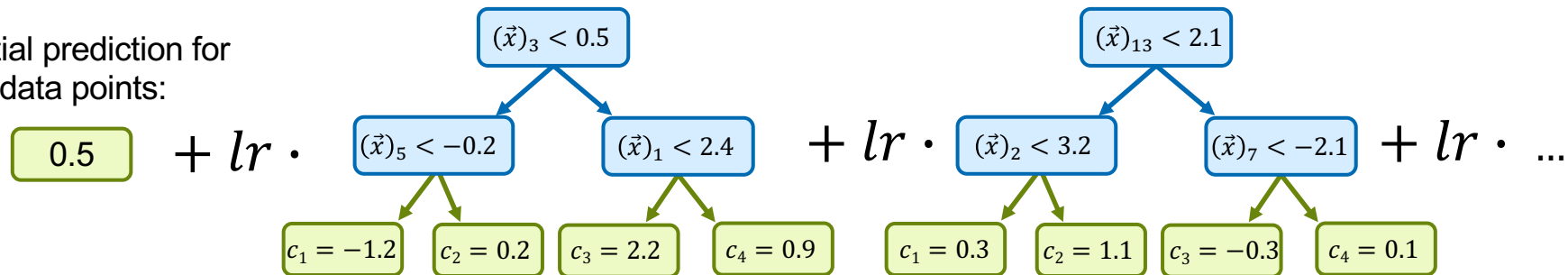


# Gradient Boosting Decision Tree Models

# What are Gradient Boosting Decision Tree Models?

- Let's assume we have input vectors  $\vec{x} \in \mathbb{R}^d$ . How do we get a prediction  $y \in \mathbb{R}$ ?

Initial prediction for  
all data points:



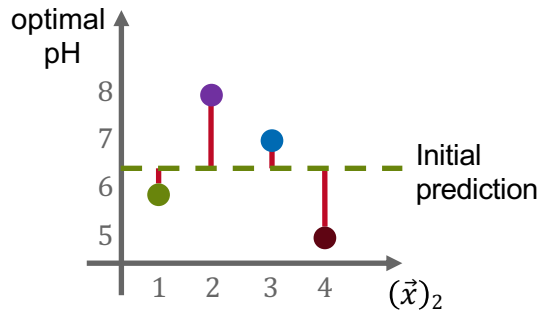
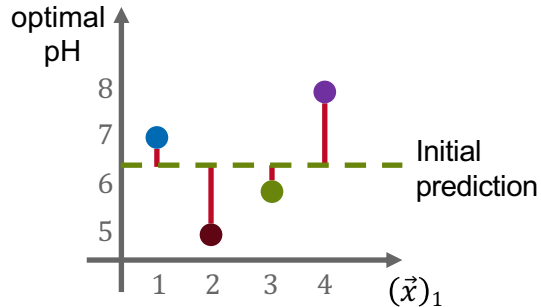
- Example:

- learning rate  $lr = 0.3$
- We are in leaf  $c_1$  for tree 1
- We are in leaf  $c_3$  for tree 2

$$y_{pred} = 0.5 + 0.3 \cdot (-1.2) + 0.3 \cdot (-0.3) + \dots$$

# Building Trees During Training

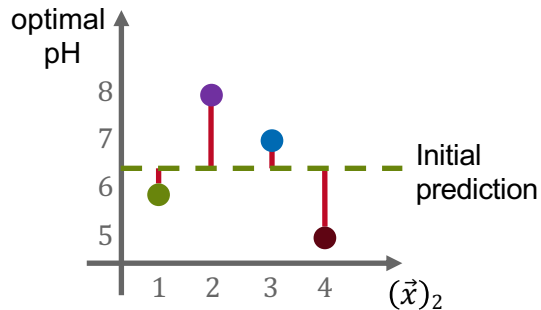
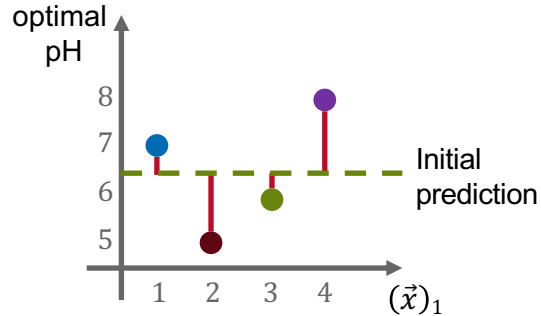
- Example training data:  $\mathcal{D} = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$ ,  $\vec{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}$  with  $n = 4, d = 2$ .



