

TREES, FORESTS AND BOOSTING  
KNOWLEDGE SHARING SESSION

**Rostislav Stoyanov**

May 28, 2025

# Part I

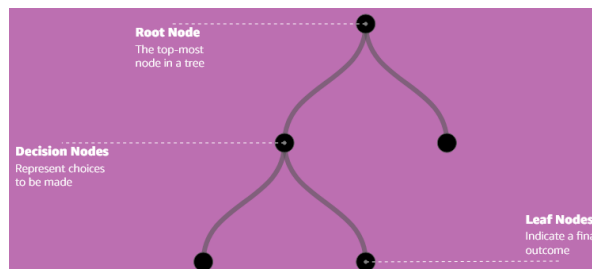
## TREES AND FORESTS

# TREES AND FORESTS

<b>1</b>	<b>Decision trees</b>	<b>3</b>
1.1	Tree training	4
1.2	Advantages of decision trees	7
1.3	Issues with decision trees	8
<b>2</b>	<b>Ensemble learning</b>	<b>10</b>
<b>3</b>	<b>Bagging</b>	<b>11</b>
3.1	Bootstrapping	11
3.2	Bagging algorithm	12
3.3	Random forest	13
3.4	Advantages of random forest	14
3.5	Disadvantages of random forest	15

# DECISION TREES

- ▶ Structure consisting of (root, internal and left) nodes connected with branches.
- ▶ Each node represents a test on a data feature and each branch from the node represent an outcome from the test.
- ▶ Nodes without children are called leafs. Each leaf represents a prediction outcome of the model.



**Figure.** Tree structure parts<sup>1</sup>

<sup>1</sup>Picture taken from: <https://mlu-explain.github.io/decision-tree/>

# DECISION TREES

## TREE TRAINING

- ▶ Trees are built in a top-down fashion by performing a greedy search. Multiple algorithms, which follow similar approach.
- ▶ Starting from the root node, a recursive algorithm is used to find the most optimal current split (test) at a node.
- ▶ The selection of a split at a node is based on a metric. Using the samples reaching the node, this metric is calculated in order to find which of all possible splits best optimizes the algorithm's objective.
- ▶ This process is repeated until some end condition is met (e.g. not enough data to split, maximum depth is reached).

# DECISION TREES

## TRAINING METRICS

### Definition 1.1 (Shannon entropy)

Given a discrete random variable  $X$ , which takes values in the alphabet  $\mathcal{X}$  and is distributed according to  $p : \mathcal{X} \rightarrow [0, 1]$ :

$$H(X) := - \sum_{x \in \mathcal{X}} p(x) \log p(x) = \mathbb{E}[-\log p(X)] \quad (1)$$

### Definition 1.2 (Information gain)

Let  $T$  denote a set of training examples and  $\alpha$  be an attribute describing the training examples. Furthermore, let  $\text{vals}(\alpha)$  contains all possible values of  $\alpha$  and  $S_\alpha(v) := \{x \in T | x_\alpha = v\}$  be the set of training inputs for which the attribute  $\alpha = v$ .

$$\text{IG}(T, \alpha) := H(T) - \sum_{v \in \text{vals}(\alpha)} \frac{|S_\alpha(v)|}{|T|} H(S_\alpha(v)) = H(T) - H(T | \alpha) \quad (2)$$

# DECISION TREES

## TRAINING VISUALIZED

---

Source: <https://mlu-explain.github.io/decision-tree/>

# DECISION TREES

## ADVANTAGES OF DECISION TREES

- ▶ Can solve non-linear problems.
- ▶ Non-parametric models.
- ▶ Can handle various both categorical and numerical data. Furthermore, missing values are also supported.
- ▶ Easy to interpret.



