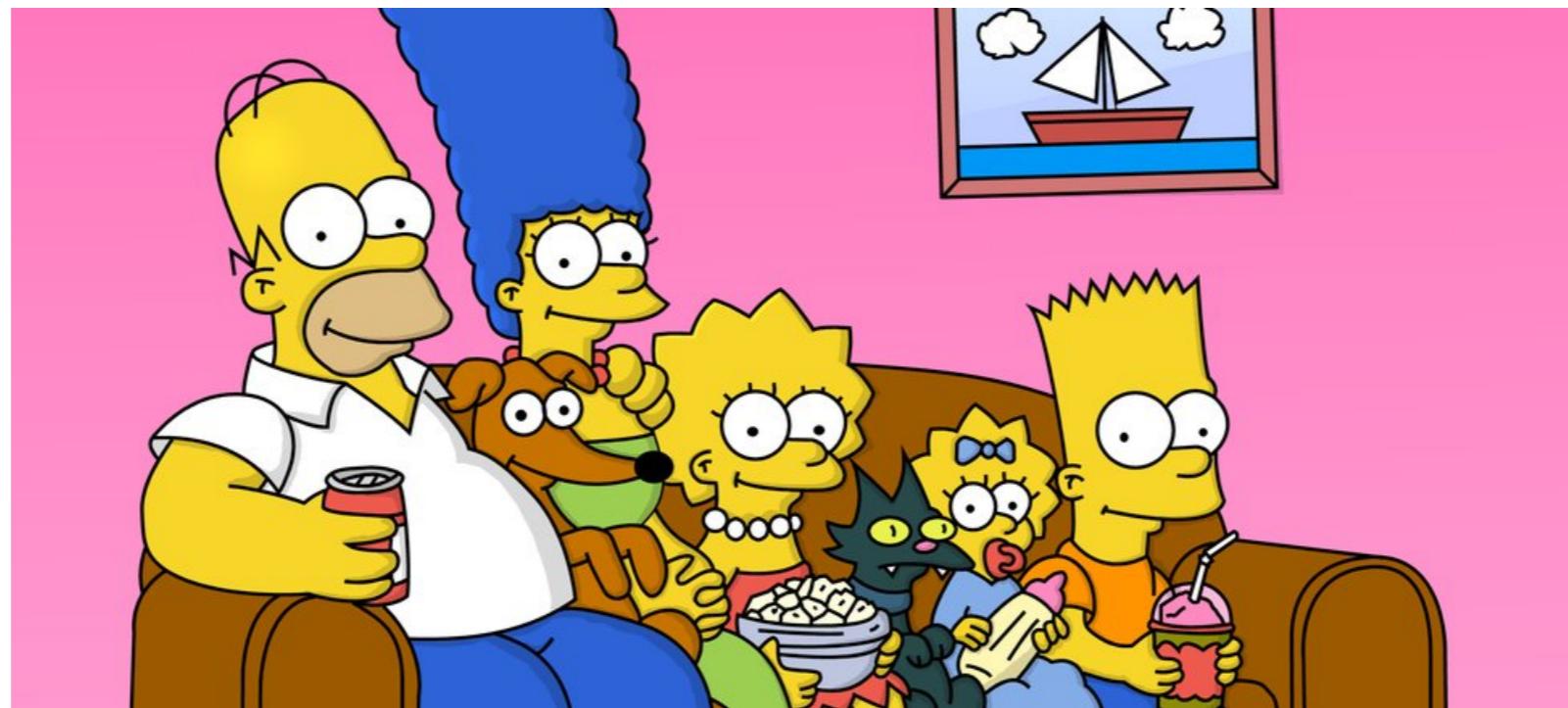


The Simpsons Challenge

Felicia Burtscher, Thomas Hadler, Hanna Wulkow



Overview

- Problem: Image Classification
- Computer Vision Features
- Decision Trees & Random Forests
- Results & Evaluation

