

INTRODUCING MACHINE LEARNING

Jan-Philipp Kolb

22 Mai, 2019

WHAT IS MACHINE LEARNING?

- ▶ Machine learning allows the user to feed a computer algorithm an immense amount of data and have the computer analyze and make data-driven recommendations and decisions based on only the input data.
- ▶ If any corrections are identified, the algorithm can incorporate that information to improve its future decision making.

CATEGORIZING MACHINE LEARNING ALGORITHMS...

- ▶ ... is tricky, and there are several approaches;
- ▶ they can be grouped into generative/discriminative, parametric/non-parametric, supervised/unsupervised, and so on.

WHAT IS SUPERVISED LEARNING?

Supervised learning includes tasks for “labeled” data (i.e. you have a target variable).

- ▶ In practice, it's often used as an advanced form of predictive modeling.
- ▶ Each observation must be labeled with a “correct answer.”
- ▶ Only then can you build a predictive model because you must tell the algorithm what's “correct” while training it (hence, “supervising” it).
- ▶ Regression is the task for modeling continuous target variables.
- ▶ Classification is the task for modeling categorical (a.k.a. “class”) target variables.

SUPERVISED VS UNSUPERVISED LEARNING

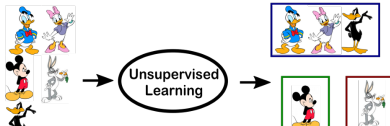
SUPERVISED LEARNING

- Prior knowledge of what output values for samples should be.



UNSUPERVISED LEARNING

- Here the most common tasks are clustering, representation learning, and density estimation - we wish to learn the inherent structure of our data without using explicitly-provided labels.



MACHINE LEARNING - COMPONENTS

- ▶ Feature Extraction + Domain knowledge
- ▶ Feature Selection
- ▶ Choice of Algorithm (Regression or classification, regularization, decision trees, k-Means clustering, ...)
- ▶ Training
- ▶ Choice of Metrics/Evaluation Criteria
- ▶ Testing

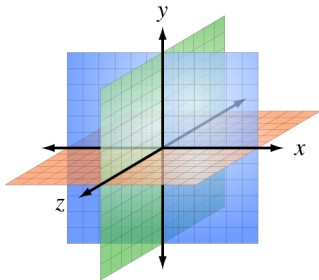
FEATURE SELECTION IN MACHINE LEARNING,...

- ▶ ... is the process of selecting a subset of relevant features (variables, predictors) for use in model construction.

FOUR REASONS FOR FEATURE SELECTION:

- 1.) simplification of models to make them easier to interpret by researchers/users,
- 2.) shorter training times,
- 3.) to avoid the curse of dimensionality,
- 4.) enhanced generalization by reducing overfitting (formally, reduction of variance)

THE CURSE OF DIMENSIONALITY



In machine learning, “dimensionality” simply refers to the number of features (i.e. input variables) in your dataset.

- ▶ When the number of features is very large relative to the number of observations, certain algorithms struggle to train effective models.
- ▶ This is called the “Curse of Dimensionality,” and it’s especially relevant for clustering algorithms that rely on distance calculations.

WHAT ARE THE ADVANTAGES AND DISADVANTAGES OF DECISION TREES?

Advantages: Decision trees are easy to interpret, nonparametric (which means they are robust to outliers), and there are relatively few parameters to tune.

Disadvantages: Decision trees are prone to be overfit.

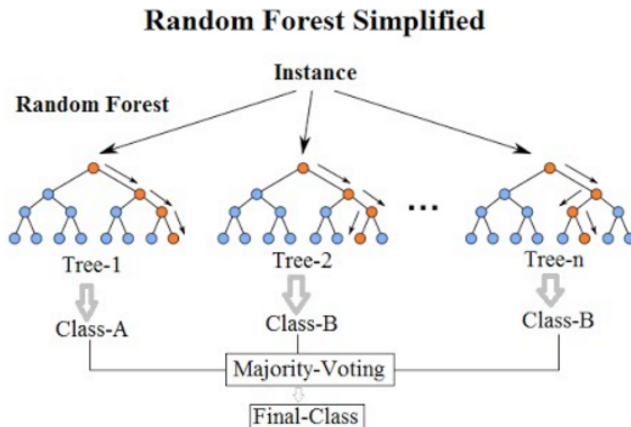
- ▶ This can be addressed by ensemble methods like random forests or boosted trees.