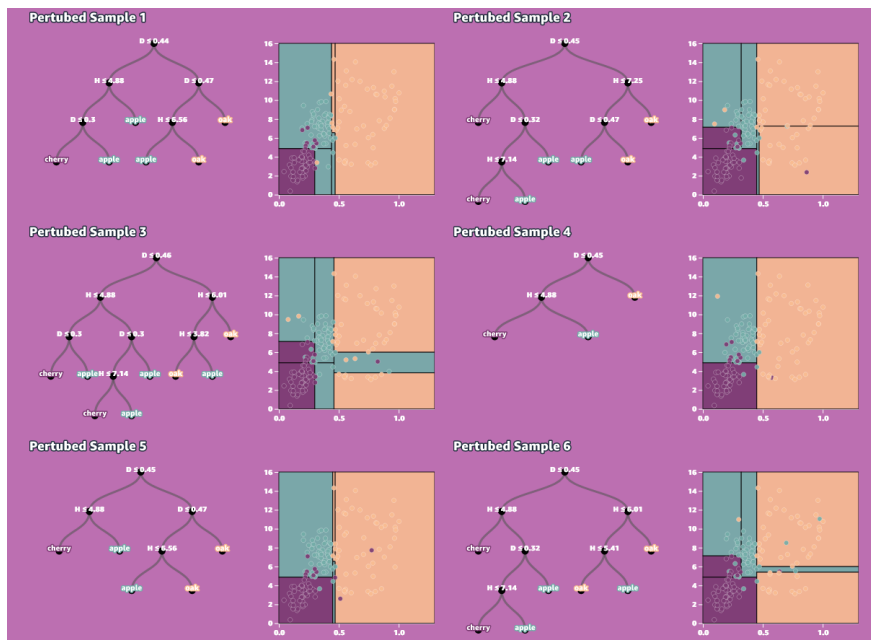
# DECISION TREES

## ISSUES WITH DECISION TREES

- ▶ Have to re-calculate the tree for new examples.
- ▶ "Vanilla" decision trees can't extrapolate.
- ▶ Overfitting. Complex trees fit the training data too well and do not generalize to new data.
  - We can use pre-pruning or post-pruning to mitigate this problem.
- ▶ Unstable. Small changes in the training samples can lead to very different decision trees.

# DECISION TREES

## HIGH VARIANCE VISUALIZED



**Figure.** Example of trees build when changing the class of 5% of the training samples <sup>2</sup>

<sup>2</sup>Source: <https://mlu-explain.github.io/decision-tree/>

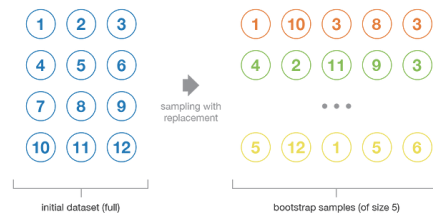
## ENSEMBLE LEARNING

- ▶ Use ensemble learning to mitigate the weaknesses of decision trees.
- ▶ Ensemble learning is a paradigm in which multiple "weak learners" are trained on the same problem to achieve better performance.
- ▶ We will explore two methods:
  - Bagging (random forest)
  - Boosting (AdaBoost, gradient boosting)

# BAGGING

## BOOTSTRAPING

- ▶ Generate new samples of size  $B$  from a dataset with  $N$  samples using uniform sampling with replacement.
- ▶ Bootstrapping needs two properties of the original dataset to work correctly (the samples are representative and independent samples of the true data distribution)<sup>1</sup>:
  - Representativity - The dataset is big enough to capture the underlying distribution.
  - Independence - The dataset is big enough so samples aren't correlated much.



**Figure. Bootstrap**<sup>1</sup>

<sup>1</sup>Source: <https://towardsdatascience.com/ensemble-methods-bagging-boosting-and-stacking-c9214a10a205>

# BAGGING

## BAGGING ALGORITHM

- ▶ Given a training set  $D$ ,  $|D| = n$ , bootstrap  $m$  new training sets  $D_i, i = 1, \dots, m, |D_i| = n'$ .
- ▶ Using each of the training sets  $D_i$  train a machine learning model  $M_i$ .
- ▶ The output of the ensemble is the average of all model outputs (regression) or voting (classification).
- ▶ The final model has less variance than a single trained model.

