# Trees & Forests

02/18/19

Andreas C. Müller
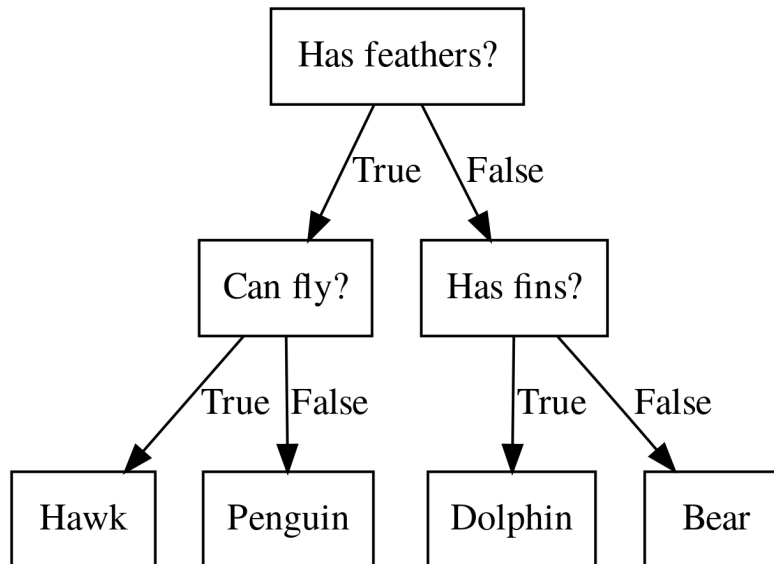
(Adapted and modified for CC 6021236 @ PCC/Ciencias/UCV by

Eugenio Scalise, September 2019)

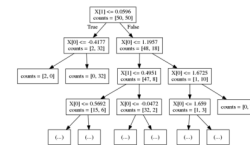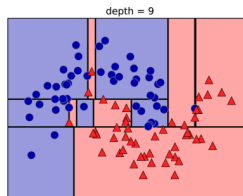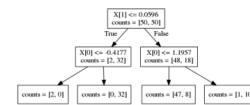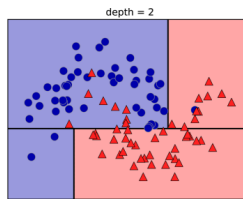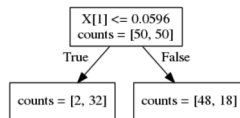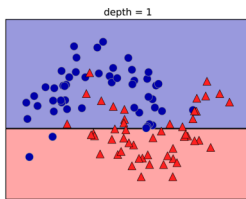# Why Trees?

- Very powerful modeling method – non-linear!

- Doesn't care about scaling of distribution of data!

- "Interpretable"

- Basis of very powerful models!

# Decision Trees for Classification

# Idea: series of binary questions

# Building Trees



depth = 1

| | |
|---|---|
| X[1] <= 0.0596 | |
| counts = [50, 50] | |
| True / | \ False |

| counts = [2, 32] | counts = [48, 18] |



depth = 2

| | |
|---|---|
| X[1] <= 0.0596 | |
| counts = [50, 50] | |
| True / | \ False |