RANDOM FOREST

Random forest aims to reduce the previously mentioned correlation issue by choosing only a subsample of the feature space at each split. Essentially, it aims to make the trees de-correlated and prune the trees by setting a stopping criteria for node splits, which I will cover in more detail later.

RANDOM FOREST

- ▶ Ensemble learning method - multitude of decision trees
- ▶ Random forests correct for decision trees' habit of overfitting to their training set.



ENSEMBLING

Ensembles are machine learning methods for combining predictions from multiple separate models.

BAGGING

attempts to reduce the chance overfitting complex models.

- ▶ It trains a large number of “strong” learners in parallel.
- ▶ A strong learner is a model that’s relatively unconstrained.
- ▶ Bagging then combines all the strong learners together in order to “smooth out” their predictions.

BOOSTING

attempts to improve the predictive flexibility of simple models.

- ▶ It trains a large number of “weak” learners in sequence.
- ▶ A weak learner is a constrained model (limit for max depth of tree).
- ▶ Each one in the sequence focuses on learning from the mistakes of the one before it.
- ▶ Boosting combines all the weak learners into a single strong learner.

BAGGING AND BOOSTING

While bagging and boosting are both ensemble methods, they approach the problem from opposite directions.

Bagging uses complex base models and tries to “smooth out” their predictions, while boosting uses simple base models and tries to “boost” their aggregate complexity.

GRADIENT BOOSTING

Gradient boosting is a machine learning technique for regression and classification problems, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees. It builds the model in a stage-wise fashion like other boosting methods do, and it generalizes them by allowing optimization of an arbitrary differentiable loss function.

The idea of gradient boosting originated in the observation by Leo Breiman that boosting can be interpreted as an optimization algorithm on a suitable cost function.

Breiman, L. (1997). "Arcing The Edge". Technical Report 486. Statistics Department, University of California, Berkeley.

Advantages of gradient boosting

- ▶ Often provides predictive accuracy that cannot be beat.
- ▶ Lots of flexibility - can optimize on different loss functions and provides several hyperparameter tuning options that make the function fit very flexible.
- ▶ No data pre-processing required - often works great with categorical and numerical values as is.
- ▶ Handles missing data - imputation not required.

Disadvantages OF GRADIENT BOOSTING

- ▶ GBMs will continue improving to minimize all errors. This can overemphasize outliers and cause overfitting. Must use cross-validation to neutralize.
- ▶ Computationally expensive - GBMs often require many trees (>1000) which can be time and memory exhaustive.
- ▶ The high flexibility results in many parameters that interact and influence heavily the behavior of the approach (number of iterations, tree depth, regularization parameters, etc.). This requires a large grid search during tuning.
- ▶ Less interpretable although this is easily addressed with various tools (variable importance, partial dependence plots, LIME, etc.).

TWO TYPES OF ERRORS FOR TREE METHODS

BIAS RELATED ERRORS

- ▶ Adaptive boosting
- ▶ Gradient boosting

VARIANCE RELATED ERRORS

- ▶ Bagging
- ▶ Random forest

LINKS

- ▶ Presentations on 'Elements of Neural Networks & Deep Learning'
- ▶ Understanding the Magic of Neural Networks
- ▶ Neural Text Modelling with R package ruimtehol
- ▶ Feature Selection using Genetic Algorithms in R
- ▶ Lecture slides: Real-World Data Science (Fraud Detection, Customer Churn & Predictive Maintenance)
- ▶ Automated Dashboard for Credit Modelling with Decision trees and Random forests in R
- ▶ Looking Back at Google's Research Efforts in 2018
- ▶ Selecting 'special' photos on your phone
- ▶ Open Source AI, ML & Data Science News
- ▶ Google's Machine Learning Crash Course
- ▶ A prelude to machine learning
- ▶ caret webinar by Max Kuhn - on youtube
- ▶ learn-math-for-data-science
- ▶ learn-statistics-for-data-science
- ▶ machine-learning-projects-for-beginners
- ▶ An Introduction to machine learning

ANNEX - PREDICTION VS. CAUSATION IN REGRESSION ANALYSIS

Paul Allison

There are two main uses of multiple regression: prediction and causal analysis. In a prediction study, the goal is to develop a formula for making predictions about the dependent variable, based on the observed values of the independent variables. . . . In a causal analysis, the independent variables are regarded as causes of the dependent variable. The aim of the study is to determine whether a particular independent variable really affects the dependent variable, and to estimate the magnitude of that effect, if any