# BAGGING

## RANDOM FOREST

- ▶ Random forest uses bagging to train multiple decision trees. We are using bagging to minimize variance, so use deep trees.
- ▶ Additionally, for each decision tree trained keep only a part of the features. This has two advantages:
  - Reduce correlation of the trees – if a feature is a very strong predictor it will be used in all trees making them correlated.
  - Increase robustness to missing data. Some of the trees will only have features without missing data.

# BAGGING

## ADVANTAGES OF RANDOM FOREST

- ▶ Similar to decision trees can handle different types of features and solve non-linear problems.
- ▶ Better performance than a single decision tree.
- ▶ Low risk of overfitting.

# BAGGING

## DISADVANTAGES OF RANDOM FOREST

- ▶ Loss of interpretability (compared to decision trees).
- ▶ Can't extrapolate.
- ▶ Typically higher bias than a single tree [Has+09].
- ▶ Very dependent on the bootstrapped datasets. Recreation of experiment would require the same sampling (keeping track of the seeds).
- ▶ Needs modifications to handle class imbalance.



## Part II

# BOOSTING

- ▶ Similarly to bagging uses many "weak" models to produce a strong ensemble model ("committee"). The two approaches, however, are fundamentally different.
- ▶ Models are fit iteratively (sequentially) so that they complement each other.
- ▶ Each new model attempts to improve prediction on samples for which the previous models have been wrong.

# BOOSTING

<b>1</b>	<b>Boosting</b>	<b>18</b>
1.1	AdaBoost	19
<b>2</b>	<b>Additive modeling</b>	<b>23</b>
2.1	Forward Stagewise Additive Modeling	24
<b>3</b>	<b>Functional gradient descent</b>	<b>25</b>
3.1	Steepest descent	26
3.2	Gradient descent	28
<b>4</b>	<b>Next time</b>	<b>31</b>
<b>5</b>	<b>Useful resources</b>	<b>32</b>

# BOOSTING

## ADABOOST

---

**Algorithm** AdaBoost [Has+09]

---

- 1: Initialize the observation weights  $w_i = 1/N, i = 1, 2, \dots, N$ .
- 2: **for**  $m = 1, \dots, M$  **do**
- 3:     Fit a classifier  $G_m(x)$  to the training data using weights  $w_i$ .
- 4:     Compute

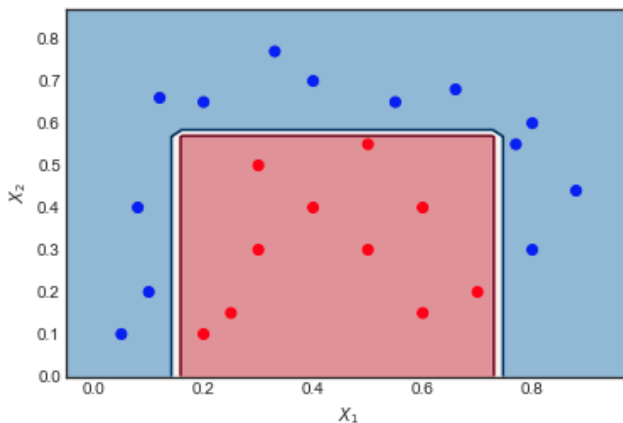
$$\text{err}_m = \frac{\sum_{i=1}^N w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i}.$$

- 5:     Compute  $\alpha_m = \log((1 - \text{err}_m) / \text{err}_m)$ .
  - 6:     Set  $w_i \leftarrow w_i \cdot \exp[\alpha_m \cdot I(y_i \neq G_m(x_i))], i = 1, 2, \dots, N$ .
  - 7: **end for**
  - 8: Output  $G(x) = \text{sign}\left[\sum_{m=1}^M \alpha_m G_m(x)\right]$ .
-

# BOOSTING

## ADABOOST

- ▶ On the right, AdaBoost fit for a toy example.
- ▶ Fit for 10 iterations.
- ▶ Weak learners are decision stumps with depth 1 and 2 leaf nodes.