| X[0] <= -0.4177 | X[0] <= 1.1957 |
| counts = [2, 32] | counts = [48, 18] |

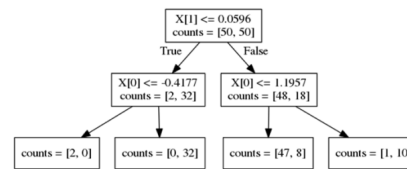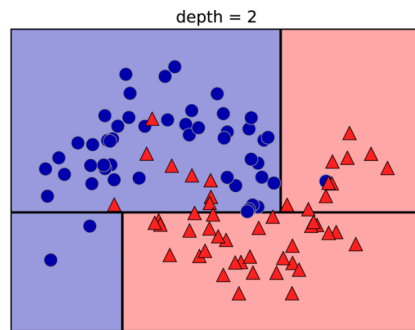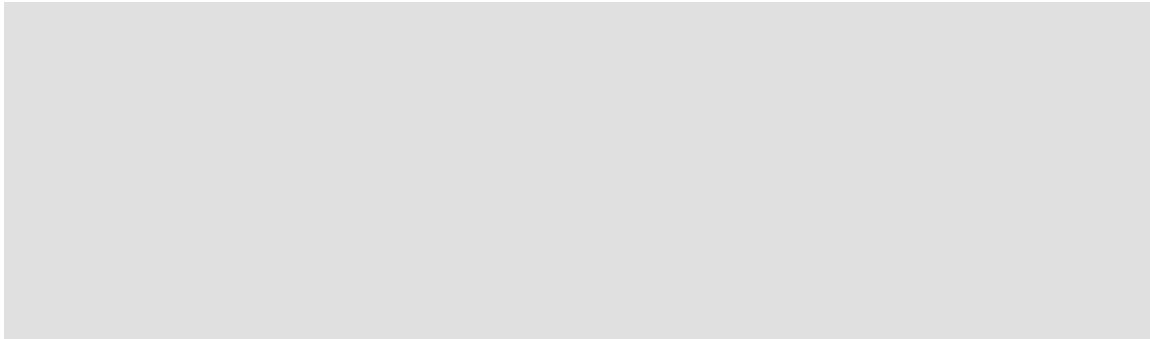| counts = [2, 0] | counts = [0, 32] | counts = [47, 8] | counts = [1, 10] |

Continuous features:

- "questions" are thresholds on single features.
- Minimize impurity



depth = 9

# Prediction

# Visualizing trees with sklearn

# Visualizing trees with sklearn

```
                        worst perimeter <= 106.1
                           entropy = 0.953
                           samples = 426
                          value = [159, 267]
              ┌──────────────────┴──────────────────┐
  worst concave points <= 0.134              mean concave points <= 0.064
        entropy = 0.254                            entropy = 0.511
        samples = 259                              samples = 167
       value = [11, 248]                          value = [148, 19]
    ┌───────┴───────┐                          ┌───────┴───────┐
entropy = 0.039   entropy = 0.998        entropy = 0.968   entropy = 0.0
samples = 240     samples = 19           samples = 48      samples = 119
value = [1, 239]  value = [10, 9]        value = [29, 19]  value = [119, 0]
```
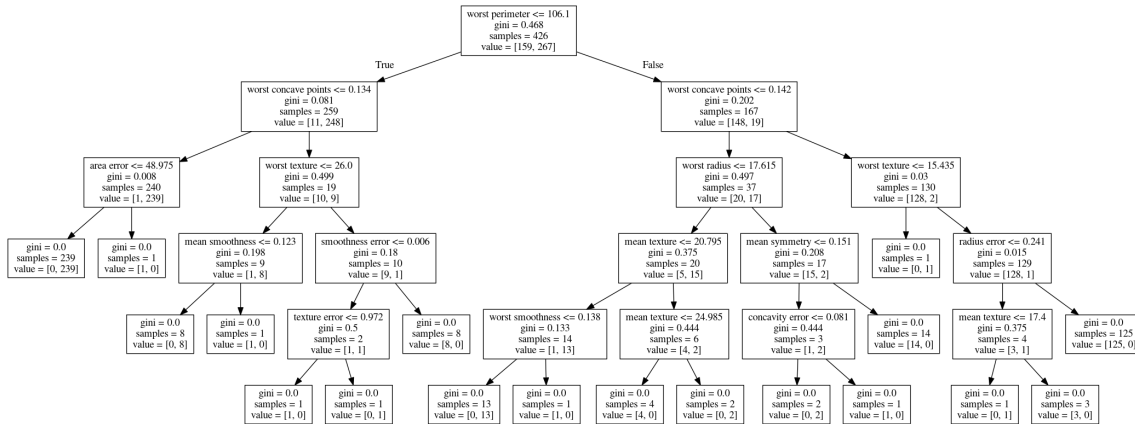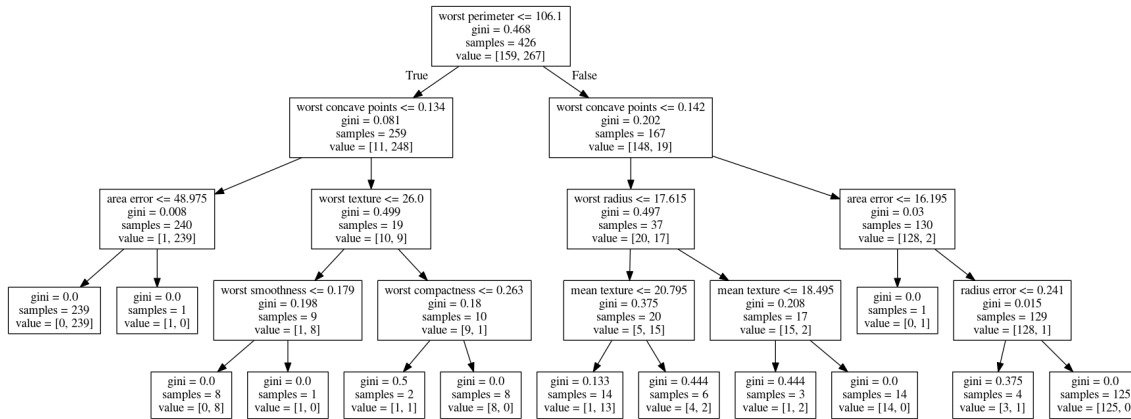
# Parameter Tuning

