

All of Machine Learning

A summary under eternal construction

Moritz Gück

github.com/MoritzGuck/All_of_ML-under_construction

Last changes: March 15, 2020

Abstract

This is a reference for machine learning approaches and methods. The topics range from basic statistics to complex machine learning models and explanation methods. For each method and model, I have provided the underlying formulas (objective functions, prediction functions, etc.) as well as code snippets from the respective python libraries. I made this reference to quickly look up things I have studied already. I published it to give data scientists a catalogue to find methods for their problem, refresh their knowledge and give references for further reading. If you find errors or unclear explanations in this text, please file an issue under: github.com/MoritzGuck/All_of_ML-under_construction

Contents

1	Probability Theory	2
1.1	Probability distributions	2
2	Classification Methods	3
2.1	Linear Classifiers	3
2.2	Support Vector Classifier (SVC)	3
3	Clustering Methods	4
3.1	KNN-Clustering	4
4	Classification methods	5
4.1	Support Vector Machines (SVM)	5
5	Kernel methods	6
5.1	Introduction to Kernels	6
5.2	Kernel SVM	6
6	Interpretability Methods	7
6.1	Local interpretability methods	7

1 Probability Theory

A probability is a measure of how frequent or likely an event will take place.

Probability Space The probability space is a triplet space containing a sample/outcome space Ω (containing all possible atomic events), a collection of events S (containing a subset of Ω to which we want to assign probabilities) and the mapping P between Ω and S .

Axioms of Probability The mapping P must fulfill the axioms of probability:

1. $P(a) \geq 0$
2. $P(\Omega) = 1$
3. $a, b \in S$ and $a \cap b = \{\}$ $\Rightarrow P(a \cup b) = P(a) + P(b)$

Random Variable A random variable is a function that maps points from the sample space Ω to some range (e.g. Real numbers or booleans). They are characterized by their distribution function. E.g. for a dice roll:

$$X(\omega) = \begin{cases} 0, & \text{if } \omega = \text{heads} \\ 1, & \text{if } \omega = \text{tails}. \end{cases}$$

Proposition A Proposition is a conclusion of a statistical inference that can be true or false (e.g. a classification of a datapoint). More formally: A disjunction of events where the logic model holds. An event can be written as a **propositional logic model**:

$A = \text{true}, B = \text{false} \Rightarrow a \wedge \neg b$. Propositions can be continuous, discrete or boolean.

1.1 Probability distributions

Probability distributions assign probabilities to all possible points in Ω (e.g. $P(\text{Weather}) = \langle 0.3, 0.4, 0.2, 0.1 \rangle$, representing Rain, sunshine, clouds and snow). Joint probability distributions give you a probability for each atomic event of the random variables (e.g. $P(\text{weather}, \text{accident})$ gives you a 2×4 matrix.)

2 Classification Methods

Cumulative Distribution Function Classification is the assignment of objects (data points) to categories (classes). It requires a data set (i.e. training set) of points with known class labels. If the class labels are not known you can instead group the data using clustering algorithms (chapter 3).

2.1 Linear Classifiers

Linear classifiers use linear decision boundaries to classify points to a respective class.

2.2 Support Vector Classifier (SVC)

SVCs use hyperplanes to separate data points according to their class label with a maximum margin (M) between the separating hyperplane ($x^T\beta + \beta_0 = 0$) and the points. If points cannot be perfectly separated by the decision boundary, a soft margin SVM is used with a slack variable ξ that punishes points in the margin or on the wrong side of the hyperplane. The optimization problem is given by [?]:

$$\begin{aligned} & \max_{\beta, \beta_0, \xi} M, \\ & \text{subject to } y_i(x_i^T\beta + \beta_0) \geq 1 - \xi_i, \quad \forall i, \\ & \xi_i \geq 0, \quad \sum \xi_i \leq \text{constant}, \quad i = 1, \dots, N, \end{aligned} \tag{1}$$

where β are the coefficients and x are the N data points. The support vectors are the points that determine the orientation of the hyperplane (i.e. the closest points). The classification function is given by:

$$G(x) = \text{sign}[x^T\beta + \beta_0] \tag{2}$$

3 Clustering Methods

3.1 KNN-Clustering

4 Classification methods

4.1 Support Vector Machines (SVM)

5 Kernel methods

Kernel methods help in using linear decision boundaries on non-linearly separable data by morphing the feature space. They can also help in disentangling data for clustering. Kernels can incorporate domain knowledge into your model.

5.1 Introduction to Kernels

Kernels can be seen in two ways: As similarity measures between data-points or transformations of the data-points into a higher dimensional space. For two points x and y a Kernel is given by:

$$K(x, y) = \sum_{i=1}^n h_i(x)h_i(y) = \langle h(x), h(y) \rangle, \quad (3)$$

where $h(x)$ is a transformation-function and $\langle \cdot \rangle$ is the inner product.

5.2 Kernel SVM

6 Interpretability Methods

Complex machine learning algorithms (e.g. NNs) are hard/impossible to interpret. Interpretability methods help with debugging, trust and taking appropriate action on the results.

6.1 Local interpretability methods

Explain, why your model made this/these exact decisions.

6.1.1 Shapley values

Find attributes that determine the deviation of your output from the *expected value*. **How:** Calculate how much each feature pushes the prediction away from the expected value, by shuffling through all combinations of features having the sample-value or expected value respectively.

SHapley Additive exPlanations (SHAP) SHAP calculates Shapley values

6.1.2 Local interpretable model-agnostic explanations (LIME)

! → LIME interpretations are not always consistent.

6.1.3 Example-based explanations

Influence Functions What would happen to the model parameters, if you would up-weight an instance? (model is function of training data.)

Index

Axioms of Probability, 2

Cumulative Distribution Function, 2

Probability distribution, 2

Probability Space, 2

Proposition, 2

Random Variable, 2

SHapley Additive exPlanations (SHAP), 7