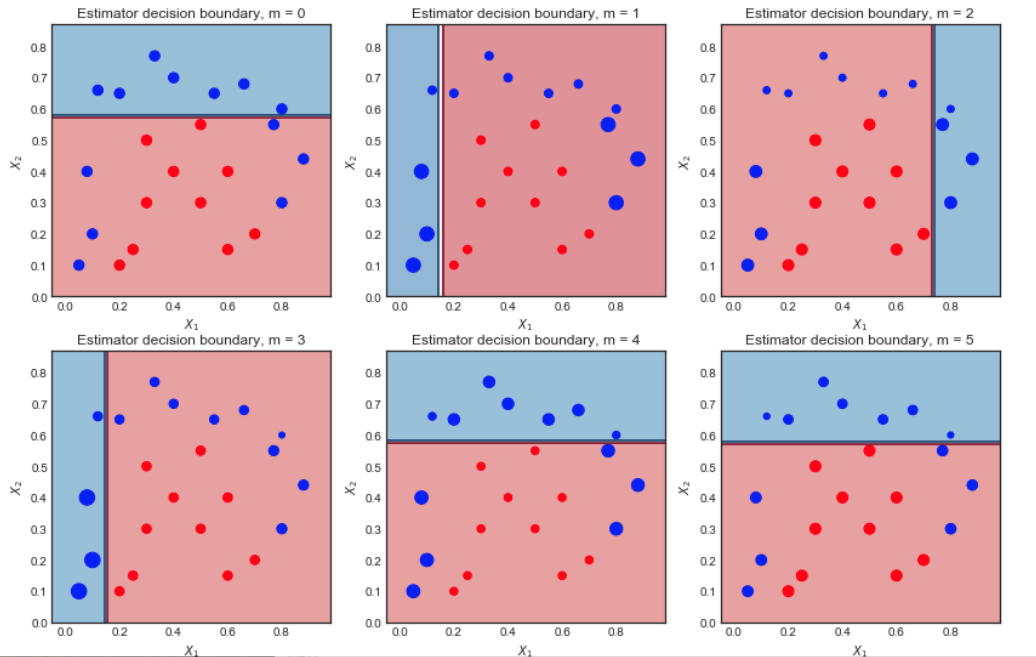


**Figure.** AdaBoost on a toy example<sup>3</sup>

<sup>1</sup>Picture taken from: <https://xavierbourretsicotte.github.io/AdaBoost.html>

# BOOSTING

## ADABOOST



**Figure.** Visualization of the sequence of weak learners and the sample weights

# BOOSTING

## ADABOOST

- ▶ Advantages:
  - Improves performance of weak classifiers fast.
  - Easy to use and tune, not a lot of parameters than a single classifier.
  - Less prone to overfitting.
- ▶ Drawbacks:
  - Sensitive to noisy data and outliers.
  - Hard to extend to other tasks.

## ADDITIVE MODELING

- ▶ Boosting success comes from fitting an additive expansion:

$$f(x) = \sum_{m=1}^M \beta_m b(x; \gamma_m) \quad (3)$$

In the case of AdaBoost:

- $b(x; \gamma) \in \mathbb{R}$  are functions of multivariate argument  $x$ , with parameters  $\gamma$ . AdaBoost uses trees as basis functions, with  $\gamma$  being the split variables and points as well as the leaf predictions.
  - The expansion coefficients  $\beta_m; m = 1, 2, \dots, M$  are the model weights  $\alpha_m$ .
- ▶ Fitting is done by minimizing a loss over all training examples:

$$\min_{\{\beta_m, \gamma_m\}_1^M} \sum_{i=1}^N L \left( y_i, \sum_{m=1}^M \beta_m b(x_i; \gamma_m) \right). \quad (4)$$



