

# Gradient Boosted Decision Trees

Jindřich Libovický (reusing materials by Milan Straka)

 December 05, 2023



Charles University in Prague  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated

After this lecture you should be able to

- Explain second-order optimization methods
- Implement gradient boosted decision trees for regression and classification
- Decide what supervised machine learning approach is suitable for particular problems

# Gradient Boosting Decision Trees

The gradient boosting decision trees also train a collection of decision trees, but unlike random forests, where the trees are trained independently, in GBDT they are trained sequentially to correct the errors of the previous trees.

If we denote  $y_t$  as the prediction function of the  $t^{\text{th}}$  tree, the prediction of the whole collection is then

$$y(\mathbf{x}_i) = \sum_{t=1}^T y_t(\mathbf{x}_i; \mathbf{w}_t),$$

where  $\mathbf{w}_t$  is a vector of parameters (leaf values, to be concrete) of the  $t^{\text{th}}$  tree.

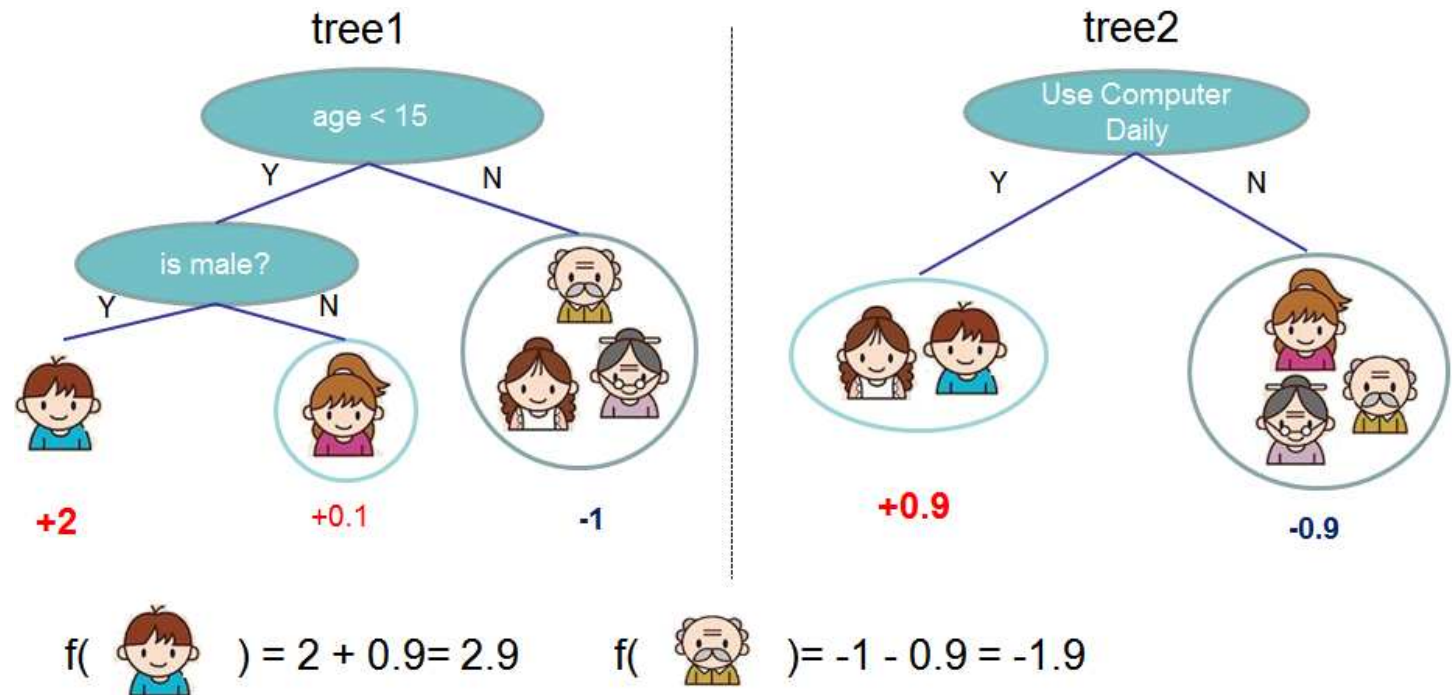


Figure 1 of "XGBoost: A Scalable Tree Boosting System", <https://arxiv.org/abs/1603.02754>

Considering a regression task first, we define the overall loss as

$$E(\mathbf{w}) = \sum_i \ell(t_i, y(\mathbf{x}_i; \mathbf{w})) + \sum_{t=1}^T \frac{1}{2} \lambda \|\mathbf{w}_t\|^2,$$

where

- $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_T)$  are the parameters (leaf values) of the trees;
- $\ell(t_i, y(\mathbf{x}_i; \mathbf{w}))$  is an per-example loss,  $(t_i - y(\mathbf{x}_i; \mathbf{w}))^2$  for regression;
- the  $\lambda$  is the usual  $L^2$ -regularization strength.