Image Classification

- The Simpsons data set: 20 000 images, 20 labels

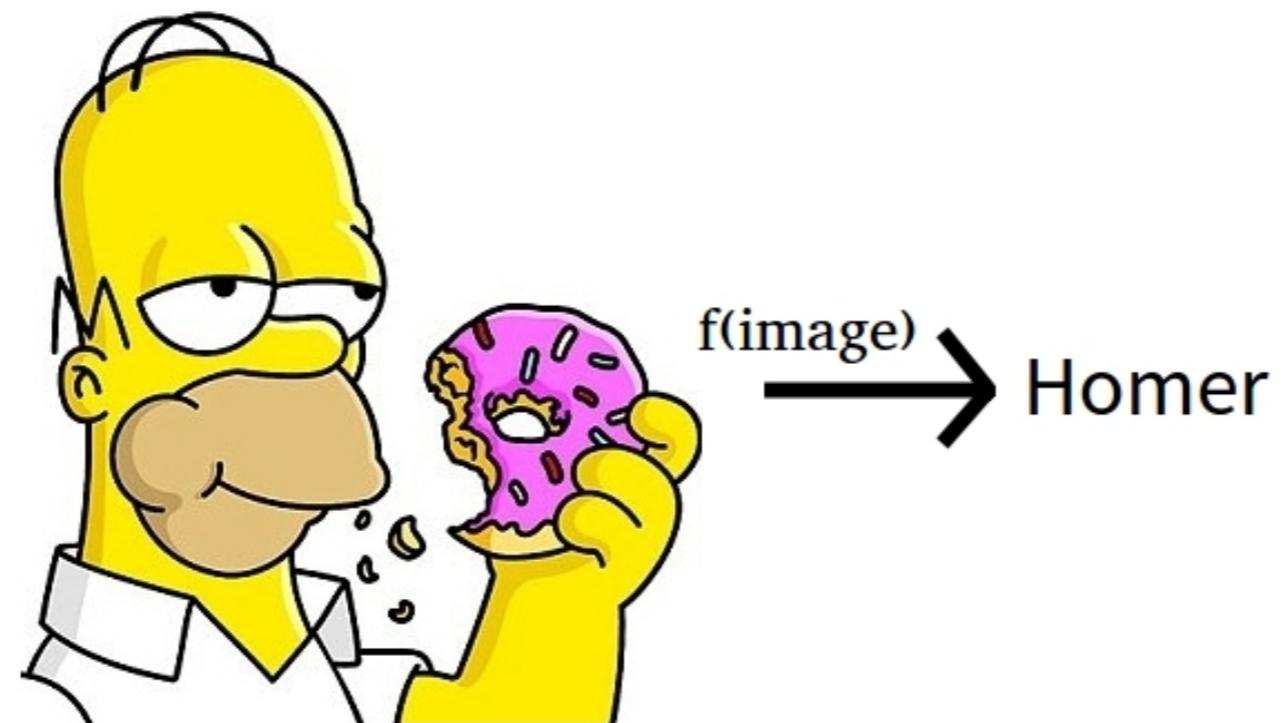
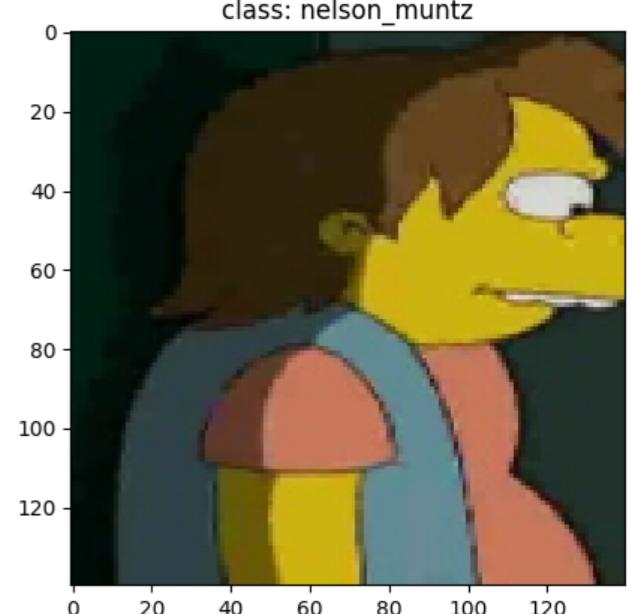


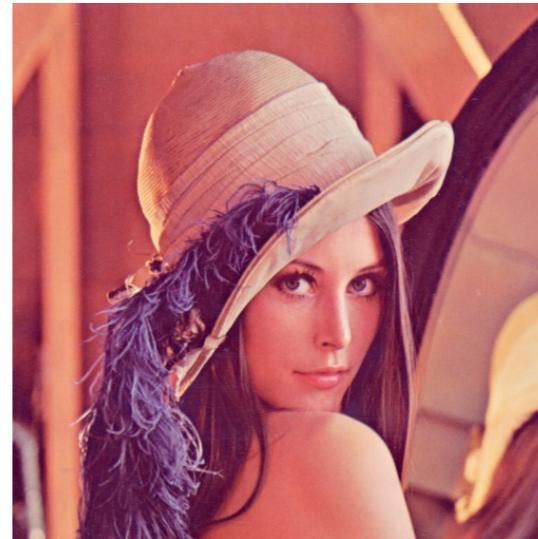
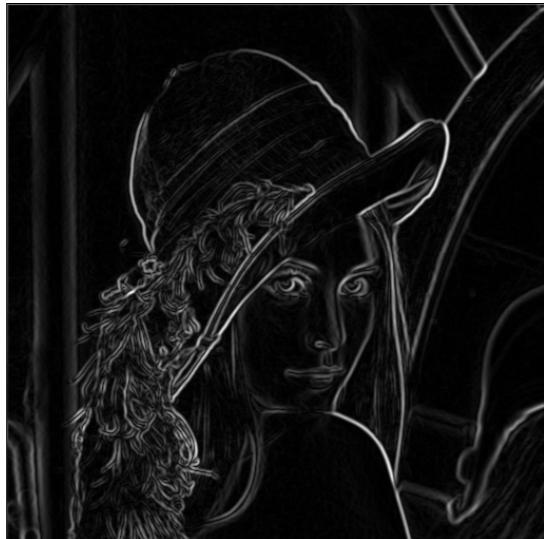
Image Classification

