

James G2m

TAE WK 9 - Bagging & Random Forest

9.1 Bootstrapping

A resampling method used to quantify the uncertainty associated in a statistical learning method.

(Other resampling method: cross-validation)

Goal: establish empirical distribution functions for a wide spread range of statistics

- How

- We obtain distinct datasets by repeatedly sampling observations from the original dataset in replacement
- Each bootstrap dataset is created by sampling in replacement, & is the same size as our original dataset.

- Terminology

- Z^{*1} : first bootstrap data set
- $\hat{\alpha}^{*1}$: bootstrap estimate for α
- B : NO. of different bootstrap datasets

Standard error
of bootstrap
estimates.

$$SE_B(\hat{\alpha}) = \sqrt{\frac{1}{B-1} \sum_{i=1}^B (\hat{\alpha}^{*i} - \bar{\hat{\alpha}})^2}$$

- Where it can be applied

It can be used to tackle the bias-variance trade-off of some statistical learning methods (such as decision tree)

8.2 Bagging (Bootstrap aggregation)

Intuition: Decision Tree suffer from high variance. (Say if we split the training data into 2 parts at random, & fit a decision tree to both halves, the results would be quite diff)

Method:

- Take repeated samples from the training dataset (Generate B diff bootstrapped training datasets)
- Train model $f^{*b}(x)$ on the b -th bootstrapped training datasets (repeat for all B datasets)
- We average all the predictions as follows

$$f_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B f^{*b}(x)$$

(take majority vote for classification tree)

Entire process
is called
Bagging

8.3 Out-Of-Bag (OOB) Error Estimation

We use OOB Error Estimation to test the error of the bagged model.

- To obtain a single prediction for the i^{th} observation, we average predicted responses (for regression) or, majority vote (for classification).

\therefore this leads to a single OOB prediction for the i^{th} observation.

- An OOB prediction can be obtained in this way for each of the n observations, from which we compute the overall OOB MSE/classification error.

- When B is sufficiently large, we have a decent alternative to cross-validation.

8.4 Random Forest

Random forest builds from the idea of considering a random sample of m predictors from the full set of predictors P .

\rightarrow this is to overcome the problem \square DT bagging where trees tend to be correlated.

Idea:

- for each split, we consider only a subset of predictors ($P \subset m$)

$\hookrightarrow \therefore$ on average $(p-m)$ predictors are not considered, so other predictors will have a chance.
(This process decorrelates a tree).

Algorithm:

1. Generate B bootstrapped training datasets.
2. For each dataset, train a decision tree. (At each split, use a subset m of p available predictors.)
3. Average the predictions from the B trees.
(For classification, use majority voting)

value of m :

- Regression, $m = \frac{p}{3}$
- Classification, $m = \sqrt{p}$

Hyperparameters Tuning.

- No. of trees, B
- No. of predictors used on each split, m

Jones 62

TAE WK 9 - Ethics of data analytics

Ethical issues in data science

① Bias, discrimination, & exclusion

Algos & AI create exclusion towards individuals & groups of people.

② Algorithmic Profiling

The model of personality can affect the reg. collection principles like democratic & cultural pluralism & risk sharing in the realm of insurance.

③ Privately massive files while enhancing AI.

Data protection laws are rooted in the belief that individuals' rights regarding their personal data must be protected.

④ Quality, quantity & relevance

The acceptance of the existence of potential bias in datasets created to train algo is of paramount importance.

How data visualization is abused

1) Omitting the baseline.

(Not starting at the origin).

2) Manipulating the y-axis

(changing the scale of y-axis to exaggerate or conceal trends)

3) Cherry picking data.

4) Using the wrong graph

(wrong type of visualization).

5) Improper intervals or units

6) Omitting data

7) Extrapolation

8) Unnecessary complexity