

Algorithmic Foundations of Datascience

A summary of the Lecture

Felix Friedberger

April 2019 - September 2019

Abstract

This summary is based on the lecture 'Algorithmic Foundations of Datascience' held by Prof. Martin Grohe in the summer term of 2019 at RWTH Aachen University.

Since this summary is made by student(s), there are definitely mistakes in the script. Should you find one feel free to contact me.

Contents

Abstract	i
1 Machine Learning Basics	3
1.1 Definitions	3
1.1.1 Domains	3
1.2 Forms of Learning	3
1.3 Hypothesis Space	4
1.4 Validation	4
1.5 Nearest Neighbour Learning	4
1.5.1 Metric	4
1.5.2 Metric Spaces	5
1.5.2.1 Common distance metrics	5
1.5.3 Nearest Neighbour Classifier	5
1.5.4 Learning Decision Trees	5
1.5.4.1 Greedy Algorithm for Building Decision Trees	6
1.6 The Perceptron	6
1.6.1 Scalar product	6
1.6.2 Euclidean norm	6
1.6.3 Cauchy-Schwarz Inequality	6
1.7 Linear Classification	6
1.7.1 Linear Classification	6

1 Machine Learning Basics

1.1 Definitions

All the following definitions are specific to 'Machine Learning' and might change throughout the summary.

Data is viewed as a collection of data items (or dataset). A **data item** is represented by a so-called feature vector. The **feature vector** basically a list of features and attributes in form of a vector

1.1.1 Domains

The **Domain** \mathbb{D} is a set containing the range of all possible values for a single feature. The **Instance Space** \mathbb{X} is the Cartesian product of the domains of all features in the feature vector. The number of features is the dimension of the instance space (or the number of entries in the feature vector). ◀ \mathbb{D}
◀ \mathbb{X}

Example: The domain space can be something like $\{true, false\}$ or \mathbb{N} , aso. The instance space of the domain spaces is then (informally) $\mathbb{D}_1 \times \mathbb{D}_2 \times \dots \times \mathbb{D}_l$

Assuming n is the number of data items, and l the number of features the whole dataset can be viewed as an $(n \times l)$ -matrix over the instance space.

1.2 Forms of Learning

Unsupervised learning tries to detect patterns in data with, hence the name, no feedback, supplied the Learning algorithm. This technique is mostly used for clustering.

Reinforcement Learning tries to find actions that maximise reward. The feedback to the learning algorithm is provided in form of a reward (or punishment)

In **Supervised learning** one tries to approximate a function by training on input-output pairs. **Classification** has a finite-valued function and tries to predict values for future inputs. **Regression** has a numerical function and tries to predict the expected values for future inputs. Supervised learning can be differentiated into two different set-ups. First, **Batch Learning**, in which all examples are supplied at once and the learner has to come up with a hypothesis. Where as in **Online Learning**, examples are given 'on-the-fly' once at a time and the initial hypothesis gets improved over time.

Semi-supervised Learning is similar to supervised learning, but takes only a few and possibly faulty examples and tries to make the best of the possibly small and faulty dataset.

1.3 Hypothesis Space

\mathcal{H} ► In Supervised learning, the unknown target function h is chosen from the **hypothesis space** \mathcal{H} , hence $h \in \mathcal{H}$. The Hypothesis space contains all linear functions, all polynomials or all polynomials of a pre-scribed degree, and all functions that can be described by a decision tree.

A learning problem is **Realisable** if the target function h is in the Hypothesis Space.

Example: Assuming a target function $h(x) := a \cdot \sin(x) + bx + c$. Then, if \mathcal{H} only contains polynomials, the learning problem is **not realizable**. But, assuming \mathcal{H} contains linear functions or polynomials in x and $\sin(x)$, the learning problem is **realisable**.

1.4 Validation

The goal of the learning algorithm is to produce a hypothesis that generalizes well, that is, approximates the target function well on all data points (and not only those in the training set). To evaluate how well a hypothesis generalizes we can evaluate it against some test set. A **test set** is generated by splitting the examples into a test set and training set.

The Empirical Observation that simpler hypothesis tend to generalise better is picked up by **Occam's Razor** which states, that the simplest hypothesis which is consistent with the data should be chosen.

Overfitting is the phenomenon, which occurs when the learner tries to match the training examples as exactly as possible which leads to complex hypotheses that generalize badly. Thus, to avoid overfitting, it is often better to choose a simple hypothesis even if it doesn't match the data exactly.