- Input image
 - Images are matrices (or tensors)
- Simpsons (red, green, blue):
 - $110 \times 110 \times 3 \rightarrow 33\ 000$ dim vector
 - Entries in $\{1, \dots, 255\}$
- Difficult problem
 - 20 000 images insufficient to approx. 33000 dimensions



Computer Vision Features

- Features are “compact information”
 - In computer vision spatially localized information
 - Example: edge detection on “Lena”



- Relevant information preserved
- Edges are local information

Computer Vision Features

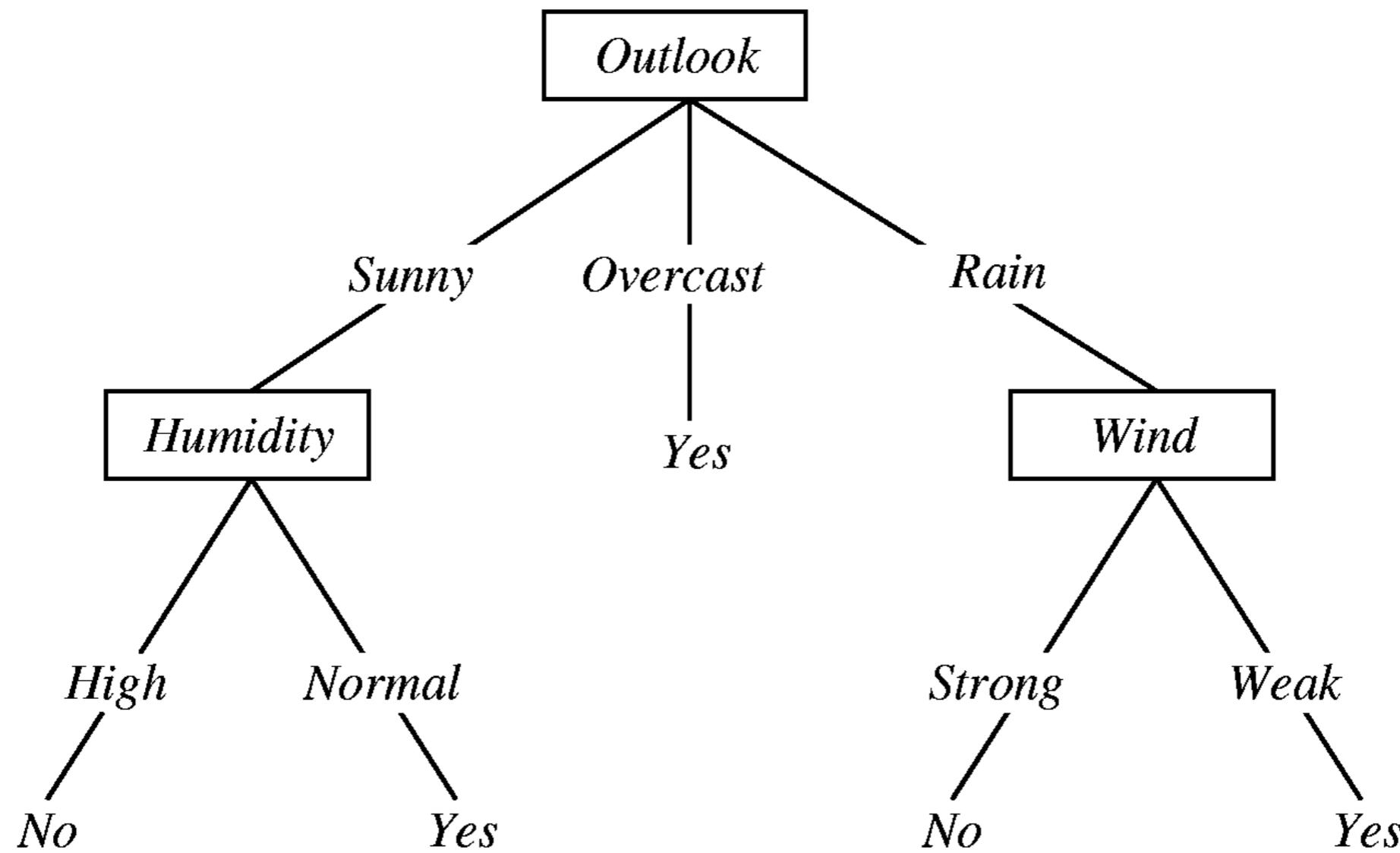
- Hand-crafted features
 - Histogram of oriented gradients (by Dalal and Triggs)
 - Means values of large patches
- HOG
 - Input image
 - Compute gradients
 - Weighted vote into spatial and orientation cells
 - Feature image
- Similar to pooled gradients
 - With smoothing interpolation

HOG Features

- Input image
- Compute gradients
- Weighted vote into spatial and orientation cells
- Feature image



Decision Trees



Random Forest

- Adaptation of decision trees

Basic principles:

