

Maschine Learning Summary WS21/22

mastermakrela

February 4, 2022

Contents

Introduction	4
Core Questions	4
Core Problems	4
Core Questions	4
Bayes Decision Theory	5
Probability Theory	5
Probability	5
Probability Densities	5
Expectations	5
Variances and Covariances	6
Bayes Decision Theory	6
Classifying with Loss Functions	6
Minimizing the Expected Loss	6
Probability Density Estimation	6
The Gaussian (or Normal) Distribution	7
Properties	7
Parametric Methods	7
Maximum Likelihood Approach	7
Non-Parametric Methods	8
Histograms	8
Kernel Density Estimation	8
K-Nearest Neighbors	9
Summary	9
Mixture Distribution	9
Mixture of Gaussians	9
Maximum Likelihood	9
K-Means Clustering	10
EM Algorithm	10
Summary: Gaussian Mixture Models	10
Applications	11
Support Vector Machines	11
Softmax Regression	11
Multi-class generalization of logistic regression	11
Optimization	12
Support Vector Machines	12
Motivation	12
SVM	12
Lagrangian Formulation	12
Solution	12

Discussion	13
Dual form formulation	13
Soft-margin classification	13
Nonlinear Support Vector Machines	14
Problem with High-dim. Basis Functions	14
Summary	14
Properties	14
Limitaitons	15
Error Function	15
Applications	15
Text Classification	15
Adaboost	15
Ensembles of Classifiers	15
Idea	15
Constructing Ensembles	15
Adaboost	17
Idea	17
Components	17
Adaboost	17
Minimizing Exponential Error	17
Final Algorithm	17
Summary	17
Deep Learning / Neural Networks	18
A Brief History of Neural Networks	18
Perceptrons	18
Standard Perceptron	18
Perceptron Learning	19
Limitations of Perceptrons	19
Multi-Layer Perceptrons	20
Learning with Hidden Units	20
Obtaining the Gradients	20
Approach 1: Naive Analytical Differentiation	20
Approach 2: Numerical Differentiation	20
Approach 3: Incremental Analytical Differentiation (Backpropagation)	21
Learning Multi-layer Networks	21
Computational Graphs	21
Approach 4: Automatic Differentiation	22
Gradient Descent	22
Stochastic vs. Batch Learning	22
Choosing the Right Learning Rate	23
Momentum	23
Better Adaptation: RMSProp	24
Summary	24
Optimization	25
Tricks of the Trade	25
Shuffling the Examples	25
Data augmentation	25
Normalization	25
Nonlinearities	25
Usage	26
Extension: ReLU	26
Initialization	26
Initializing the Weights	27
Glorot Initialization	27

Advanced techniques	27
Batch Normalization	27
Dropout	27

Introduction

Machine Learning Principles, methods, and algorithms for learning and prediction on the basis of past evidence

Goal of Machine Learning:

Machines that *learn* to *perform* a *task* from *experience*.

Learning most important aspect

We provide *data* and *goal* - machine figures rest out.

Learning Tools:

- statistics
- probability theory
- decision theory
- information theory
- optimization theory

Core Questions

Task: $y = f(x; w)$

Where:

- x Input
- y Output
- w Learned parameters

Regression	Classification
Continuous output	Discrete output

Core Problems

1. How to input data / how to interpret inputted data
2. Features
 - Invariance to irrelevant input variations
 - Selecting the “right” features is crucial
 - Encoding and use of “domain knowledge”
 - Higher-dimensional features are more discriminative.
3. Curse of Dimensionality
 - complexity increases exponentially with number of dimensions

Core Questions

1. Measuring performance of a model
 - eg. % of correct classifications
2. Generalization performance
 - performing on test data is not enough
 - model has to perform on new data
3. What data is available?
 - Supervised vs unsupervised learning
 - mix: semi-supervised
 - reinforcement learning - with feedback

Most often learning is an optimization problem

I.e., maximize $y = f(x; w)$ with regard to performance.

Bayes Decision Theory

Probability Theory

Probability

$$P(X = x) \in [0, 1] \quad (1)$$

Where $X \in \{x_1, \dots, x_N\}$ an Occurrence of something.

Assuming two random variables $X \in \{x_i\}$ and $Y \in \{y_i\}$. Consider N trials and let:

$$n_{ij} = \text{count}\{X = x_i \wedge Y = y_j\}$$

$$c_i = \text{count}\{X = x_i\}$$

$$r_i = \text{count}\{Y = y_j\}$$

Then we can derive *The Rules of Probability*:

Joint Probability	$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$
Marginal Probability	$p(X = x_i) = \frac{c_i}{N}$
Conditional Probability	$p(Y = y_j X = x_i) = \frac{n_{ij}}{c_i}$
Sum Rule	$p(X) = \sum_Y p(X, Y)$
Product Rule	$p(X, Y) = p(X, Y)P(X)$

And *Bayes' Theorem*:

$$p(X|Y) = \frac{p(X|Y)P(X)}{P(X)} \quad (2)$$

where: $p(X) = \sum_Y p(X, Y)p(Y)$

Probability Densities

If the variable is continuous we can't just look at probability for x . We have to look for the interval the x is in, using *Probability Density Function* $p(x)$.

