

Introduction to Machine Learning

ML-Basics: Losses & Risk Minimization

HOW TO EVALUATE MODELS

OVERVIEW

No Free Lunch In machine learning, there's something called the "No Free Lunch" theorem. In a nutshell, it states that no one algorithm works best for every problem, and it's especially relevant for supervised learning (i.e. predictive modeling).

For example, you can't say that neural networks are always better than decision trees or vice-versa. There are many factors at play, such as the size and structure of your dataset.

As a result, you should try many different algorithms for your problem, while using a hold-out "test set" of data to evaluate performance and select the winner. Hypothesis space + Risk + Optimization

CART FUNCTIONALITY

SUPERVISED

NON-PARAMETRIC

WHITE-BOX

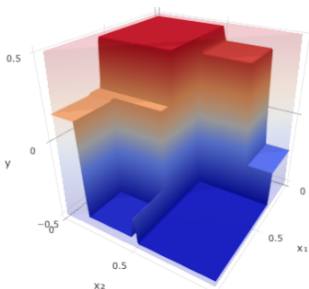
FEATURE SELECTION

General idea Starting from a root node, **classification & regression trees (CART)** perform repeated **binary splits** of the data according to feature values, thereby subsequently dividing the input space \mathcal{X} into T **rectangular partitions** Q_t .

Hypothesis space

$$\mathcal{H} = \left\{ f(\mathbf{x}) : f(\mathbf{x}) = \sum_{t=1}^T c_t \mathbb{I}(\mathbf{x} \in Q_t) \right\}$$

- Pass observations along until each ends up in exactly one leaf node
- In each step, find the optimal feature-threshold combination to split by
- Assign response c_t to leaf node m



CART FUNCTIONALITY

Empirical risk for typical loss functions

- Empirical risk is calculated per terminal node \mathcal{N}_m .
- Classification (for g classes):

- Using **Brier score**

$$\mathcal{R}(\mathcal{N}_t) = \sum_{(\mathbf{x}, y) \in \mathcal{N}_t} \sum_{k=1}^g (\mathbb{I}(y = k) - \pi_k(\mathbf{x}))^2$$

- Using **Bernoulli loss**

$$\mathcal{R}(\mathcal{N}_t) = \sum_{(\mathbf{x}, y) \in \mathcal{N}_t} \sum_{k=1}^g \mathbb{I}(y = k) \cdot \log(\pi_k(\mathbf{x}))$$

- Regression: Using **quadratic loss**

$$\mathcal{R}(\mathcal{N}_t) = \sum_{(\mathbf{x}, y) \in \mathcal{N}_t} (y - c_t)^2$$

Optimization Exhaustive search for optimal splitting criterion (greedy optimization)

Hyperparameters Tree depth, minimum number of observations per node, ...

CART PRO'S & CON'S

Advantages

- + **Easy** to understand, interpret & visualize
- + Automatic handling of **non-numerical** features
- + Built-in **feature selection**
- + Automatic handling of **missings**
- + **Interaction** effects between features easily possible, even of higher orders
- + **Fast** computation and good scalability
- + High **flexibility** (custom split criteria or leaf-node prediction rules)

Disadvantages

- Rather **low accuracy** (at least, without bagging or boosting)
- High **variance/instability**: strong dependence on training data
- Therefore, poor generalization & risk of **overfitting**
- Several steps required for modeling linear relationships
- In presence of categorical features, **bias** towards features with **many categories**

Simple and good with feature selection, but not the best predictor

CART APPLICATION

For applications of CART, note the following:

Pruning / early stopping

Unless interrupted, splitting will go on until each leaf node contains a single observation (expensive + overfitting!)

→ Use **pruning** and **stopping criteria** to limit complexity.

Implementation

R: package `rpart`

Python: `DecisionTreeClassifier` from package `scikit-learn`

Bagging

Since CART are instable predictors on their own, they are typically ensembled to form a **random forest**.

RANDOM FORESTS FUNCTIONALITY

SUPERVISED

NON-PARAMETRIC

BLACK-BOX

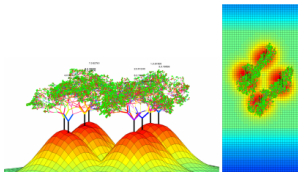
FEATURE SELECTION

General idea Random forests are **bagging ensembles**: they combine multiple CART (weak learners) to form a strong learner. They use **complex**, fully grown trees with low bias and compensate for single trees' high variance by aggregating M of them in a **decorrelated** manner.

Hypothesis space

$$\mathcal{H} = \left\{ f(\mathbf{x}) : f(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \sum_{t=1}^{T^{[m]}} c_t^{[m]} \mathbb{I}(\mathbf{x} \in Q_t^{[m]}) \right\}$$

- Training on bootstrap samples of the data and only on a random subset of features to incur variability
- Aggregation via averaging (regression) or majority voting (classification)



RANDOM FORESTS FUNCTIONALITY

Empirical risk for typical loss functions

- Empirical risk is calculated per terminal node \mathcal{N}_m .
- Classification (for g classes):

- Using **Brier score**

$$\mathcal{R}(\mathcal{N}_m) = \sum_{(\mathbf{x}, y) \in \mathcal{N}_m} \sum_{k=1}^g (\mathbb{I}(y = k) - \pi_k(\mathbf{x}))^2$$

- Using **Bernoulli loss**

$$\mathcal{R}(\mathcal{N}_m) = \sum_{(\mathbf{x}, y) \in \mathcal{N}_m} \sum_{k=1}^g \mathbb{I}(y = k) \cdot \log(\pi_k(\mathbf{x}))$$

- Regression: Using **quadratic loss**

$$\mathcal{R}(\mathcal{N}_m) = \sum_{(\mathbf{x}, y) \in \mathcal{N}_m} (y - c_m)^2$$

Optimization Exhaustive search for optimal splitting criterion (greedy optimization)

Hyperparameters Tree depth, minimum number of observations per node, ...

RANDOM FORESTS PRO'S & CON'S

Advantages

+ Foo

Disadvantages

- Foo

Foo

RANDOM FORESTS APPLICATION

For applications of CART, note the following:

Foo

Foo