- Grows various trees
- Chooses a random collection of samples from the training data set
- Chooses a random set of attributes from the original data set to test for splitting

Splitting Criterion

Information gain

- Entropy
- Gini index

Entropy

- Measures the uncertainty of a class in a subset of examples S
- (Binary) entropy is defined as

$$H(S) = -p_+ \cdot \log_2 p_+ - p_- \cdot \log_2 p_-$$

S ...the subset of training examples

p_+ , p_- ...the positive and negative examples in S

Information gain

- Goal: items in pure sets
- Expected drop in entropy after the split:

$$Gain(S, a) = H(S) - \sum_{V \in values(a)} \frac{|S_V|}{|S|} H(S_V)$$

V ... possible values of a

S ... set of examples (X)

S_V ... subset where X_a = V

Results

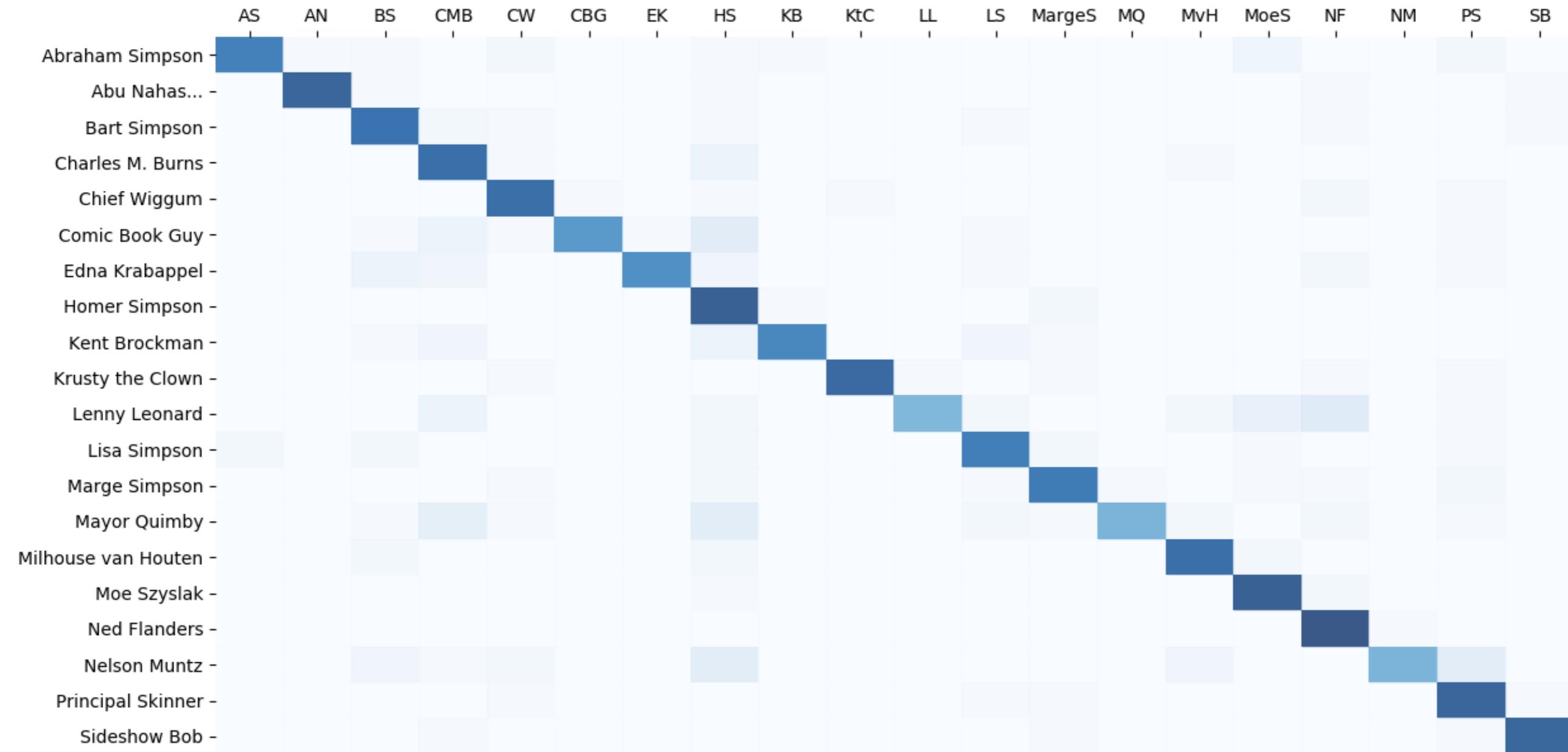
Accuracy & Confusion Matrix

- Accuracy:
 - $(\text{True Pos} + \text{True Neg}) / (\text{Pos} + \text{Neg})$
- Confusion Matrix C:
 - $C(i,j)$: how often character i is classified as character j