$$p(x) = \int_a^b p(x)dx$$

The probability that x lies in the interval is given by $(-\infty, z)$ the *cumulative distribution function*:

$$P(z) = \int_{-\infty}^z p(x)dx$$

Expectations

Expectation average value of some function $f(x)$ under a probability distribution $p(x)$

$$\text{Discrete: } \mathbb{E}[f] = \sum_x p(x)f(x)$$

$$\text{Continuous: } \mathbb{E}[f] = \int p(x)f(x)dx$$

For finite N samples we can approximate:

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

Conditional expectation:

$$\mathbb{E}_x[f|y] = \sum_x p(x|y)f(x)$$

Variances and Covariances

Variance measures how much variability there is in $f(x)$ around its mean value $\mathbb{E}[f(x)]$

$$var[f] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

Covariance for two random variables x and y

$$cov[x, y] = \mathbb{E}_{x,y}[xy] - \mathbb{E}[x]\mathbb{E}[y]$$

Bayes Decision Theory

Priors a priori probabilities

what can we tell about probability before seeing the data

sum of all priors is 1

Conditional probabilities $p(x|C_k)$ is **likelihood** for class C_k

where x measures/describes certain properties of input

Posterior probabilities $p(C_k|x)$

probability of class C_k given the measurement vector x

$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)} = \frac{p(x|C_k)p(C_k)}{\sum_i p(x|C_i)p(C_i)}$$

Goal: Minimize the probability of a misclassification.

Decide for C_k when:

$$p(C_k|x) > p(C_j|x) \forall j \neq k$$

$$p(x|C_k)p(C_k) > p(x|C_j)p(C_j) \forall j \neq k$$

Classifying with Loss Functions

Motivation: Decide if it's better to choose wrong or nothing.

In general formalized as a matrix L_{jk}

$$L_{jk} = \text{loss for decision } C_j \text{ if } C_k \text{ is correct selection}$$

Minimizing the Expected Loss

Optimal solution requires knowing which class is correct - *this is unknown*. So we **minimize the expected loss**:

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

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

$$\mathbb{E}[L] = \sum_k L_{kj} p(C_k|x)$$

Probability Density Estimation

How can we estimate (= learn) those probability densities?

In Supervised training case: data and class labels are known. So we can estimate the probability density for each class separately.

The Gaussian (or Normal) Distribution

One-dimensional

- 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

- Mean μ
- Variance Σ

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right\} \quad (3)$$

Properties

Central Limit Theorem “The distribution of the sum of N i.i.d. random variables becomes increasingly Gaussian as N grows.”

[There was some more stuff in slides but it was boring math]

The marginals of a Gaussian are again Gaussians

Parametric Methods

Given some data X with parameters $\theta = (\mu, \sigma)$. What is the probability that X has been generated from a probability density with parameters θ .

$$L(\theta) = p(X|\theta)$$

Maximum Likelihood Approach

Single data point:

$$p(x_n|\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x_n-\mu)^2}{2\sigma^2}\right\}$$

Assumption all data points are independent:

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

$$\text{Log-Likelihood: } E(\theta) = -\ln L(\theta) = -\sum_{n=1}^N \ln p(x_n|\theta)$$

Goal: Minimize $E(\theta)$.

How to:

1. Take the derivate of $E(\theta)$
2. Set it to zero
3. Quic Mafs

Warning

MLA is *biased* - it underestimates the true variance. I.e., *overfits to the observed data*.

Frequentist vs Bayesian

Frequentist	Bayesian
probabilities are frequencies of random, repeatable events	quantify the uncertainty about certain states or events

Frequentist	Bayesian
fixed, but can be estimated more precisely when more data is available	uncertainty can be revised in the light of new evidence

Non-Parametric Methods

Often the functional form of the distribution is unknown. So we have to estimate probability from data.

Histograms

Partition data into bins with widths Δ_i and count number of observations in each bin.

$$p_i = \frac{n_i}{N\Delta_i}$$

Often: $\Delta_i = \Delta_j \forall i, j$

Advantages	Disadvantages
works in any dimension D	curse of dimensionality
no need to store data after computation	rather brute force

Bin size:

- too large: too much smoothing
- too small: too much noise

Kernel Density Estimation

Parzen Window

Idea: Hypercube of dimension D with width edge length h .

We place a kernel window k at location x and count how many data points fall inside it. Crude solution because the chosen function k creates hard cuts around Hypercubes.

$$\text{Kernel Function: } k(u) = \begin{cases} 1, & |u_i| \leq \frac{1}{2}, i = 1, \dots, D \\ 0, & \text{else} \end{cases}$$

$$K = \sum_{i=1}^N k\left(\frac{x-x_n}{h}\right)$$

$$V = \int k(u) du = h^D$$

Then probability density is:

$$p(x) \simeq \frac{K}{NV} = \frac{1}{Nh^D} \sum_{i=1}^N k\left(\frac{x-x_n}{h}\right)$$

Gaussian Kernel

Similar to Parzen Window, but with Gaussian kernel function k .

$$k(u) = \frac{1}{(2\pi h^2)^{D/2}} \exp\left\{-\frac{u^2}{2h^2}\right\}$$

$$K = \sum_{i=1}^N k(x - x_n)$$

$$V = \int k(u) du = 1$$

Then probability density is:

$$p(x) \simeq \frac{K}{NV} = \frac{1}{N} \sum_{i=1}^N \frac{1}{(2\pi h^2)^{D/2}} \exp\left\{-\frac{\|x-x_n\|^2}{2h^2}\right\}$$

K-Nearest Neighbors

Similar to above but we fix K (the number of neighbors) and calculate V (size of the neighbourhood).