- Limit tree size (pick one, maybe two):

    - max_depth

    - max_leaf_nodes

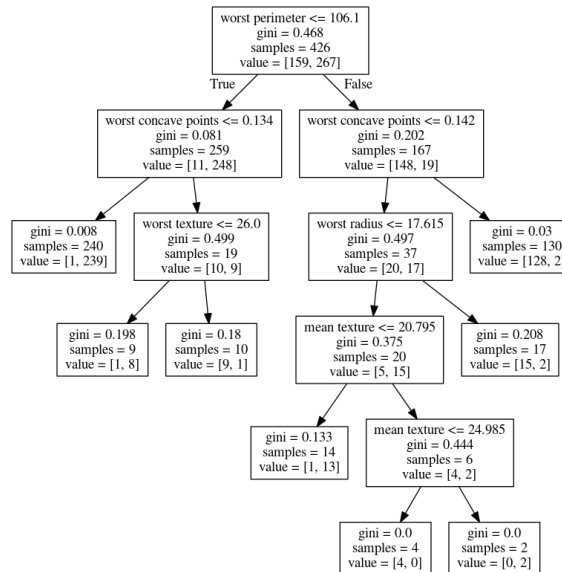    - min_samples_split

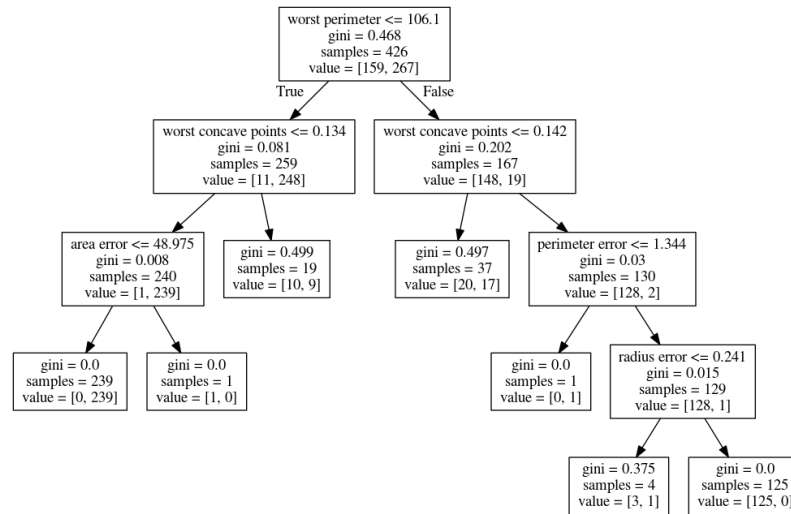    - min_impurity_decrease

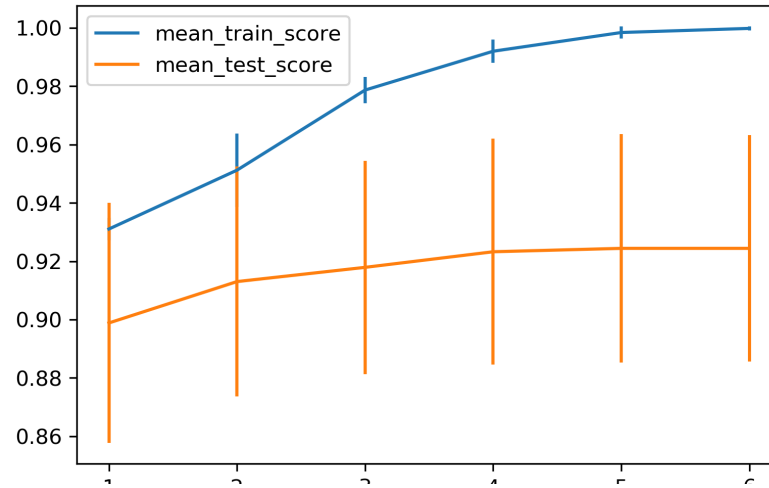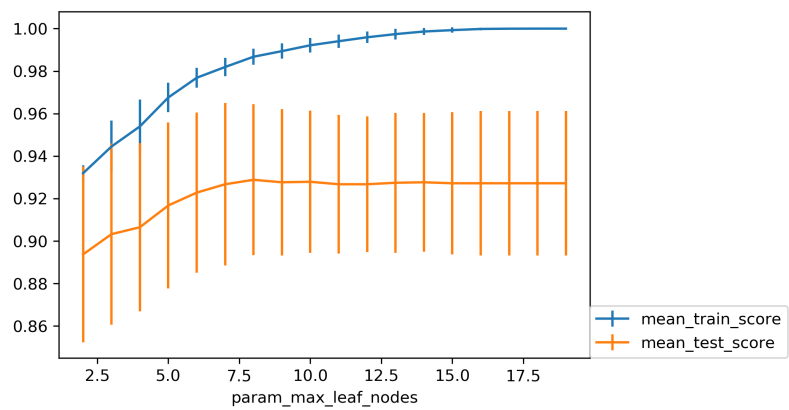    - ...

# No pruning

# max_depth = 4

# max_leaf_nodes = 8

worst perimeter <= 106.1
gini = 0.468
samples = 426
value = [159, 267]

True — False

worst concave points <= 0.134
gini = 0.081
samples = 259
value = [11, 248]

worst concave points <= 0.142
gini = 0.202
samples = 167
value = [148, 19]

gini = 0.008
samples = 240
value = [1, 239]

worst texture <= 26.0
gini = 0.499
samples = 19