1.5 Nearest Neighbour Learning

The underlying idea of this simple learning algorithm is to predict the value of a function at point x by looking at known values of points close to the given one and assume the value of x is similar.

1.5.1 Metric

d ► A **metric** d on \mathbb{X} is a function $d : \mathbb{X}^2 \longrightarrow \mathbb{R}$ such that for all $x, y, z \in \mathbb{X}$

Nonnegativity	Symmetry	Triangle Inequality
$d(x, y) \geq 0$ and $d(x, y) = 0 \iff x = y;$	$d(x, y) = d(y, x)$	$d(x, z) \leq d(x, y) + d(y, z)$

1.5.2 Metric Spaces

Suppose \mathbb{X} is the instance space, and d be some metric on \mathbb{X} . Then we have a **Metric Space** (\mathbb{X}, d) .

◀ (\mathbb{X}, d)

1.5.2.1 Common distance metrics

The distance between points x and y is denoted as $d(x, y)$.

◀ $d(x, y)$

Euclidean distance with $\mathbb{X} = \mathbb{R}^l$:

$$d(x, y) = \sqrt{\sum_{i=1}^l (x_i - y_i)^2} \quad (1.1)$$

Hamming distance with $\mathbb{X} = \mathbb{D}_1 \times \mathbb{D}_2 \times \dots \times \mathbb{D}_l$

$$d(x, y) = |\{i \mid x_i \neq y_i\}| \quad (1.3)$$

Manhattan distance with $\mathbb{X} = \mathbb{R}^l$:

$$d(x, y) = \sum_{i=1}^l |x_i - y_i| \quad (1.2)$$

1.5.3 Nearest Neighbour Classifier

Problem description: 'Learn' a function f , that associates a class $f(x) \in \mathbb{Y}$ with every $x \in \mathbb{X}$. Basically: $f : \mathbb{X} \rightarrow \mathbb{Y}, x \mapsto y$. Function f is provided with a dataset to learn from: $(x_1, y_1), \dots, (x_m, y_m)$ with m entries where each entry corresponds to $(x_i, f(x_i))$. x_i is called data item and y_i is called class x_i .

k-Nearest Neighbour Classifier uses $x \in \mathbb{X}$ as input.

1. Find the k nearest neighbours x_{i_1}, \dots, x_{i_k} of x in $\{x_1, \dots, x_m\}$
2. Take a 'majority vote' of the class that appears the most among y_{i_1}, \dots, y_{i_k}

Choosing a appropriate parameter k . When choosing $k = 1$ the classifier tends to overhit. While with larger k one might start to oversimplify. Obviously, the right k is difficult to compute and depends heavily on the application.

1.5.4 Learning Decision Trees

A decision tree is a rooted tree with labelled nodes and edges

- Every internal node of the tree is labelled by an (input) feature
- Every edge is labeled by a value or range of values for the feature labelling its source node.
- Every leaf is labeled with an output value

Semantics: The function value for an input vector \mathbf{x} is the value at the leaf of the unique path in the tree whose edges are labelled by the feature values in x .

Note: Computing a smallest decision tree for a given sets of examples is NP-Complete

1.5.4.1 Greedy Algorithm for Building Decision Trees

TODO: Greedy Algorithm for Building Decision Trees

1.6 The Perceptron

The perceptron algorithm attempts to find a linear classifier for a Boolean classification problem.

1.6.1 Scalar product

$\langle \mathbf{x}, \mathbf{y} \rangle$ ► Two vectors, $\mathbf{x} = (x_1, \dots, x_l)^T$, $\mathbf{y} = (y_1, \dots, y_l)^T \in \mathbb{R}^l$

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^l x_i y_i$$

1.6.2 Euclidean norm

$\|\mathbf{x}\|$ ► Vector $\mathbf{x} \in \mathbb{R}^l$

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\sum_{i=1}^l x_i^2}$$

1.6.3 Cauchy-Schwarz Inequality

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\| \quad (1.4)$$

1.6.4 Hyperplanes

1.6.5 Halfspaces

1.6.6 Linear Classification

$$S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)) \in \mathbb{R}^\ell \times \{+1, -1\} \quad (1.5)$$

$$h(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle - b) = \begin{cases} +1 & \text{if } \langle \mathbf{w}, \mathbf{x} \rangle - b > 0 \\ 0 & \text{if } \langle \mathbf{w}, \mathbf{x} \rangle - b = 0 \\ -1 & \text{if } \langle \mathbf{w}, \mathbf{x} \rangle - b < 0 \end{cases} \quad (1.6)$$