Then: $p(x) \simeq \frac{K}{NV}$

Warning: Strictly speaking, the model produced by K-NN is not a true density model, because the integral over all space diverges.

Bayesian Classification

$$p(C_j|x) = \frac{p(x|C_j)p(C_j)}{p(x)}$$

Summary

- Very General
- Training requires no computation
- Requires storing and computing the entire dataset -> cost linear in the number of data points -> can be saved in implementation

Kernel size K in K-NN?

- Too large: too much smoothing
- Too small: too much noise

Bias-Variance Tradeoff

	Histograms	Kernel methods	K-Nearest Neighbors
	<i>bin size?</i> Δ	<i>kernel size?</i> h	K ?
too large	Δ : too smooth	h : too smooth	K : too smooth
too small	Δ : not smooth enough	h : not smooth enough	K : not smooth enough

Mixture Distribution

Motivation: single parametric distribution is often not sufficient, e.g., for multimodal data.

Mixture of Gaussians

Idea: we take multiple Gaussians and combine them, giving each one a weight that also depends on input.

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

Where:

- $p(x|\theta_j)$ is the probability of measurement x given the j -th mixture component
- $p(j)$ is the prior of component j

Note:

$$\int p(x)dx = 1$$

Maximum Likelihood

$$E = -\ln L(\theta) = -\sum_{n=1}^N \ln p(x_n|\theta)$$

$$\ln p(X|\pi, \mu, \Sigma) = \sum_{n=1}^N \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k) \right\}$$

-> leads to infinite loop, because gaussians depend on gaussians -> It is possible to apply iterative numerical optimization here, but in the following, we will see a simpler method.

K-Means Clustering

Iterative procedure:

1. Pick K random data points as initial cluster centres
2. Assign each data point to the closest cluster centre
3. Compute the new cluster centres as the mean of the assigned data points
4. Repeat until convergence

Guaranteed to converge after finite number of iterations

But:

- local optimum
- depends on the initial cluster centres

Note:

Can be used e.g. for image compression

Advantages	Disadvantages
simple — fast to compute converges to local minimum	how to select k sensitive to initial cluster centres sensitive to outliers only spherical clusters

EM Algorithm

Expectation-maximization algorithm

E-Step	softly assign samples to mixture components
M-Step	update the parameters of the mixture components

Technical Advice

- When implementing EM, we need to take care to avoid singularities in the estimation!
=> Enforce minimum width for the Gaussians
- EM is very sensitive to the initialization - will converge to local minimum of E => Initialize with k-Means to get better results
- Typical procedure
 - Run k-Means M times (e.g. $M = 10-100$)
 - Pick the best result (lowest error J).
 - Use this result to initialize EM
 - * Set μ_j to the corresponding cluster mean from k-Means.
 - * Initialize Σ_j to the sample covariance of the associated data points.

Summary: Gaussian Mixture Models

Properties:

- Very general, can represent any (continuous) distribution.
- Once trained, very fast to evaluate.
- Can be updated online.

Problems:

- Need to apply regularization in order to avoid singularities
- EM for MoG is computationally expensive

- Need to select the number of mixture components K
 - Model selection problem

Applications

Computer Vision

- Model distributions of pixel colors.
- Each pixel is one data point in, e.g., RGB space.
- Learn a MoG to represent the class-conditional densities.
- Use the learned models to classify other pixels.

Background Model for Tracking

Train background MoG for each pixel

- Model “common” appearance variation for each background pixel
- Initialization with an empty scene.
- Update the mixtures over time
 - Adapt to lighting changes, etc.

Anything that cannot be explained by the background model is labelled as foreground

Image Segmentation

User assisted image segmentation

- User marks two regions for foreground and background.
- Learn a MoG model for the color values in each region.
- Use those models to classify all other pixels.
 - Simple segmentation procedure (building block for more complex applications)

Support Vector Machines

Softmax Regression

Multi-class generalization of logistic regression

Assume binary labels $t_n \in \{0, 1\}$.

Softmax generalizes this to K values in 1-of- K notation.

$$y(x; w) = \begin{bmatrix} P(y = 1|x; u) \\ P(y = 2|x; w) \\ \vdots \\ P(y = k|x; W) \end{bmatrix} \quad (4)$$

With *softmax* function: $\frac{\exp(a_k)}{\sum_j \exp(a_j)}$

Logistic Regression

Alternative way of writing the cost function.

$$E(w) = \sum_{n=1}^N \sum_{k=0}^1 \{\mathbb{I}(t_n = k)\} \ln P(y_n = k|x_n; w)$$

Softmax Regression

Generalization to K classes using indicator functions.

$$E(w) = \sum_{n=1}^N \sum_{k=0}^1 \{\mathbb{I}(t_n = k)\} \ln \frac{\exp(w_k^T x)}{\sum_{j=1}^K \exp(w_j^T x)}$$

Optimization

Again, no closed-form solution is available — Resort again to Gradient Descent.

$$\nabla_{w_k} E(w) = - \sum_{n=1}^N [\mathbb{I}(t_n = k) \ln P(y_n = k | x_n; w)]$$

We can now plug this into a standard optimization package.

Support Vector Machines

Motivation

Goal: predict class labels of new observations.