value = [10, 9]

worst radius <= 17.615
gini = 0.497
samples = 37
value = [20, 17]

gini = 0.03
samples = 130
value = [128, 2]

gini = 0.198
samples = 9
value = [1, 8]

gini = 0.18
samples = 10
value = [9, 1]

mean texture <= 20.795
gini = 0.375
samples = 20
value = [5, 15]

gini = 0.208
samples = 17
value = [15, 2]

gini = 0.133
samples = 14
value = [1, 13]

mean texture <= 24.985
gini = 0.444
samples = 6
value = [4, 2]

gini = 0.0
samples = 4
value = [4, 0]

gini = 0.0
samples = 2
value = [0, 2]

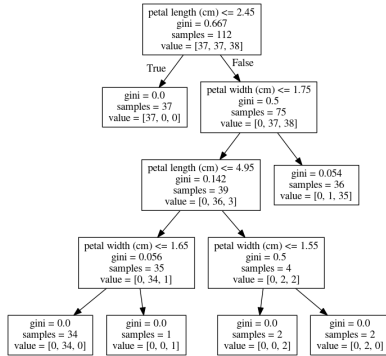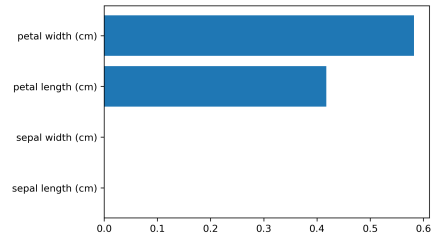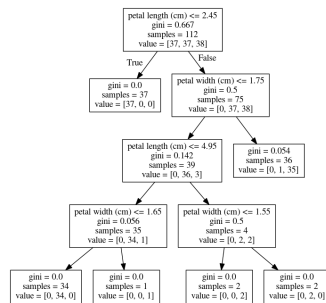# min_samples_split = 50
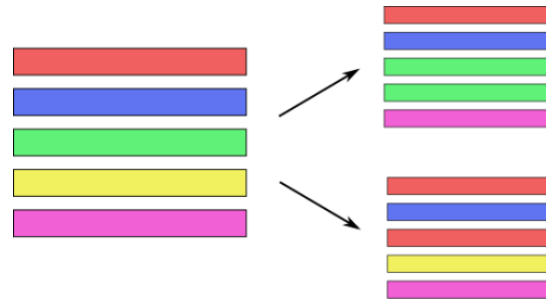
# Instability