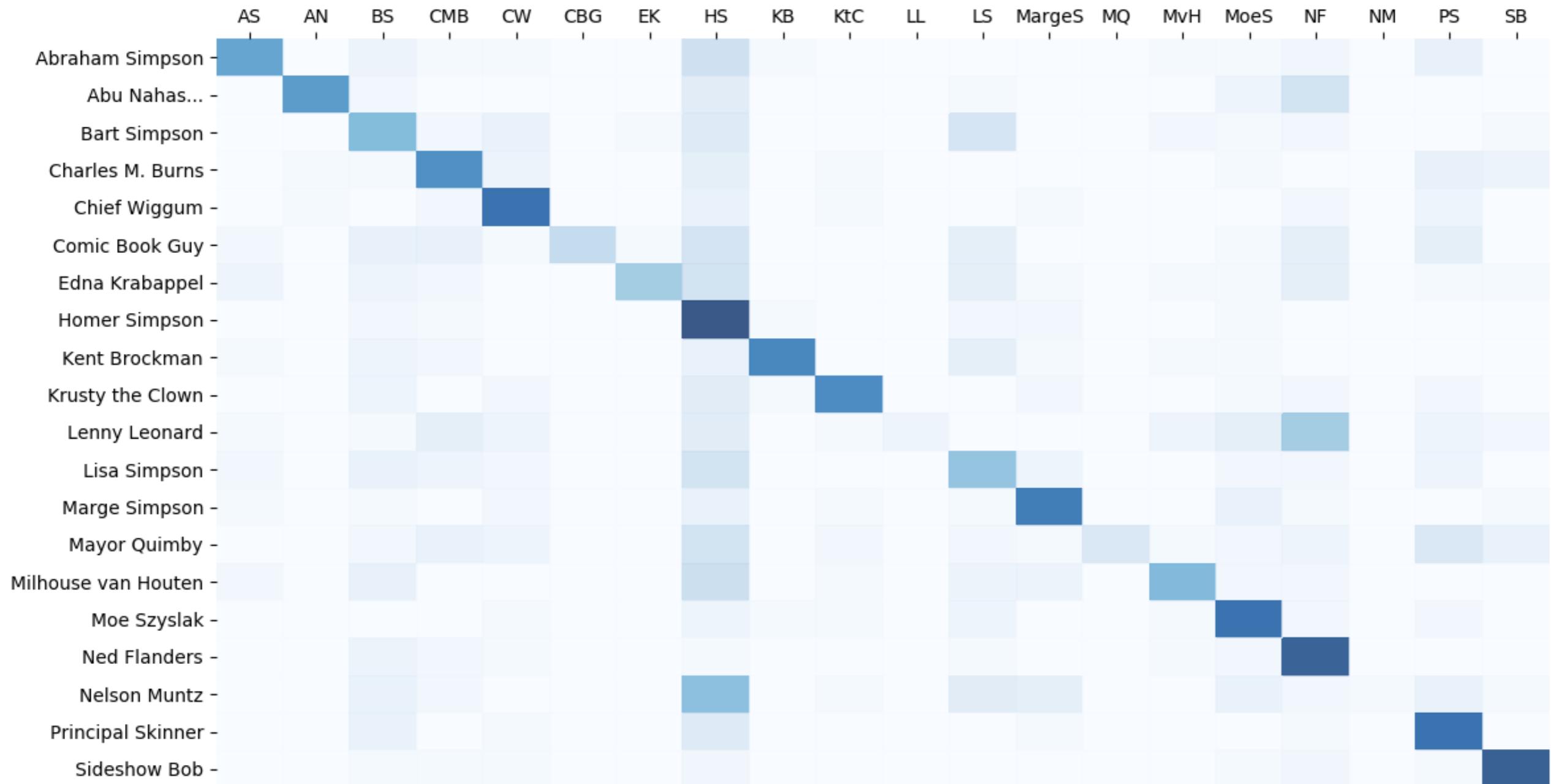
Accuracy of Test Data

	1	5	10	15	20	50
1	0.07	0.07	0.08	0.06	0.07	0.07
5	0.17	0.25	0.26	0.27	0.27	0.27
10	0.25	0.43	0.5	0.53	0.54	0.58
15	0.33	0.55	0.63	0.66	0.72	0.76
20	0.31	0.56	0.66	0.74	0.76	0.83
50	0.36	0.61	0.68	0.75	0.79	0.85