But: as training progresses we start to overfit to training data.

Popular solution: Cross-validation.

- Split the available data into training and validation sets.
- Estimate the generalization error based on the error on the validation set.
- Choose the model with minimal validation error.

With linearly separable data there are many ways to split the data. The problem is how to choose the best split?

Intuitively, we would like to select the classifier which leaves maximal “safety room” for future data points.

This can be obtained by maximizing the margin between positive and negative data points.

=> The SVM formulates this problem as a convex optimization problem. I.e., we can find optimal solution.

SVM

Consider linearly separable data:

- N training points $\{(x_i, y_i)\}_{i=1}^N$ $x_i \in \mathbb{R}^d$
- Target values $t_i \in \{-1, 1\}$
- Hyperplane separating the two classes: $w^T x + b = 0$

Then Canonical representation of the decision hyperplane:

$$t_n(w^T x_n + b) \geq 1 \forall n$$

Optimization problem

Find the hyperplane satisfying: $\arg_{w,b} \min \frac{1}{2} \|w\|^2$

under constraints: $t_n(w^T x_n + b) \geq 1 \forall n$

- Quadratic programming problem with linear constraints.
- Can be formulated using Lagrange multipliers (see slides).

Lagrangian Formulation

$$L_p = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N a_n \{t_n y(x_n) - 1\}$$

Under conditions:

- $a_n \geq 0$
- $t_n y(x_n) - 1 \geq 0$
- $a_n \{t_n y(x_n) - 1\} = 0$

Solution

Computed as a linear combination of the training examples: $w = \sum_{n=1}^N a_n t_n x_n$.

Because of the KKT conditions, the following must also hold: $a_n(t_n(w^T x_n + b) - 1) = 0$.

This implies that $a_n > 0$ only for training data points for which $t_n(w^T x_n + b) - 1 = 0$.

=> Only some of the data points actually influence the decision boundary!

To define the decision boundary, we still need to know b :

$$b = \frac{1}{N_S} \sum_{n \in S} (t_n - \sum_{m \in S} a_m t_m x_m^T x_n)$$

Discussion

Linear SVM

- Linear classifier
- SVMs have a “guaranteed” generalization capability.
- Formulation as convex optimization problem.

=> Globally optimal solution!

Primal form formulation

- Solution to quadratic prog. problem in M variables is in $\mathcal{O}(M^3)$.
- Here: D variables -> $\mathcal{O}(D^3)$.
- Problem: scaling with high-dim. data (“curse of dimensionality”)

Dual form formulation

Does something that I don't understand.

At the end we should maximize:

$$L_d(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m (x_m^T x_n)$$

Under conditions: $a_n \geq 0 \forall n$ and $\sum_{n=1}^N a_n t_n = 0$.

Then the hyperplane is given by the N_S support vectors:

$$w = \sum_{n=1}^S a_n t_n x_n$$

Soft-margin classification

Solution above works only if data linearly separable. But we can add some tolerance to this division to make it work even if outliers present.

We add “*slack variable*” ξ to the formulation:

$$w^T x_n + b \geq 1 - \xi \text{ for } t_n = 1 \quad w^T x_n + b \leq -1 + \xi \text{ for } t_n = -1$$

where $\xi_n \geq 0 \forall n$.

Interpretation

- $\xi_n = 0$: point on correct side of the margin
- $\xi_n = |t_n - y(x_n)|$: otherwise
 - $\xi_n > 1$: misclassified point

Note:

We do not have to set the slack variables ourselves! They are jointly optimized together with w .

New Dual Formulation

$$L_d(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m (x_m^T x_n)$$

Under conditions: $a_n \geq 0 \forall n$ and $\sum_{n=1}^N a_n t_n = 0$.

Nonlinear Support Vector Machines

Not everything can be linearly separated, so we need nonlinear classifiers.

General idea:

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

Nonlinear transformation ϕ of the data points x_n :

$$x \in \mathbb{R}^D \quad \phi : \mathbb{R}^D \rightarrow \mathcal{H}$$

Hyperplane in higher dimensional space \mathcal{H} :

$$w^T \phi(x) + b = 0$$

Problem with High-dim. Basis Functions

In order to apply the SVM, we need to evaluate the function: $y(x) = w^T \phi(x) + b$.

Using hyperplane $w = \sum_{n=1}^N a_n t_n \phi(x_n)$.

Which leads to Problems in high dimensional feature space.

Solution:

We can replace dot product $\phi(x)^T \phi(x)$ by a kernel function: $k(x, y)$.

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

The kernel function implicitly maps the data to the higher dimensional space (without having to compute $\phi(x)$ explicitly)!

But it only works for specific kernel functions.

“Every positive definite symmetric function is a kernel.” Mercer’s theorem (modernized version)

(positive definite = all eigenvalues are > 0)

Example kernels:

- Polynomial kernel: $k(x, y) = (x^T y + 1)^p$
- Radial Basis Function kernel (e.g. Gaussian): $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$

New New Dual Formulation

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

Under conditions: $\sum_{n=1}^N a_n t_n = 0$ and $\sum_{n=1}^N a_n = C$.