- Initial prediction for all data points: 6.5
- Find decision rule  $(\vec{x})_i < a$  that maximizes similarity of prediction errors  $e_i$  in two new subgroups
- Calculate similarity score for the two subgroups:  $S = \frac{(\sum_i e_i)^2}{n}$
- $S_{Root} = \frac{(0.5 - 1.5 - 0.5 + 1.5)^2}{4} = 0$

Data point	Prediction error $e_i$
<span style="color: blue;">●</span>	0.5
<span style="color: darkred;">●</span>	-1.5
<span style="color: green;">●</span>	-0.5
<span style="color: purple;">●</span>	1.5

# Building Trees During Training (2)



- Similarity score:

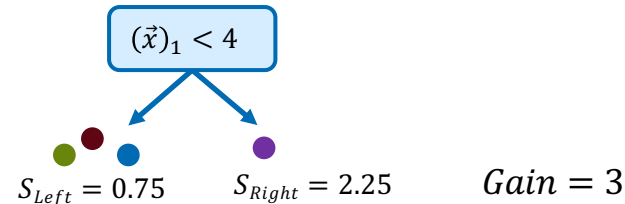
- $$S = \frac{(\sum_i e_i)^2}{n}$$
- $$S_{Root} = \frac{(0.5 - 1.5 - 0.5 + 1.5)^2}{4} = 0$$

- Aim:** Find decision rule  $(\vec{x})_i < a$  that splits the dataset in a left and a right group that maximizes the gain:

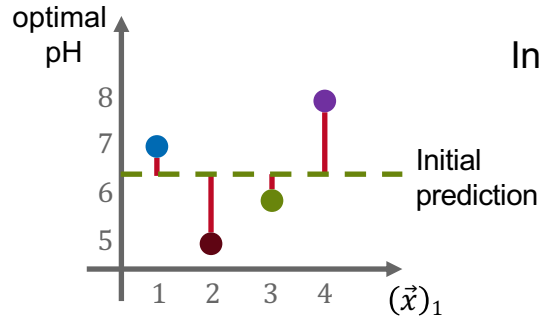
$$Gain = S_{Left} + S_{Right} - S_{Root}$$

- We test all possible values for  $i$  and  $a$  and select the combination with the highest *Gain*

Data point	Prediction error $e_i$
●	0.5
●	-1.5
●	-0.5
●	1.5

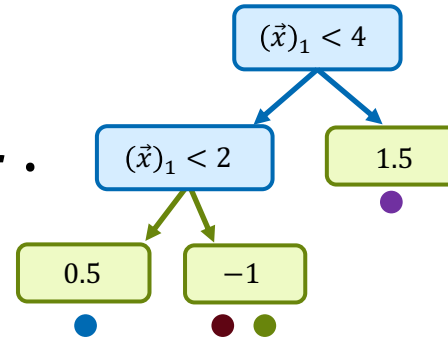


# Building Trees During Training (3)

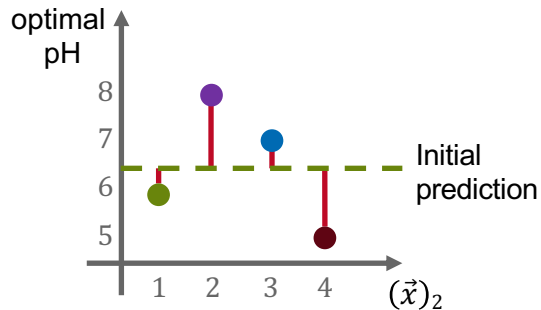


Initial prediction:

$$6.5 + lr \cdot$$



Data point	Prediction error $e_i$
<span style="color:blue">●</span>	0.5
<span style="color:red">●</span>	-1.5
<span style="color:green">●</span>	-0.5
<span style="color:purple">●</span>	1.5



For  $lr = 0.5$ :