Confusion Matrix: 50 trees with depth 50



Confusion Matrix: 15 trees with depth 10



The best



The worst



Confusion Matrix: 15 trees with depth 10

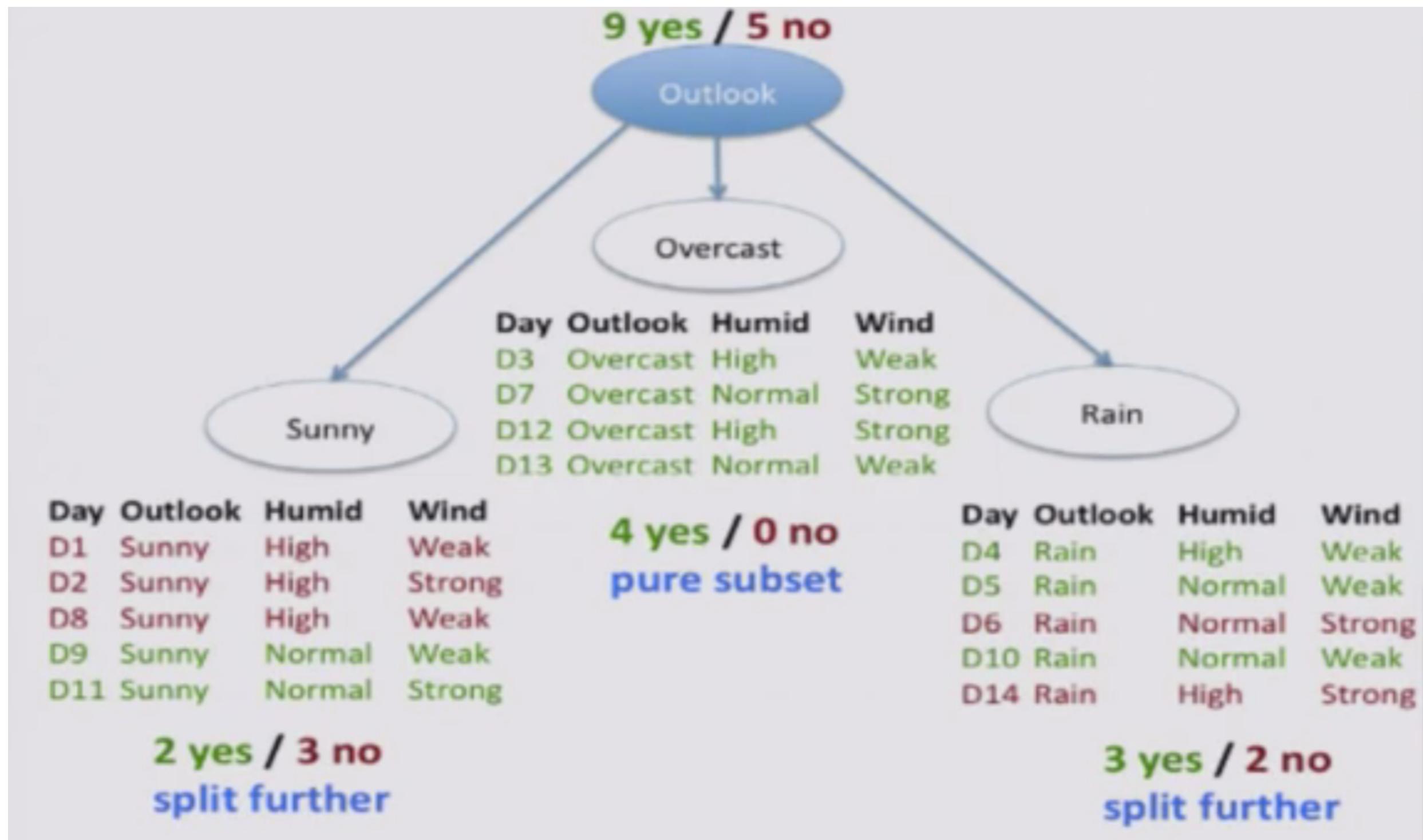
- more than 50 % of characters are correctly classified more than 50 % of the time
- worst: Nelson, only 2 %!
- best: Homer Simpson, 82 %
- average correct classification: 50.25 %