Summary

Properties

- in practice work very well
- among the best performers for a number of classification tasks ranging from text to genomic data
- can be applied to complex data types by designing kernel functions for such data
- The kernel trick has been used for a wide variety of applications. It can be applied wherever dot products are in use

Limitaitons

- How to select the right kernel?
 - Best practice guidelines are available for many applications
- How to select the kernel parameters?
 - (Massive) cross-validation.
 - Usually, several parameters are optimized together in a grid search.
- Solving the quadratic programming problem
 - Standard QP solvers do not perform too well on SVM task.
 - Dedicated methods have been developed for this, e.g. SMO.
- Speed of evaluation
 - Evaluating $y(x)$ scales linearly in the number of SVs.
 - Too expensive if we have a large number of support vectors. -> There are techniques to reduce the effective SV set.
- Training for very large datasets (millions of data points)
 - Stochastic gradient descent and other approximations can be used

Error Function

SVMs result in so-called *hinge error*. Where the error is minimized and correct classification is constant.

This leads to:

- sparsity - Zero error for points outside the margin
- robustness - Linear penalty for misclassified points

Applications

Text Classification

Problem: Classify a document in a number of categories.

Representation:

- “Bag-of-words” approach
- Histogram of word counts (on learned dictionary)

Usage:

- spam filters
- ocr - optical character recognition
- object detection

Adaboost

Ensembles of Classifiers

Idea

- Assume K classifiers.
- They are independent.
- Each has error probability $p < 0.5$ on training data.

Then: Majority vote of all classifiers should have a lower error than each individual classifier

Constructing Ensembles

How do we get different classifiers?

Simplest Case

1. Subsample the training data

Reuse the same training algorithm several times on different subsets of the training data.

2. Train same classifier on different data.

=> Well-suited for “unstable” learning algorithms.

Where:

Unstable	Stable
Decision trees	Nearest neighbor
neural networks	linear regression
rule learning algorithms	SVMs

Unstable := small differences in training data can produce very different classifiers

Bagging

Bagging := “Bootstrap aggregation”

1. In each run of the training algorithm, randomly select M samples with replacement from the full set of N training data points.
2. If $M = N$, then on average, 63.2% of the training points will be represented. The rest are duplicates.

Injecting randomness

- Many (iterative) learning algorithms need a random initialization (e.g. k-means, EM)
- Perform multiple runs of the learning algorithm with different random initializations.

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: $p(X) = \sum_{h=1}^H p(X|h)p(h)$

Interpretation

- Just one model is responsible for generating the entire data set.
- The probability distribution over h just reflects our uncertainty which model that is.
- As the size of the data set increases, this uncertainty reduces, and $p(X|h)$ becomes focused on just one of the models.

Model Combination

- (e.g., Mixtures of Gaussians)
- Different data points generated by different model components.
- Uncertainty is about which component created which data point.

-> One latent variable z_n , for each data point: $p(X) = \prod_{n=1}^N p(x_n) = \prod_{n=1}^N \sum_{z_n} p(x_n, z_n)$

Model Averaging: Expected Error

Some math I don't understand. See slides 19-21 in VL 12.

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!

In reality the errors are **not** uncorrelated - usually highly correlated.

Adaboost

Idea

- Iteratively select an ensemble of component classifiers
- After each iteration, 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

$h_m(x)$ “weak” / base classifier

$H(x)$ “strong” / final classifier

Adaboost

Construct a strong classifier as a thresholded linear combination of the weighted weak classifiers:

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

Minimizing Exponential Error

Exponential error function: $E = \sum_{n=1}^N \exp\{-t_n f_m(x_n)\}$

where $f_m(x)$ is a classifier defined as a linear combination of base classifiers $h_l(x)$:

$$f_m(x) = \frac{1}{2} \sum_{l=1}^m \alpha_l h_l(x)$$

Goal:

Minimize E with respect to both the weighting coefficients α_l and the parameters of the base classifiers $h(x)$.

Here some step by step math that minimizes this.

Final Algorithm

1. Initialization: $w_n^{(1)} = \frac{1}{N}$ for $n = 1, \dots, N$.
2. For $m = 1, \dots, M$ iterations
 - a) Train a new weak classifier $h_m(x)$ using the current weighting coefficients $W(m)$ by minimizing the weighted error function:
$$J_m = \sum_{n=1}^N w_n^{(m)} I(h_m(x) \neq t_n)$$
 - b) Estimate the weighted error of this classifier on X :
TODO
 - c) TODO
 - d) TODO

Summary

Properties

- Simple combination of multiple classifiers.
- Easy to implement.
- Can be used with many different types of classifiers.
 - None of them needs to be too good on its own.
 - In fact, they only have to be slightly better than chance.
- Commonly used in many areas.
- Empirically good generalization capabilities.

Limitations

- Original AdaBoost sensitive to misclassified training data points.
 - Because of exponential error function.
 - Improvement by GentleBoost
- Single-class classifier
 - Multiclass extensions available

Deep Learning / Neural Networks

A Brief History of Neural Networks

1957 Rosenblatt invents the Perceptron

- And a cool learning algorithm: “Perceptron Learning”
- Hardware implementation “Mark | Perceptron” for 20x20 pixel image analysis

1969 Minsky & Papert

- Showed that (single-layer) Perceptrons cannot solve all problems.
- This was misunderstood by many that they were worthless.

1980s Resurgence of Neural Networks

- Some notable successes with multi layer perceptrons.
- Backpropagation learning algorithm
- But they are hard to train, tend to overfit, and have unintuitive parameters.
- So, the excitement fades again...

