

Machine Learning

Bagging and Boosting

S. Nõmm

¹Department of Software Science, Tallinn University of Technology

18.04.2019

Lectures and practices in May

- May the 2nd: Lecture as usual. Defend Home assignment 3 during the practice.
- May the 9th Closed book test 2. Practice is reserved for those who need to defend there home assignments 1 , 2 and 3
- May the 16th. Make up tests during the lecture. Practice time is reserved for consultation.

Bootstrap

- Remind what is the main goal of cross validation.
- Let $Z = (z_1, \dots, z_n)$ is the training set.
- Draw randomly data sets with replacement (the samples are independent) from Z . This will result in B *bootstrap* data sets.
- Fit the model for each of B data sets. Examine behaviour over B replacements.
- This approach allows to estimate any aspect of distribution $S(Z)$.

Bagging

- Induced from the bootstrap technique (which is used to assess accuracy of estimate).
- Draw B samples with replacements and train the model on each sample.
- The bagging estimate then is defined by:

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

Random Forests

The idea is to build large collection of de-correlated trees, and then average them.

- For $b = 1$ to B :
 - ▶ Draw a bootstrap sample Z^* of size N from the available training data.
 - ▶ Grow tree T_b . Repeat recursively for each terminal node until minimum node size is reached.
 - ★ Select m variables from p .
 - ★ Pick the best variable among m .
 - ★ Split the node.
- Output the ensemble of trees $\{T_b\}_1^B$.
- Prediction:
 - ▶ Regression: $\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$.
 - ▶ Classification: $\hat{C}_{\text{rf}}^B(x) = \text{mode}\{\hat{C}_b(x)\}_1^B$.

Committee learning

- Some times referred as ensemble learning.
- The idea is to combine a number of weak (accuracy is slightly larger than of random guessing) classifiers into a powerful committee.
- Motivation is to improve estimate by reducing variance and sometimes bias.

Boosting

- The final prediction is given by:

$$G(x) = \text{sign}\left(\sum_{m=1}^M \alpha_m G_m(x)\right).$$

which is weighted majority vote of classifiers $G_m(x)$. Here α_m are weights describing contribution of each classifier.

- While on the first view result is very similar to the bagging, there are some major differences.
- Two class problem where output variable coded as $Y \in \{-1, 1\}$.
- For the classifier $G(X)$ error rate is given by:

$$\overline{\text{err}} = \frac{1}{N} \sum_{i=1}^N I(y_i \neq G(x_i)),$$

where N is the power of training data set.

Ada Boost

AdaBoost.M1. by Freund and Shapire (1997).

- Initialize observation weights $w_i = 1/N$, $i = 1, \dots, N$.
- For $m = 1$ to M :
 - ▶ Fit weak classifier G_m that minimizes the weighted sum error for misclassified points.

$$\epsilon_m = \frac{\sum_{i=1}^N w_i I(G_m(x_i) \neq y_i)}{\sum_{i=1}^N w_i}$$

- ▶ Compute $\alpha_m = \log((1 - \epsilon_m)/\epsilon_m)$.
- ▶ Update weights w_i as

$$w_i = w_i * \exp(\alpha_m * I(y_i \neq G_m(x_i))), \quad i = 1, \dots, N.$$

- Output classifier:

$$G(x) = \text{sign}\left(\sum_{m=1}^M \alpha_m G_m(x)\right).$$