To construct the trees sequentially, we extend the definition to

$$E^{(t)}(\mathbf{w}_t; \mathbf{w}_{1..t-1}) = \sum_i \left[ \ell(t_i, y^{(t-1)}(\mathbf{x}_i; \mathbf{w}_{1..t-1}) + y_t(\mathbf{x}_i; \mathbf{w}_t)) \right] + \frac{1}{2} \lambda \|\mathbf{w}_t\|^2.$$

In the following text, we drop the parameters of  $y^{(t-1)}$  and  $y_t$  for brevity.

The original idea of gradient boosting was to set

$$y_t(\mathbf{x}_i) \leftarrow - \frac{\partial \ell(t_i, y^{(t-1)}(\mathbf{x}_i))}{\partial y^{(t-1)}(\mathbf{x}_i)} = - \frac{\partial \ell(t_i, y)}{\partial y} \Bigg|_{y=y^{(t-1)}(\mathbf{x}_i)}$$

as a direction minimizing the residual loss and then finding a suitable constant  $\gamma_t$ , which would minimize the loss

$$\sum_i \left[ \ell(t_i, y^{(t-1)}(\mathbf{x}_i) + \gamma_t y_t(\mathbf{x}_i)) \right] + \frac{1}{2} \lambda \|\mathbf{w}_t\|^2.$$

# First-order and Second-order Methods

Until now, we used mostly SGD for finding a minimum, by performing

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla E(\mathbf{w}).$$

A disadvantage of this (so-called **first-order method**) is that we need to specify the learning rates by ourselves, usually using quite a small one, and perform the update many times.

However, in some situations, we can do better.

# Newton's Root-Finding Method

Assume we have a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  and we want to find its root. An SGD-like algorithm would always move “towards” zero by taking small steps.

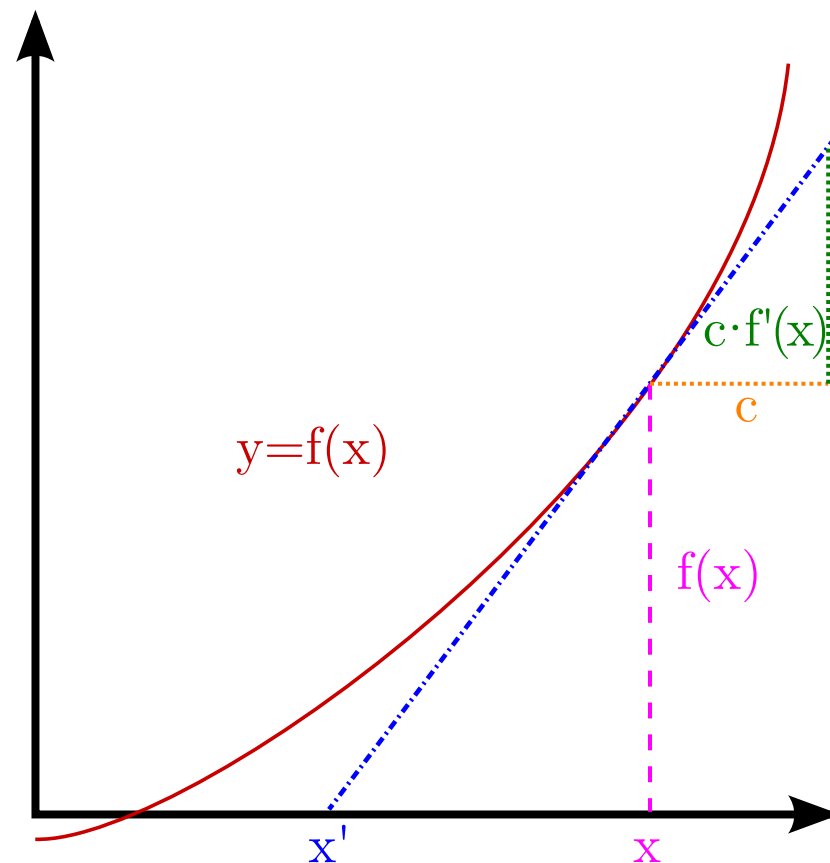
Instead, we could consider the linear local approximation (i.e., consider a line “touching” the function in a given point) and perform a step so that our linear local approximation has a value 0:

$$x' \leftarrow x - \frac{f(x)}{f'(x)}.$$

## Finding Minima

The same method can be used to find minima, because a minimum is just a root of a derivative, resulting in:

$$x' \leftarrow x - \frac{f'(x)}{f''(x)}.$$



Modification of [https://commons.wikimedia.org/wiki/File:Newton-Raphson\\_method.svg](https://commons.wikimedia.org/wiki/File:Newton-Raphson_method.svg)

The following update is the Newton's method of searching for extremes:  $x' \leftarrow x - \frac{f'(x)}{f''(x)}$ .

It is a so-called **second-order** method, but it is just an SGD update with a learning rate  $\frac{1}{f''(x)}$ .

## Derivation from Taylor's Expansion

The same update can be derived also from the Taylor's expansion ( $x$  is a fixed point and  $\epsilon$  is now the variable that moves)

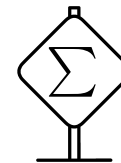
$$f(x + \epsilon) \approx f(x) + \epsilon f'(x) + \frac{1}{2} \epsilon^2 f''(x) + \mathcal{O}(\epsilon^3),$$

which we can minimize for  $\epsilon$  by (i.e., the minimum of the approximation)

$$0 = \frac{\partial f(x + \epsilon)}{\partial \epsilon} \approx f'(x) + \epsilon f''(x), \text{ obtaining } x + \epsilon = x - \frac{f'(x)}{f''(x)}.$$



Note that the second-order methods (methods utilizing second derivatives) are impractical when training MLPs (and GLMs) with many parameters. The problem is that there are too many second derivatives – if we consider weights  $\mathbf{w} \in \mathbb{R}^D$ ,



- the gradient  $\nabla E(\mathbf{w})$  has  $D$  elements;
- however, we have a  $D \times D$  matrix with all second derivatives, called the **Hessian**  $\mathbf{H}$ :

$$H_{i,j} \stackrel{\text{def}}{=} \frac{\partial^2 E(\mathbf{w})}{\partial w_i \partial w_j}.$$

For completeness, the Taylor expansion of a multivariate function then has the following form:

$$f(\mathbf{x} + \boldsymbol{\varepsilon}) = f(\mathbf{x}) + \boldsymbol{\varepsilon}^T \nabla f(\mathbf{x}) + \frac{1}{2} \boldsymbol{\varepsilon}^T \mathbf{H} \boldsymbol{\varepsilon},$$

from which we obtain the following second-order method update:

$$\mathbf{x} \leftarrow \mathbf{x} - \mathbf{H}^{-1} \nabla f(\mathbf{x}).$$

Returning to the gradient boosting decision trees, instead of using a first-order method, it was later suggested that a second-order method could be used. Denoting

$$g_i = \frac{\partial \ell(t_i, y^{(t-1)}(\mathbf{x}_i))}{\partial y^{(t-1)}(\mathbf{x}_i)} = \left. \frac{\partial \ell(t_i, y)}{\partial y} \right|_{y=y^{(t-1)}(\mathbf{x}_i)}$$

and

$$h_i = \frac{\partial^2 \ell(t_i, y^{(t-1)}(\mathbf{x}_i))}{\partial y^{(t-1)}(\mathbf{x}_i)^2} = \left. \frac{\partial^2 \ell(t_i, y)}{\partial y^2} \right|_{y=y^{(t-1)}(\mathbf{x}_i)},$$

we can expand the objective  $E^{(t)}$  using a second-order approximation to

$$E^{(t)}(\mathbf{w}_t; \mathbf{w}_{1..t-1}) \approx \sum_i \left[ \ell(t_i, y^{(t-1)}(\mathbf{x}_i)) + g_i y_t(\mathbf{x}_i) + \frac{1}{2} h_i y_t^2(\mathbf{x}_i) \right] + \frac{1}{2} \lambda \|\mathbf{w}_t\|^2.$$