1995+ Interest shifts to other learning methods

- Notably Support Vector Machines
- Machine Learning becomes a discipline of its own.
- The general public and the press still love Neural Networks.

2005+ Gradual progress

- Better understanding how to successfully train deep networks
- Availability of large datasets and powerful GPUs
- Still largely under the radar for many disciplines applying ML

2012 Breakthrough results

- ImageNet Large Scale Visual Recognition Challenge
- A ConvNet halves the error rate of dedicated vision approaches.
- Deep Learning is widely adopted.

Perceptrons

Standard Perceptron

Construction: Input layer -> Weights -> Output layer

Input Layer := Hand-designed features based on common sense

Outputs:

- linear: $y(x) = w^T x + w_0$
- logistic: $y(x) = \sigma(w^T x + w_0)$

Learning := Determining the weights w

Multi-Class Networks

Construction: Input layer -> Weights -> Output layer

One output node per class

Outputs:

- linear: $y_k(x) = \sum_{i=0}^d W_{ki}x_i$
- logistic: $y_k(x) = \sigma(\sum_{i=0}^d W_{ki}x_i)$

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

Non-Linear Basis Functions

Construction: Input layer -> Mapping (fixed) -> Feature Layer -> Weights -> Output layer

Outputs:

- linear: $y_k(x) = \sum_{i=0}^d W_{ki}\phi(x_i)$
- logistic: $y_k(x) = \sigma(\sum_{i=0}^d W_{ki}\phi(x_i))$

Notes:

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

Perceptron Learning

- Very simple algorithm
- 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.
- This is guaranteed to converge to a correct solution if such a solution exists.

In this process we can use many loss functions:

- L2 loss
- L1 loss
- (Binary) cross-entropy loss
- Hinge loss
- Softmax cross-entropy loss

Regularization

In addition, we can apply regularizers.

$$E(w) = \sum_n L(t_n; y(x_n; w)) + \lambda \|w\|^2$$

- This is known as weight decay in Neural Networks.
- We can also apply other regularizers, e.g. L1 => sparsity
- Since Neural Networks often have many parameters, regularization becomes very important in practice.

Limitations of Perceptrons

What makes the task difficult?

Perceptrons with fixed, hand-coded input features can model any separable function perfectly... given the right input features.

Once the hand-coded features have been determined, there are very strong limitations on what a perceptron can learn. -> Classic example: XOR function.

A linear classifier cannot solve certain problems

However, with a non-linear classifier based on the right kind of features, the problem becomes solvable.

Multi-Layer Perceptrons

Construction: Input layer -> Mapping (learned) -> Hidden layer -> Output layer

Output:

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

With activation functions $g^{(k)}$.

E.g., $g^{(1)}(a) = a$, $g^{(2)}(a) = \sigma(a)$

The hidden layer can have an arbitrary number of nodes

Universal approximators A 2-layer network (1 hidden layer) can approximate any continuous function of a compact domain arbitrarily well!

Learning with Hidden Units

Networks without hidden units are very limited in what they can learn

- More layers of linear units do not help = still linear
- Fixed output non-linearities are not enough.

We need multiple layers of *adaptive* non-linear hidden units. But how can we train such nets?

- Need an efficient way of adapting all weights, not just the last layer.
- Learning the weights to the hidden units = learning features
- This is difficult, because nobody tells us what the hidden units should do.

-> Main challenge in deep learning.

Gradient Descent

Two main steps:

1. Computing the gradients for each weight
2. Adjusting the weights in the direction of the gradient

Set up Error function $E(W) = \sum_n L(t_n; y(x_n; W)) + \lambda \Omega(W)$ with a loss L and regularizer Ω .

Then update each weight $W_{ij}^{(k)}$ in the direction of the gradient: $\frac{\partial E(W)}{\partial W_{ij}^{(k)}}$

Obtaining the Gradients

Approach 1: Naive Analytical Differentiation

Compute the gradients for each variable analytically.

Multi-dimensional case: Total derivative -> Need to sum over all paths that lead to the target variable x .

What is the problem when doing this?

- With increasing depth, there will be exponentially many paths!
- Infeasible to compute this way.

Approach 2: Numerical Differentiation

Given the current state $W^{(\tau)}$, we can evaluate $E(W^{(\tau)})$.

Idea: Make small changes to $W^{(\tau)}$ and accept those that improve $E(W^{(\tau)})$.

=> Horribly inefficient! Need several forward passes for each weight. Each forward pass is one run over the entire dataset!

Approach 3: Incremental Analytical Differentiation (Backpropagation)

Idea:

- Compute the gradients layer by layer.
- Each layer below builds upon the results of the layer above.

The gradient is propagated backwards through the layers.

Backpropagation Algorithm

1. Convert the discrepancy between each output and its target value into an error derivative.
2. Compute error derivatives in each hidden layer from error derivatives in the layer above.
3. Use error derivatives w.r.t. activities to get error derivatives w.r.t. the incoming weights.

Notation:

- $y_j^{(k)}$ Output of layer k Connections: $z_j^{(k)} = \sum_i w_{ij}^{(k-1)} y_i^{(k-1)}$
- $z_j^{(k)}$ Input of layer k $y_j^{(k)} = g(z_j^{(k)})$

Efficient propagation scheme:

- $y_i^{(k-1)}$ is already known from forward pass! (Dynamic Programming)

-> Propagate back the gradient from layer k and multiply with $y_i^{(k-1)}$

Here the algorithms in pseudocode

Analysis: Backpropagation

Backpropagation is the key to make deep NNs tractable.

