITY

Recap: Fully Convolutional Networks

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

Machine Learning Winter '18

Image source: Long, Shelhamer, Darrell

B. Leibe

105

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

Machine Learning Winter '18

Image source: Newell et al

106

RWTH AACHEN UNIVERSITY

Course Outline

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

Machine Learning Winter '18

B. Leibe

107

RWTH AACHEN UNIVERSITY

Recap: Neural Probabilistic Language Model

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

Y. Bengio, R. Ducharme, P. Vincent, C. Jauvin, [A Neural Probabilistic Language Model](#), In JMLR, Vol. 3, pp. 1137-1155, 2003.

Machine Learning Winter '18

Slide adapted from Geoff Hinton

B. Leibe

108

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.

109
B. Leibe
Image source: Mikolov et al., 2013

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)

110
B. Leibe
Image source: Yin Rong, 2015

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

111
B. Leibe
Image source: Yin Rong, 2015

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.
 - ⇒ $W'_{N \times V}$ has $300 \times 1M$ entries
- Softmax gets expensive!
 - Need to compute normalization over 100k-1M outputs

112
B. Leibe
Image source: Yin Rong, 2015

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.
 - ⇒ Decision based on output vector of hidden units directly.

113
B. Leibe
Image source: Yin Rong, 2015

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

114
B. Leibe
Image source: Andrej Karpathy

Machine Learning Winter '18

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 115
Image source: Andrej Karpathy

Machine Learning Winter '18

Recap: Backpropagation Through Time (BPTT)

- Configuration

$$\mathbf{h}_t = \sigma(\mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + b)$$
- Backpropagated gradient
 - For weight w_{ij} :

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

116

Machine Learning Winter '18

Recap: Backpropagation Through Time (BPTT)

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

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

117

Machine Learning Winter '18

Recap: Backpropagation Through Time (BPTT)

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

$$\frac{\partial E_t}{\partial w_{ij}} = \sum_{1 \leq k \leq t} \left(\frac{\partial E_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} \frac{\partial^+ h_k}{\partial w_{ij}} \right)$$
 - Propagation term:

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} = \prod_{t \geq i > k} \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}}$$

118

Machine Learning Winter '18

Recap: Backpropagation Through Time (BPTT)

- Summary
 - Backpropagation equations

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

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

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

see Exercise 6.1

B. Leibe 119

Machine Learning Winter '18

Recap: Exploding / Vanishing Gradient Problem

- BPTT equations:

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

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

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

(if l 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 120

Machine Learning Winter '18

Recap: Gradient Clipping

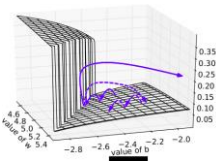
- Trick to handle exploding gradients
 - If the gradient is larger than a threshold, clip it to that threshold.

Algorithm 1 Pseudo-code

```

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

```



- This makes a big difference in RNNs

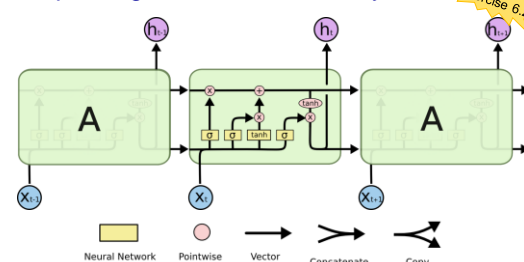
Slide adapted from Richard Socher. B. Leibe

121

Machine Learning Winter '18

Recap: Long Short-Term Memory

see Exercise 6.2



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

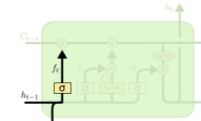
Image source: Christopher Olah, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

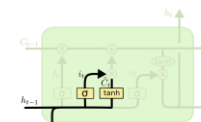
122

Machine Learning Winter '18

Recap: Elements of LSTMs

- Forget gate layer**
 - Look at h_{t-1} and x_t and output a number between 0 and 1 for each dimension in the cell state 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.



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


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

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

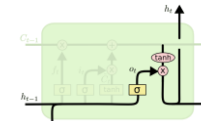
Source: Christopher Olah, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

124

Machine Learning Winter '18

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.



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

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

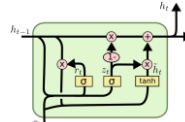
Source: Christopher Olah, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

124

Machine Learning Winter '18

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.



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

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

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

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

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

125

Machine Learning Winter '18

Any More Questions?

Good luck for the exam!

126