Recall that we denote the indices of instances belonging to a leaf  $\mathcal{T}$  as  $I_{\mathcal{T}}$ , and let us denote the prediction for the leaf  $\mathcal{T}$  as  $w_{\mathcal{T}}$ . Then we can rewrite

$$\begin{aligned} E^{(t)}(\mathbf{w}_t; \mathbf{w}_{1..t-1}) &\approx \sum_i \left[ g_i y_t(\mathbf{x}_i) + \frac{1}{2} h_i y_t^2(\mathbf{x}_i) \right] + \frac{1}{2} \lambda \|\mathbf{w}_t\|^2 + \text{const} \\ &\approx \sum_{\mathcal{T}} \left[ \left( \sum_{i \in I_{\mathcal{T}}} g_i \right) w_{\mathcal{T}} + \frac{1}{2} \left( \lambda + \sum_{i \in I_{\mathcal{T}}} h_i \right) w_{\mathcal{T}}^2 \right] + \text{const}. \end{aligned}$$

By setting a derivative with respect to  $w_{\mathcal{T}}$  to zero, we get

$$0 = \frac{\partial E^{(t)}}{\partial w_{\mathcal{T}}} = \sum_{i \in I_{\mathcal{T}}} g_i + \left( \lambda + \sum_{i \in I_{\mathcal{T}}} h_i \right) w_{\mathcal{T}}.$$

Therefore, the optimal weight for a node  $\mathcal{T}$  is

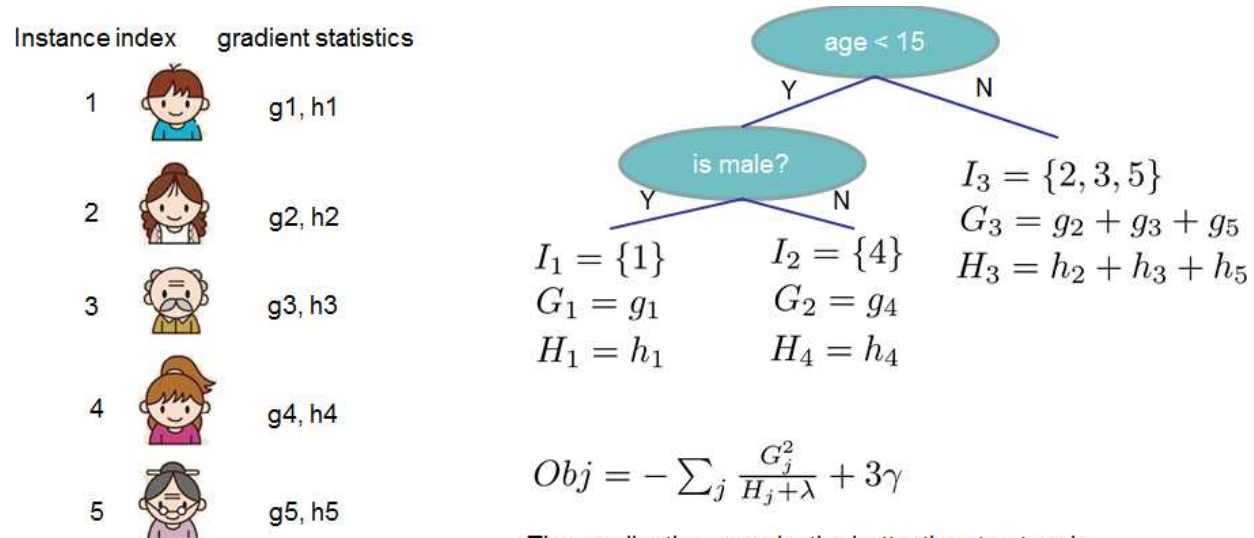
$$w_{\mathcal{T}}^* = - \frac{\sum_{i \in I_{\mathcal{T}}} g_i}{\lambda + \sum_{i \in I_{\mathcal{T}}} h_i}.$$

# Gradient Boosting

Substituting the optimum weights to the loss, we get

$$E^{(t)}(\mathbf{w}^*) \approx -\frac{1}{2} \sum_{\mathcal{T}} \frac{(\sum_{i \in I_{\mathcal{T}}} g_i)^2}{\lambda + \sum_{i \in I_{\mathcal{T}}} h_i} + \text{const},$$

which can be used as a *splitting criterion*.