However the Backprop algorithm given here is specific to MLPs.

- It does not work with more complex architectures, e.g. skip connections or recurrent networks!
- Whenever a new connection function induces a different functional form of the chain rule, you have to derive a new Backprop algorithm for it.

Learning Multi-layer Networks

Computational Graphs

We can think of mathematical expressions as graphs. I.e. we can divide every equation into a set of simple sub-equations.

From: $e = (a + b) * (b + 1)$

To:

$c = a + b$ $d = b + 1$ $e = c * d$

Which makes it easier to take derivatives of those equations.

Problem: Combinatorial explosion

Solution:

Efficient algorithms for computing the sum.

Instead of summing over all of the paths explicitly, compute the sum more efficiently by merging paths back together at every node.

Approach 4: Automatic Differentiation

- 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.

Modular Implementation

Solution in many current Deep Learning libraries

- Provide a limited form of automatic differentiation
- Restricted to “programs” composed of “modules” with predefined set of operations.

Module is defined by two main functions:

1. $y = \text{module.fprop}(x)$ computes outputs y given inputs x , where x, y and intermediate results stored in the module
2. $\frac{\partial E}{\partial x} = \text{module.bprop}(\frac{\partial E}{\partial y})$ computing the gradient $\partial E / \partial x$ of a scalar cost w.r.t. the inputs x given the gradient $\partial E / \partial y$ w.r.t. the outputs y

Here some stuff about Implementation Problems.

Gradient Descent

Two main steps:

1. Computing the gradients for each weight
2. Adjusting the weights in the direction of the gradient

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \frac{\partial E(w)}{\partial w_{kj}} \Big|_{w^{(\tau)}}$$

Stochastic vs. Batch Learning

Batch learning

Process the full dataset at once to compute the gradient.

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \frac{\partial E(w)}{\partial w_{kj}} \Big|_{w^{(\tau)}}$$

Advantages:

- Conditions of convergence are well understood.
- Many acceleration techniques (e.g., conjugate gradients) only operate in batch learning.
- Theoretical analysis of the weight dynamics and convergence rates are simpler.

Stochastic learning

- Choose a single example from the training set.
- Compute the gradient only based on this example
- This estimate will generally be noisy, which has some advantages.

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \frac{\partial E_n(w)}{\partial w_{kj}} \Big|_{w^{(\tau)}}$$

Advantages:

- Usually much faster than batch learning.
- Often results in better solutions.
- Can be used for tracking changes.

Minibatches

Middle ground between batch and stochastic learning.

Idea:

- Process only a small batch of training examples together
- Start with a small batch size & increase it as training proceeds.

Advantages

- Gradients will be more stable than for stochastic gradient descent, but still faster to compute than with batch learning.
- Take advantage of redundancies in the training set.
- Matrix operations are more efficient than vector operations.

Caveat

Error function should be normalized by the minibatch size, s.t. we can keep the same learning rate between minibatches

$$E(W) = \frac{1}{N} \sum_n L(t_n, y(x_n; W)) + \frac{\lambda}{N} \Omega(W)$$

Choosing the Right Learning Rate

Considering a simple 1D example. If E is quadratic, the optimal learning rate is given by the inverse of the Hessian.

What happens if we exceed this learning rate?

Instead of getting closer to minimum, we start moving away from it fast.

Momentum

Batch Learning

- Simplest case: steepest descent on the error surface.
- Updates perpendicular to contour lines

Stochastic Learning

- Simplest case: zig-zag around the direction of steepest descent.
- Updates perpendicular to constraints from training examples.

If the inputs are correlated, the ellipse will be elongated and direction of steepest descent is almost perpendicular to the direction towards the minimum!

The Momentum Method

Idea

Instead of using the gradient to change the position of the weight “particle”, use it to change the velocity.

Intuition

- Example: Ball rolling on the error surface
- It starts off by following the error surface, but once it has accumulated momentum, it no longer does steepest descent.

Effect

- Dampen oscillations in directions of high curvature by combining gradients with opposite signs.
- Build up speed in directions with a gentle but consistent gradient.

Behavior

- If the error surface is a tilted plane, the ball reaches a terminal velocity
 - If the momentum α is close to 1, this is much faster than simple gradient descent.

- At the beginning of learning, there may be very large gradients.
 - Use a small momentum initially (e.g., $\alpha = 0.5$).
 - Once the large gradients have disappeared and the weights are stuck in a ravine, the momentum can be smoothly raised to its final value (e.g., $\alpha = 0.90$ or even $\alpha = 0.99$).
 - This allows us to learn at a rate that would cause divergent oscillations without the momentum.

Separate, Adaptive Learning Rates

Problem

- In multilayer nets, the appropriate learning rates can vary widely between weights.
- The magnitudes of the gradients are often very different for the different layers, especially if the initial weights are small.
 - Gradients can get very small in the early layers of deep nets.
- The fan-in of a unit determines the size of the “overshoot” effect when changing multiple weights simultaneously to correct the same error.
 - The fan-in often varies widely between layers

Solution

- Use a global learning rate, multiplied by a local gain per weight (determined empirically)

Better Adaptation: RMSProp

Motivation

- The magnitude of the gradient can be very different for different weights and can change during learning.
- This makes it hard to choose a single global learning rate.
- For batch learning, we can deal with this by only using the sign of the gradient, but we need to generalize this for minibatches.