# ADDITIVE MODELING

## FORWARD STAGewise ADDITIVE MODELING

- ▶ Fitting 4 is very expensive, so instead we can approximate by adding new functions without modifying the previous ones:

---

**Algorithm** Forward Stagewise Additive Modeling [Has+09]

---

- 1: Initialize  $f_0(x) = 0$ .
- 2: **for**  $m = 1, \dots, M$  **do**
- 3:     Compute

$$(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma)).$$

- 4:     Set  $f_m(x) = f_{m-1}(x) + \beta_m b(x; \gamma_m)$ .
  - 5: **end for**
- 

- ▶ It can be proven that AdaBoost is a forward stagewise additive modeling with loss function:  
 $L(y, f(x)) = \exp(-y f(x))$
- ▶ We need to be able to solve step 3 for general loss criteria.

## FUNCTIONAL GRADIENT DESCENT

- ▶ Explore numerical optimization approaches.
- ▶ Want to minimize:

$$J(f) = \sum_{i=1}^N L(y_i, f(x_i))$$

- ▶ However  $J(f)$  depends on  $f$  only on  $N$  training points (vector of model predictions on training points):

$$\mathbf{f} = \{f(x_1), f(x_2), \dots, f(x_N)\}^T.$$

- ▶ So we can rewrite  $J(f)$ :

$$J(f) = \sum_{i=1}^N L(y_i, \mathbf{f})$$

# FUNCTIONAL GRADIENT DESCENT

## STEEPEST DESCENT

- Perform gradient descent on:

$$J(f) = \sum_{i=1}^N L(y_i, \mathbf{f}) ,$$

where  $\mathbf{f}_0 = \mathbf{h}_0$  is an initial guess, and each successive  $\mathbf{f}_m$  is induced based on the current parameter vector  $\mathbf{f}_{m-1}$ , which is the sum of the previously induced updates.

- The negative gradient direction at  $\mathbf{f} = \mathbf{f}_{m-1}$  is  $\mathbf{g}_m$ :

$$g_{im} = \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i)=f_{m-1}(x_i)}$$

- We also need a weight  $\rho_m$  to perform an update. Two options:

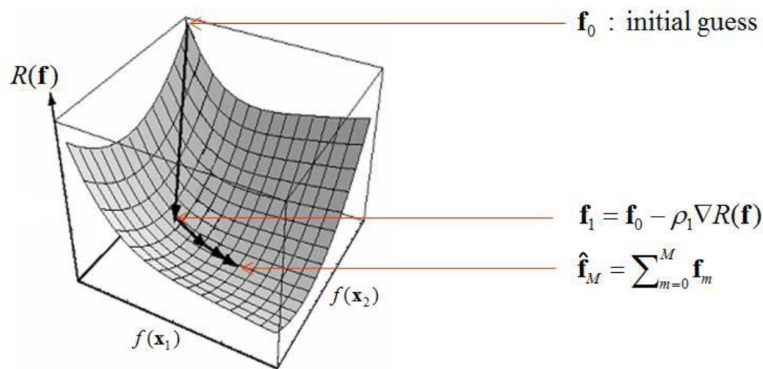
- $\rho_m = \arg \min_{\rho} L(\mathbf{f}_{m-1} - \rho \mathbf{g}_m)$ .
- Some constant value (e.g. 0.01).

- $\mathbf{f}_m = \mathbf{f}_{m-1} - \rho_m \mathbf{g}_m$

# FUNCTIONAL GRADIENT DESCENT

## STEEPEST DESCENT

- ▶  $-g_m \in \mathbb{R}$  is the direction in which we want to update our training data predictions.
- ▶ We need more than  $f$ , which is predictions on training data.
  - Might not be defined for test points.
  - Might not exist a hypothesis (tree) that generates this outcome on the training data.