Homer Simpson

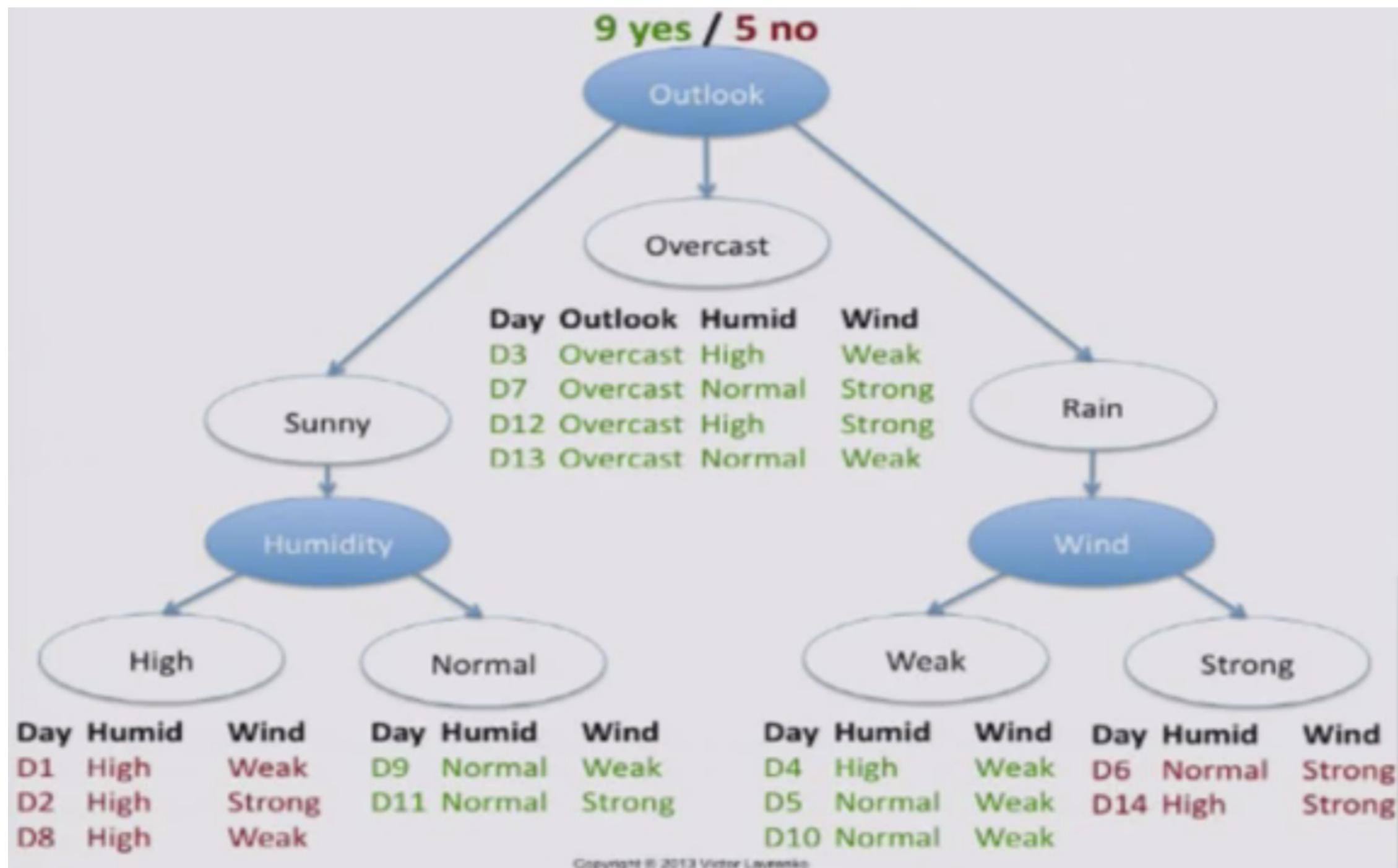
- Why are all characters classified as Homer so often?