The smaller the score is, the better the structure is

Figure 2 of "XGBoost: A Scalable Tree Boosting System", <https://arxiv.org/abs/1603.02754>

When splitting a node, the criteria of all possible splits can be effectively computed using the following algorithm:

---

**Algorithm 1: Exact Greedy Algorithm for Split Finding**

---

**Input:**  $I$ , instance set of current node

**Input:**  $D$ , feature dimension

$score \leftarrow -\infty$

$G \leftarrow \sum_{i \in I} g_i, H \leftarrow \sum_{i \in I} h_i$

**for**  $k = 1$  **to**  $D$  **do**

$G_L \leftarrow 0, H_L \leftarrow 0$

**for**  $j$  **in**  $sorted(I, \text{by } \mathbf{x}_{jk})$  **do**

$G_L \leftarrow G_L + g_j, H_L \leftarrow H_L + h_j$

$G_R \leftarrow G - G_L, H_R \leftarrow H - H_L$

**if**  $\mathbf{x}_{j_{next\ k}} \neq \mathbf{x}_{jk}$  **then**

$score \leftarrow \max(score, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda})$

**end**

**end**

**Output:** Split with max score

---

*Modified from Algorithm 1 of "XGBoost: A Scalable Tree Boosting System", <https://arxiv.org/abs/1603.02754>*

Furthermore, gradient boosting trees frequently use:

- data subsampling: either bagging or (even more commonly) only a fraction of the original training data is utilized for training of a single tree (with 0.5 being a common value),
- feature subsampling;
- shrinkage: multiply each trained tree by a learning rate  $\alpha$ , which reduces the influence of each individual tree and leaves space for future optimization.

To perform classification, we train the trees to perform the linear part of a generalized linear model.

Specifically, for a binary classification, we perform prediction by

$$\sigma(y(\mathbf{x}_i)) = \sigma\left(\sum_{t=1}^T y_t(\mathbf{x}_i; \mathbf{w}_t)\right),$$

and the per-example loss is defined as

$$\ell(t_i, y(\mathbf{x}_i)) = -\log\left[\sigma(y(\mathbf{x}_i))^{t_i} (1 - \sigma(y(\mathbf{x}_i)))^{1-t_i}\right].$$

For multiclass classification, we need to model the full categorical output distribution.

Therefore, for each “timestep”  $t$ , we train  $K$  trees  $\mathbf{w}_{t,k}$ , each predicting a single value of the linear part of a generalized linear model.

Then, we perform prediction by

$$\text{softmax}(\mathbf{y}(\mathbf{x}_i)) = \text{softmax}\left(\sum_{t=1}^T y_{t,1}(\mathbf{x}_i; \mathbf{w}_{t,1}), \dots, \sum_{t=1}^T y_{t,K}(\mathbf{x}_i; \mathbf{w}_{t,K})\right),$$

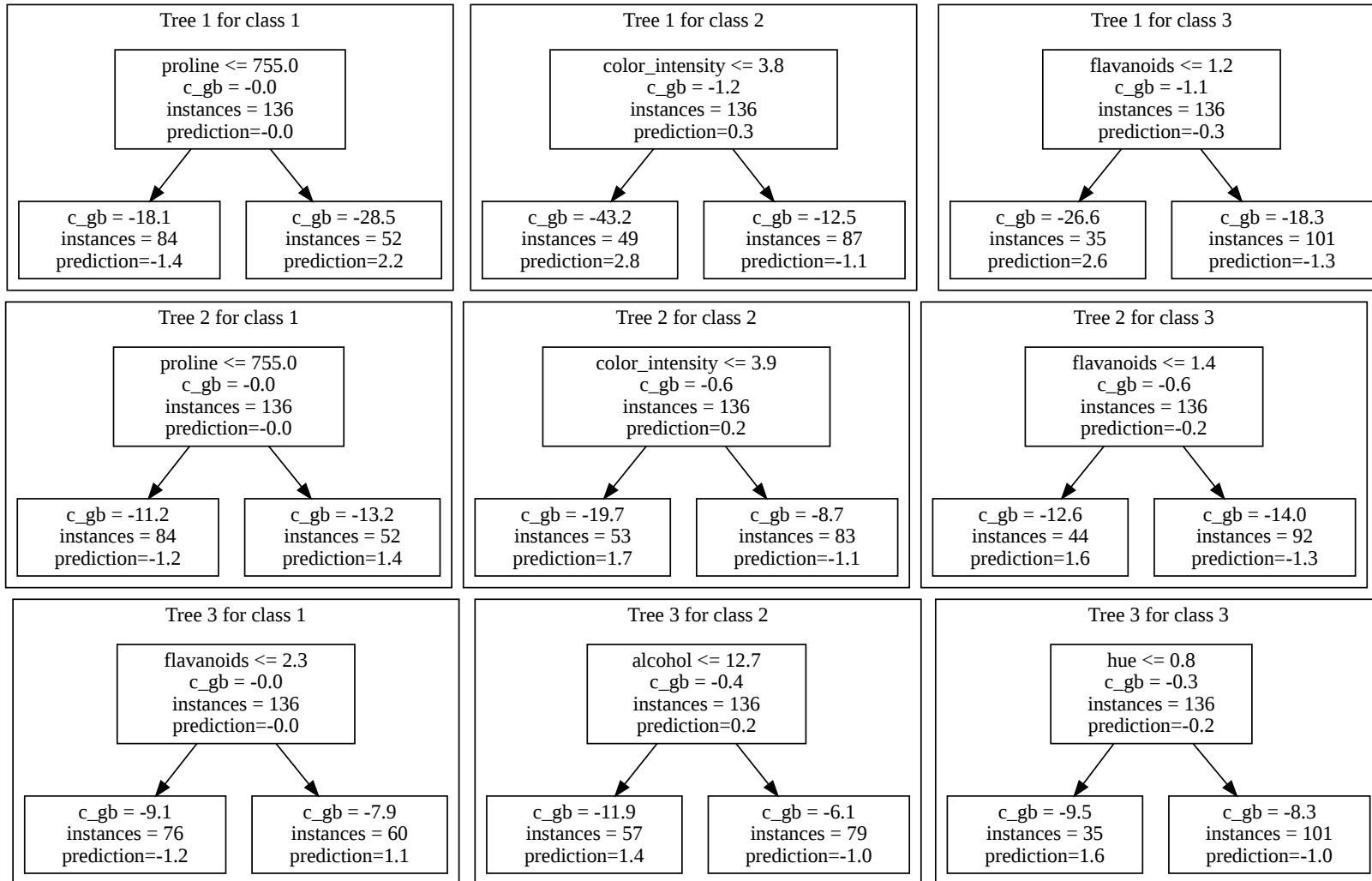
and the per-example loss for all  $K$  trees is defined analogously as