**Figure.** Steepest descent illustration. The “risk” criterion  $R(\mathbf{f})$ , is plotted as a function of  $\mathbf{f}$  evaluated at two training data points  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . [SE10]

# FUNCTIONAL GRADIENT DESCENT

## GRADIENT DESCENT

- ▶ The gradient is defined only on the  $N$  training points and we want to generalize  $\mathbf{f}_m$ .
- ▶ A solution is to find the hypothesis  $h \in \mathcal{H}$  that best approximates  $-g$  as out step direction.
- ▶ Approach this a least squares regression problem over  $\mathcal{H}$ :

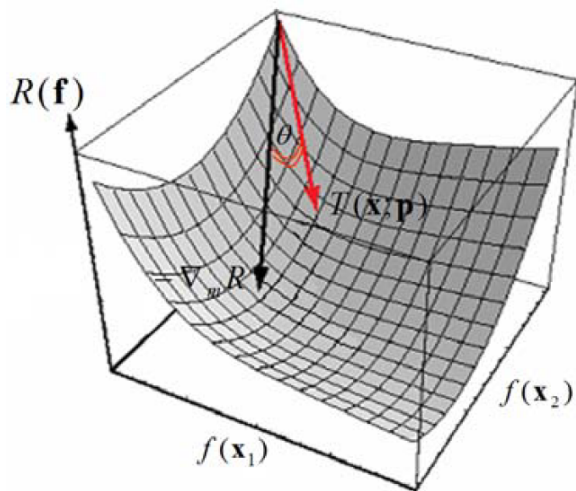
$$\min_{h \in \mathcal{H}} \sum_{i=1}^N (-g - h(x_i))^2$$

In other words, we are fitting our model over the pseudo-residuals.

# FUNCTIONAL GRADIENT DESCENT

## GRADIENT DESCENT

- ▶  $T(x; p) \in \mathcal{H}$  is the actual step direction (model) we are using.
- ▶ Can think of it as projection of  $-g$  onto  $\mathcal{H}$ .



**Figure.** Surrogate-loss illustration. The base-learner most parallel to the negative gradient vector is chosen at every step of the gradient-descent algorithm [SE10]

# FUNCTIONAL GRADIENT DESCENT

## GRADIENT DESCENT

---

**Algorithm** Gradient boosting [Fri01]

---

- 1:  $F_0(\mathbf{x}) = \arg \min_{\rho} \sum_{i=1}^N L(y_i, \rho)$
  - 2: **for**  $m = 1, \dots, M$  **do**
  - 3:    $\tilde{y}_i = - \left[ \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}, i = 1, \dots, N$
  - 4:    $\mathbf{a}_m = \arg \min_{\mathbf{a}, \beta} \sum_{i=1}^N [\tilde{y}_i - \beta h(\mathbf{x}_i; \mathbf{a})]^2$
  - 5:    $\rho_m = \arg \min_{\rho} \sum_{i=1}^N L(y_i, F_{m-1}(\mathbf{x}_i) + \rho h(\mathbf{x}_i; \mathbf{a}_m))$
  - 6:    $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$
  - 7: **end for**
-

## NEXT TIME

- ▶ XGBoost and friends.
- ▶ Handling categorical features.



## USEFUL RESOURCES

- ▶ For trees [MM97]
- ▶ Random forest + boosting <https://www.youtube.com/watch?v=fz1H03ZKvLM>
- ▶ Gradient boosting <https://www.youtube.com/watch?v=wPqtzj5VZus>

## REFERENCES I

- [Fri01] Jerome H Friedman. “Greedy function approximation: a gradient boosting machine”. In: *Annals of statistics* (2001), pp. 1189–1232.
- [Has+09] Trevor Hastie et al. *The elements of statistical learning: data mining, inference, and prediction*. Vol. 2. Springer, 2009.
- [MM97] Tom M Mitchell and Tom M Mitchell. *Machine learning*. Vol. 1. 9. McGraw-hill New York, 1997.
- [SE10] Giovanni Seni and John F Elder. “Ensemble methods in data mining: improving accuracy through combining predictions”. In: *Synthesis lectures on data mining and knowledge discovery* 2.1 (2010), pp. 1–126.