Idea

Divide the gradient by a running average of its recent magnitude

$$MeanSq(w_{ij}, t) = 0.9MeanSq(w_{ij}, t-1) + 0.1\left(\frac{\partial E}{\partial w_{ij}}(t)\right)^2$$

Divide the gradient by $\sqrt{MeanSq(w_{ij}, t)}$.

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.

Summary

Deep multi-layer networks are very powerful.

But training them is hard!

- Complex, non-convex learning problem
- Local optimization with stochastic gradient descent

Main issue: getting good gradient updates for the early layers of the network

- Many seemingly small details matter!
- Weight initialization, normalization, data augmentation, choice of nonlinearities, choice of learning rate, choice of optimizer, . . .

Optimization

Tricks of the Trade

Shuffling the Examples

Idea: Network learns faster from most unexpected sample. It is advisable to choose a sample at each iteration that is most unfamiliar to the system.

I.e., **do not present all samples of class A, then all of class B.**

A large relative error indicates that an input has not been learned by the network yet, so it contains a lot of information. It can make sense to present such inputs more frequently. **But:** be careful, this can be disastrous when the data are outliers.

When working with stochastic gradient descent or minibatches, make use of shuffling.

Data augmentation

Idea: Augment original data with synthetic variations to reduce overfitting.

Examples of augmentation:

- Cropping
- Zooming
- Flipping
- Color PCA (color space transformation)

Effect:

- Much larger training set
- Robustness against expected variations

Normalization

Motivation:

Consider the Gradient Descent update steps. When all of the components of the input vector y_i are positive, all of the updates of weights that feed into a node will be of the same sign.

Leads to:

- Weights can only all increase or decrease together.
- Slow convergence

Convergence is the 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.

Nonlinearities

Sigmoid	$g(a) = \sigma(a) = \frac{1}{1+\exp\{-a\}}$
Hyperbolic tangent	$g(a) = \tanh(a) = 2\sigma(2a) - 1$
Softmax	$g(a) = \frac{\exp\{-a_i\}}{\sum_j \exp\{-a_j\}}$

Normalization is also important for intermediate layers: Symmetric sigmoids, such as tanh, often converge faster than the standard logistic sigmoid.

Recommended sigmoid:

$$f(x) = 1.7159 \tanh(\frac{2}{3}x)$$

When used with transformed inputs, the variance of the outputs will be close to 1.

Usage

Output nodes - Typically, a sigmoid or tanh function is used here:

- Sigmoid for nice probabilistic interpretation (range [0,1]).
- tanh for regression tasks

Internal nodes:

- Historically, tanh was most often used.
- tanh is better than sigmoid for internal nodes, since it is already centered.
- Internally, tanh is often implemented as piecewise linear function (similar to hard tanh and maxout).
- More recently: ReLU often used for classification tasks.

Effects of sigmoid/tanh function:

- Linear behavior around 0
- Saturation for large inputs

Extension: ReLU

Improvement for learning deep models. That uses Rectified Linear Units (ReLU) instead of sigmoid: $g(a) = \max\{0, a\}$.

Effect gradient is propagated with a constant factor: $\frac{\partial g(a)}{\partial a} = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$

Advantages:

- Much easier to propagate gradients through deep networks.
- We do not need to store the ReLU output separately
 - Reduction of the required memory by half compared to tanh!

Disadvantages / Limitations:

- A certain fraction of units will remain “stuck at zero”.
 - If the initial weights are chosen such that the ReLU output is 0 for the entire training set, the unit will never pass through a gradient to change those weights.
- ReLU has an *offset bias*, since its outputs will always be positive

Further Extensions

ReLU	$g(a) = \max\{0, a\}$
Leaky ReLU	$g(a) = \max\{\beta a, a\}$
	Avoids stuck-at-zero units Weaker offset bias
ELU	$g(a) = \begin{cases} a & x < 0 \\ e^a - 1 & x \geq 0 \end{cases}$
	No offset bias anymore BUT: need to store activations

Initialization

Initializing the Weights

Motivation:

- The starting values of the weights can have a significant effect on the training process.
- Weights should be chosen randomly, but in a way that the sigmoid is primarily activated in its linear region.

Assuming that:

- The training set has been normalized
- The recommended sigmoid $f(x) = 1.7159 \tanh(\frac{2}{3})$ is used

the initial weights should be randomly drawn from a distribution (e.g., uniform or Normal) with mean zero and variance:

$$\sigma_w^2 = \frac{1}{n_{in}}$$

where n_{in} is the fan-in (#connections into the node).

Glorot Initialization

In 2010, Xavier Glorot published an analysis of what went wrong in the initialization and derived a more general method for automatic initialization.

This new initialization massively improved results and made direct learning of deep networks possible overnight.

Analysis

Some maths

Advanced techniques

Batch Normalization

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
- **Complication:** centering + normalization also needs to be done at test time, but minibatches are no longer available at that point.
 - Learn the normalization parameters to compensate for the expected bias of the previous layer (usually a simple moving average)

Effect

- Much improved convergence (but parameter values are important!)
- Widely used in practice

Dropout

Idea:

- Randomly switch off units during training (a form of regularization).
- Change network architecture for each minibatch, 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 during training.

=> Greatly improved performance