$$\ell(t_i, \mathbf{y}(\mathbf{x}_i)) = -\log\left(\text{softmax}(\mathbf{y}(\mathbf{x}_i))_{t_i}\right),$$

so that for a tree  $k$  at time  $t$ ,

$$\frac{\partial \ell(t_i, \mathbf{y}^{(t-1)}(\mathbf{x}_i))}{\partial \mathbf{y}^{(t-1)}(\mathbf{x}_i)_k} = \left(\text{softmax}(\mathbf{y}^{(t-1)}(\mathbf{x}_i)) - \mathbf{1}_{t_i}\right)_k.$$





## Playground

You can explore the [Gradient Boosting Trees playground](#) and [Gradient Boosting Trees explained](#).

## Implementations

Scikit-learn offers an implementation of gradient boosting decision trees, `sklearn.ensemble.GradientBoostingClassifier` for classification and `sklearn.ensemble.GradientBoostingRegressor` for regression.

- Furthermore, `sklearn.ensemble.HistGradientBoosting{Classifier/Regressor}` provide histogram-based splitting (which can be much faster for larger datasets – tens of thousands of examples and more) and efficient categorical feature splitting.

There are additional efficient implementations, capable of distributed processing of data larger than available memory (both offering also scikit-learn interface):

- XGBoost,
- LightGBM (which is the inspiration for the `HistGradientBoosting*` implementation).

This concludes the **supervised machine learning** part of our course.

We have encountered:

- parametric models
  - generalized linear models: perceptron algorithm, linear regression, logistic regression, multinomial (softmax) logistic regression
    - linear models, but manual feature engineering allows solving nonlinear problems
  - multilayer perceptron: nonlinear, perfect approximator – Universal approx. theorem
- nonparametric models
  - k-nearest neighbors
  - support vector machines not in this course, but in the state exam
- decision trees
  - can be both parametric or nonparametric depending on the constraints
- generative models
  - naive Bayes

When training a model for a new dataset, a good start is evaluating two models:

- an **MLP** with one/two hidden layers
  - works best for high-dimensional data (images, speech, text), where an individual single dimension (feature) does not convey much meaning; use pre-trained representation if possible;
- **gradient boosted decision tree**
  - works best for lower-dimensional data (“tabular data”), where the input features have interpretations on their own.

If there are only a few training examples with a lot of features, **naive Bayes** might also work well.

Finally, if your goal is to reach the highest possible performance and you have a lot of resources, definitely use **ensembling**.

After this lecture you should be able to

- Explain second-order optimization methods
- Implement gradient boosted decision trees for regression and classification
- Decide what supervised machine learning approach is suitable for particular problems