# Feature importance

# Ensemble Models
# (Random Forests)
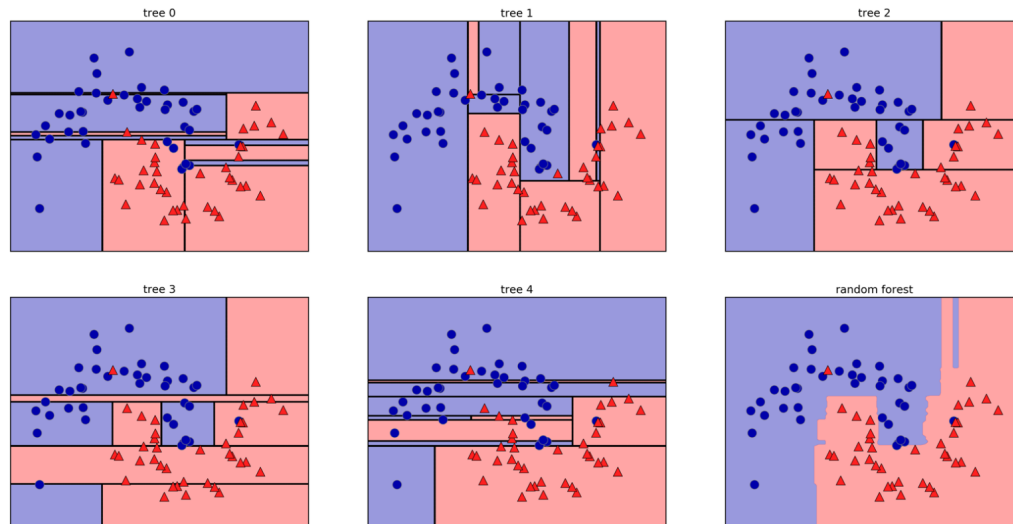
# Poor man's ensembles

- Build different models
- Average the result
- More models are better – if they are not correlated.
- Also works with neural networks
- You can average any models as long as they provide calibrated ("good") probabilities.
- Scikit-learn: VotingClassifier

# Bagging (Bootstrap AGGregation)

- Generic way to build "slightly different" models
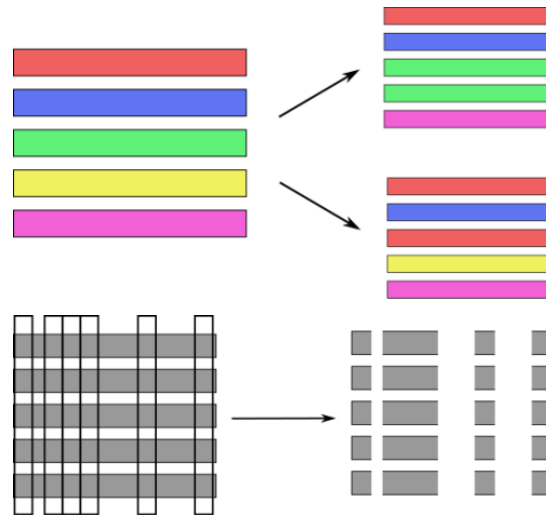
- BaggingClassifier, BaggingRegressor

# Random Forests

# Randomize in two ways

- For each tree:
  - Pick bootstrap sample of data

- For each split:
  - Pick random sample of features
- More trees are always better

# Tuning Random Forests

- Main parameter: max_features

  - around sqrt(n_features) for classification

  - Around n_features for regression

- n_estimators > 100

- max_depth, max_leaf_nodes, min_samples_split again