The End

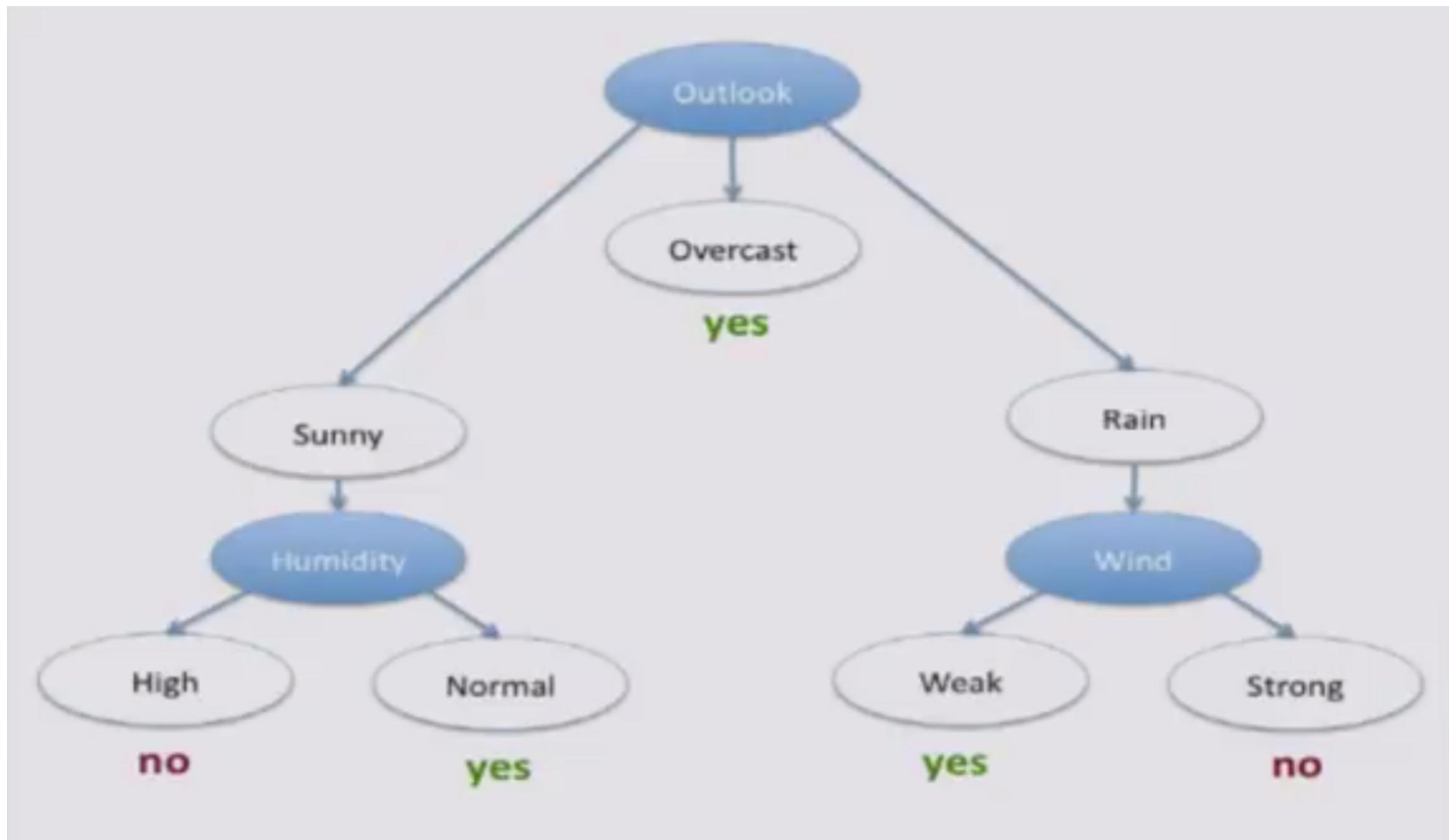
Decision Trees



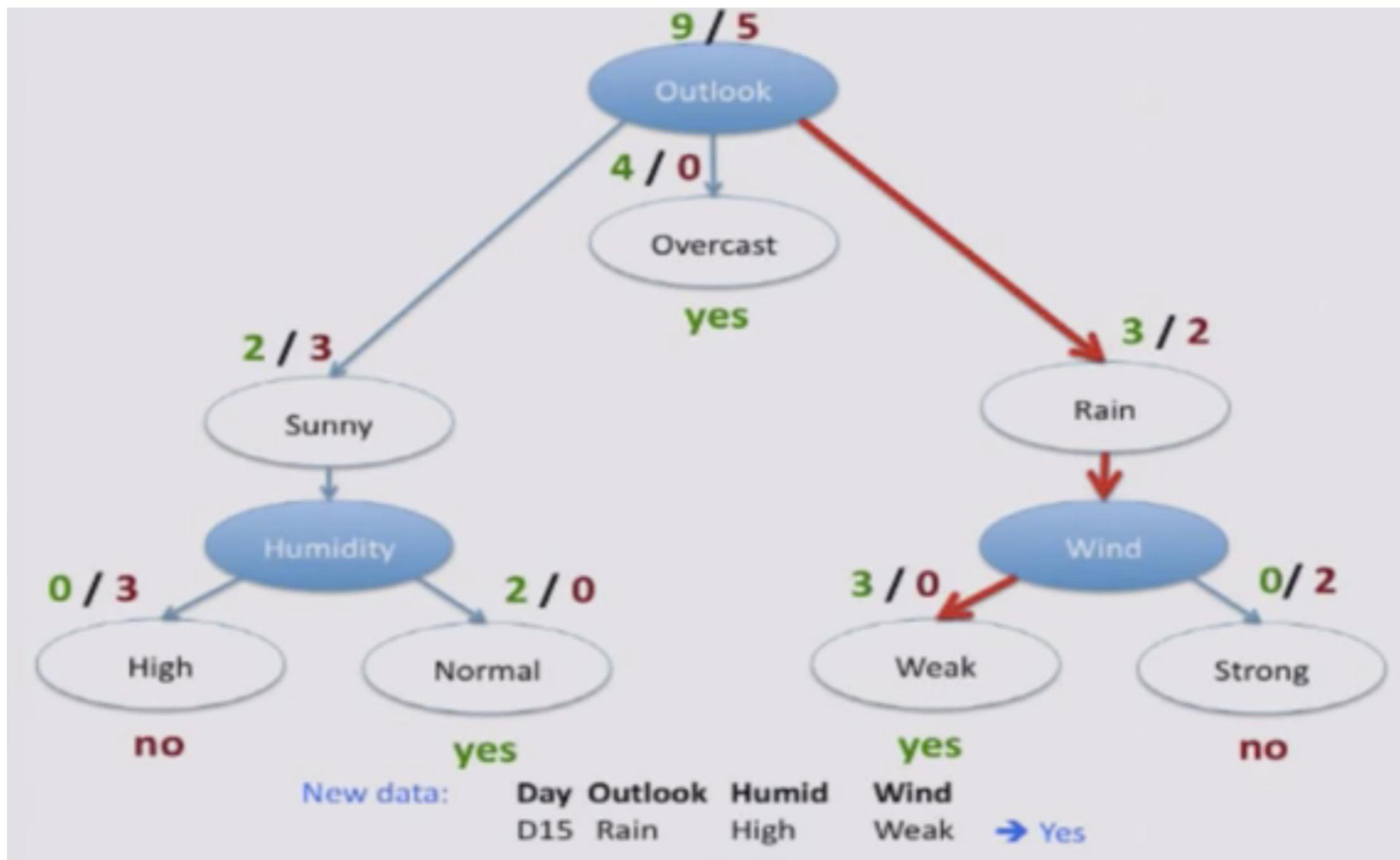
Decision Trees



Decision Trees



Decision Trees



New data: Day Outlook Humid Wind
D15 Rain High Weak → Yes