Data point	New prediction	New prediction error
<span style="color:blue">●</span>	$6.5 + 0.5 \cdot 0.5 = 6.75$	0.25
<span style="color:red">●</span>	$6.5 + 0.5 \cdot (-1) = 6.00$	-1.00
<span style="color:green">●</span>	$6.5 + 0.5 \cdot (-1) = 6.00$	0.00
<span style="color:purple">●</span>	$6.5 + 0.5 \cdot 1.5 = 7.25$	0.75

# Most important hyperparameters

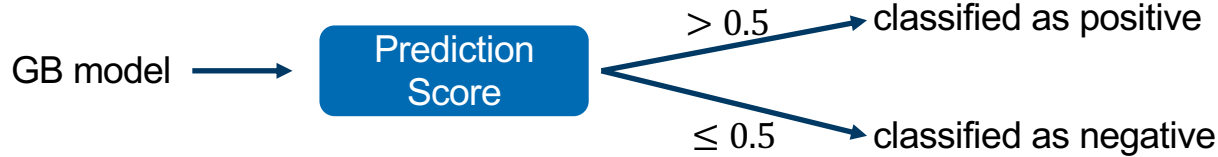
- *learning\_rate*
  - Specifies how much weight the prediction of each new tree gets
  - Range: (0,1]
  - Default: 0.3
- *reg\_lambda* and *reg\_alpha*
  - Regularization coefficients that
    - penalize large predictions in the leaves
    - weaken the influence of branches that contain only few data points
  - Range:  $[0, \infty]$  for both
  - Default: 1 for *reg\_lambda* and 0 for *reg\_alpha*
- *max\_depth*
  - Defines to maximum allowed depth for each tree.
  - Default: 6

# Most important hyperparameters (2)

- *n\_estimators*
  - Defines the numbers of tree that will be build in the gradient boosting model
  - Range: Any positive integer
- *min\_child\_weight*
  - Tells the model when to stop splitting tree nodes when:
    - number of data points in this node is too small (regression task)
    - most data points in this node belong to one class (classification task)
  - Range: Any positive integer
  - Default: 1
- For additional hyperparameters see the XGBoost documentation:  
<https://xgboost.readthedocs.io/en/stable/parameter.html>

# Gradient Boosting for binary Classification

- Gradient Boosting for binary classification tasks works very similar as for regression tasks
- Model output is a prediction score between 0 and 1:



- With each tree we try to predict the prediction error:
  - $Error = True\ class\ label - Prediction\ score$
- Differences to regression prediction task:
  - Calculation of similarity scores that are required for splitting tree nodes
  - Calculating the prediction values in the leaves



# Gradient Boosting for binary Classification (2)

- The model predicts a  $\log(\text{odds})$  score  $s$

- $\text{odds} = \frac{\text{likelihood of positive class}}{\text{likelihood of negative class}} = \frac{p}{1-p}$

- Converting  $\log(\text{odds})$  score  $s$  to a probability score by applying the logistic function

$$\text{sigmoid}(s) = \frac{e^s}{1 + e^s} = p$$

- Initial prediction:  $\log(\text{odds})$  based on the training data positive:

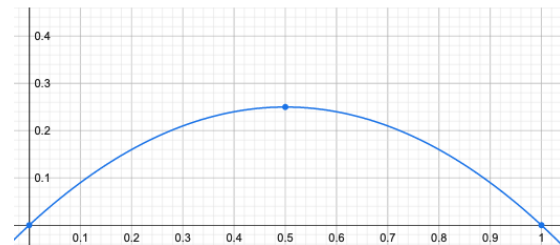
$$\frac{\# \text{positive class labels}}{\# \text{negative class labels}}$$

- With each tree we try improve the prediction error:

- $\text{Error} = \text{True class label} - \text{predicted probability}$

- Computing predictions for all data points in the same leaf:

$$\frac{\sum_i \text{Error}_i}{\sum_i \text{predicted probability}_i \times (1 - \text{predicted probability}_i)}$$



## Python

## Python

## Python

## Python

# hhu.de

# Plotting XGBoost trees

```
from xgboost import plot_tree
import matplotlib.pyplot as plt
plot_tree(model, num_trees=0)
plot_tree(model, num_trees=1)
plt.show()
```

✓ 0.5s

Python

