

W4995 Applied Machine Learning

Trees, Forests & Ensembles

02/18/19

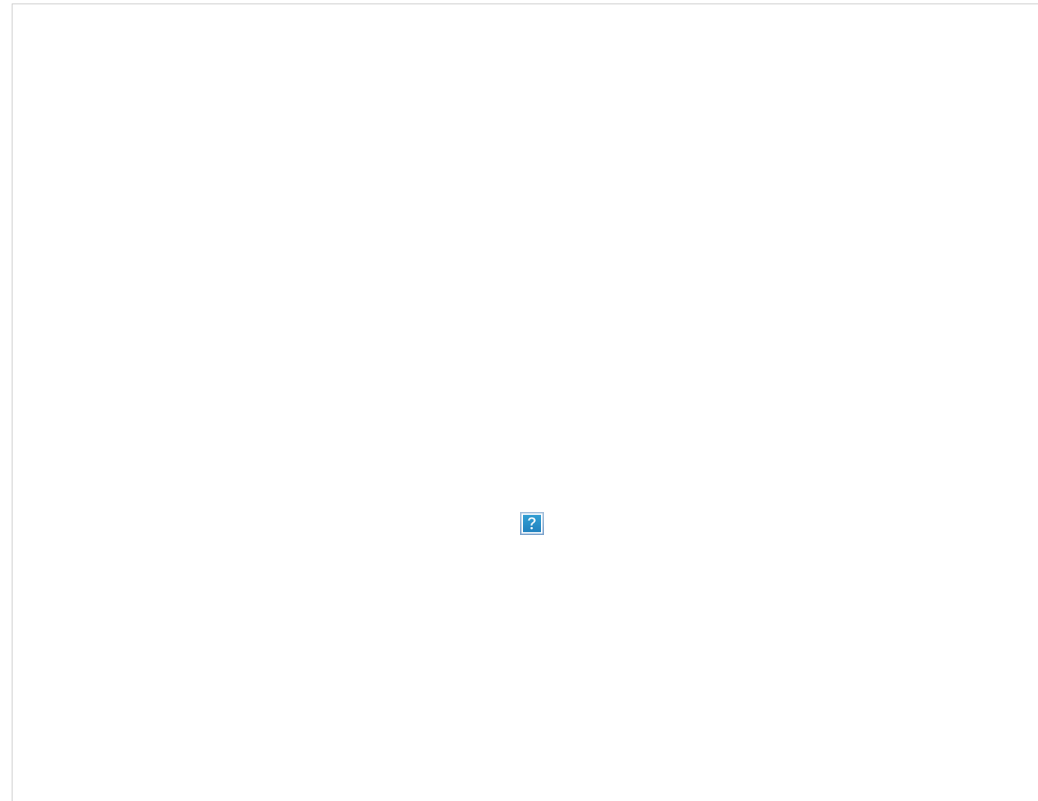
Andreas C. Müller

1 / 43

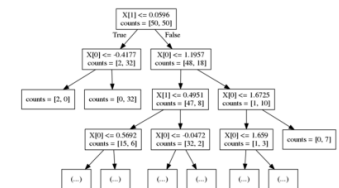
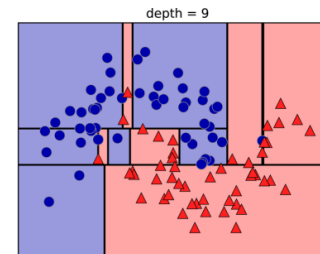
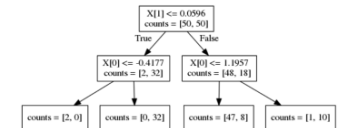
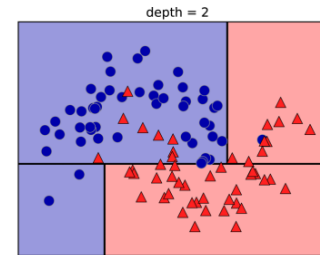
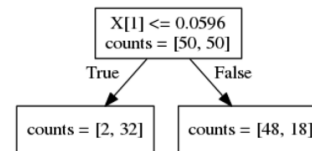
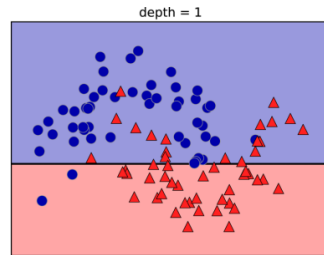
Why Trees?

Decision Trees for Classification

Idea: series of binary questions



Building Trees



Continuous features:

- “questions” are thresholds on single features.
- Minimize impurity

Criteria (for classification)

- Gini Index:

$$H_{\text{gini}}(X_m) = \sum_{k \in \mathcal{Y}} p_{mk}(1 - p_{mk})$$

- Cross-Entropy:

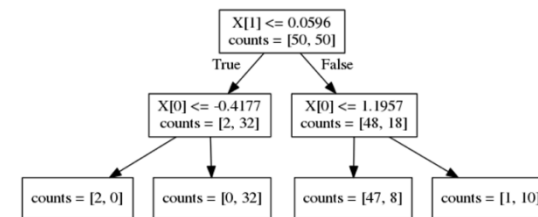
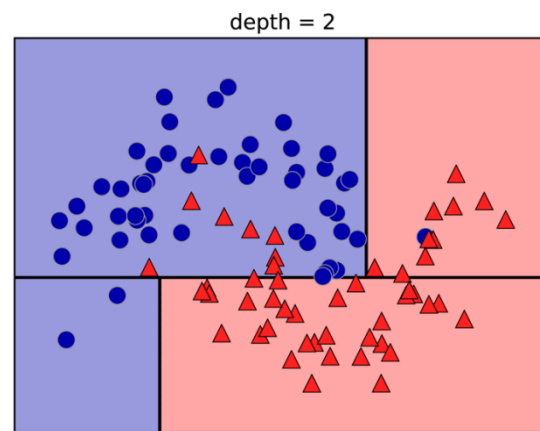
$$H_{\text{CE}}(X_m) = - \sum_{k \in \mathcal{Y}} p_{mk} \log(p_{mk})$$

X_m observations in node m

\mathcal{Y} classes

p_m distribution over classes in node m

Prediction



Regression trees

$$\text{Prediction: } \bar{y}_m = \frac{1}{N_m} \sum_{i \in N_m} y_i$$

Mean Squared Error:

$$H(X_m) = \frac{1}{N_m} \sum_{i \in N_m} (y_i - \bar{y}_m)^2$$

Mean Absolute Error:

$$H(X_m) = \frac{1}{N} \sum |y_i - \bar{y}_m|$$

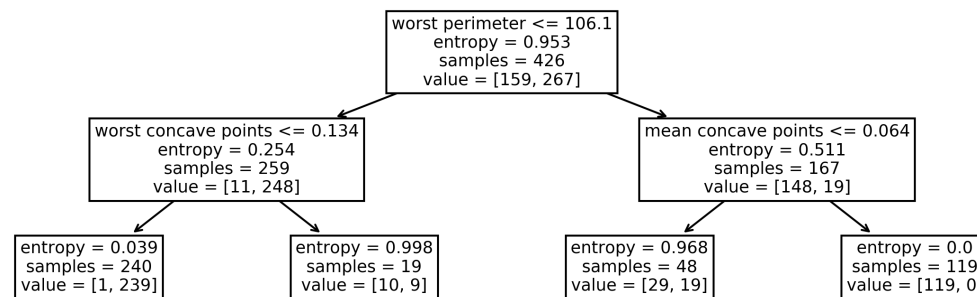
Visualizing trees with sklearn

```
from sklearn.datasets import load_breast_cancer
cancer = load_breast_cancer()
X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target,
                                                    stratify=cancer.target,
                                                    random_state=0)

from sklearn.tree import DecisionTreeClassifier, export_graphviz
tree = DecisionTreeClassifier(max_depth=2)
tree.fit(X_train, y_train)
```

Visualizing trees with sklearn

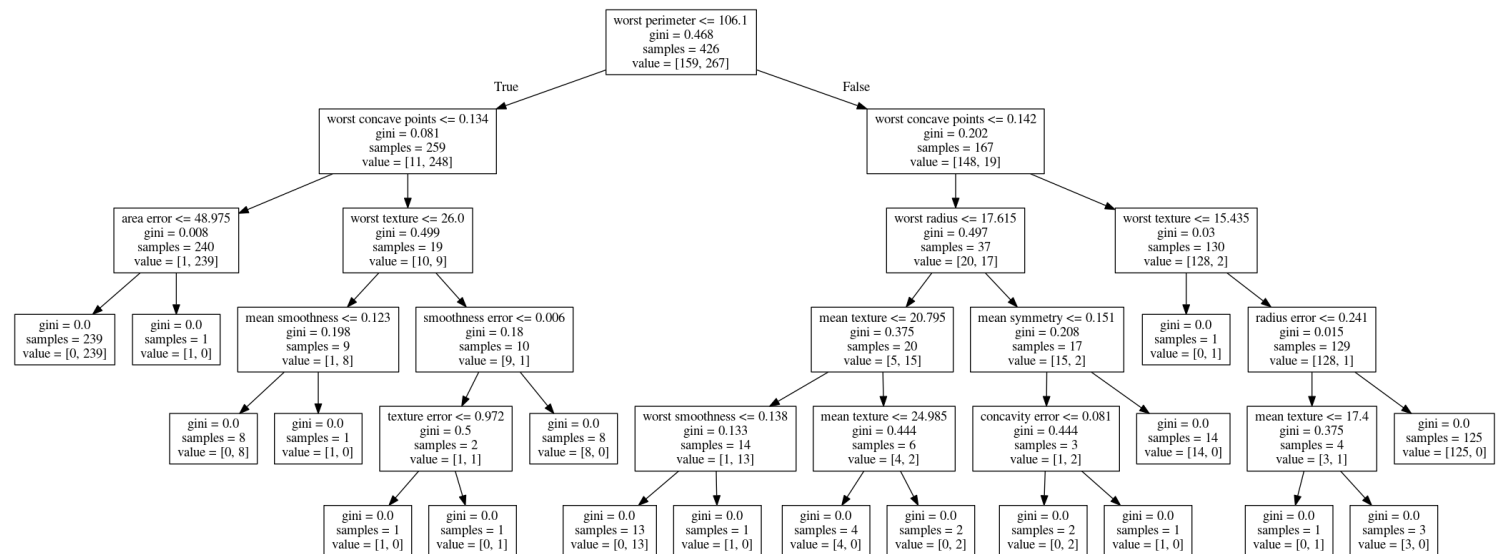
```
from sklearn.tree import plot_tree  
tree_dot = plot_tree(tree, feature_names=cancer.feature_names)
```



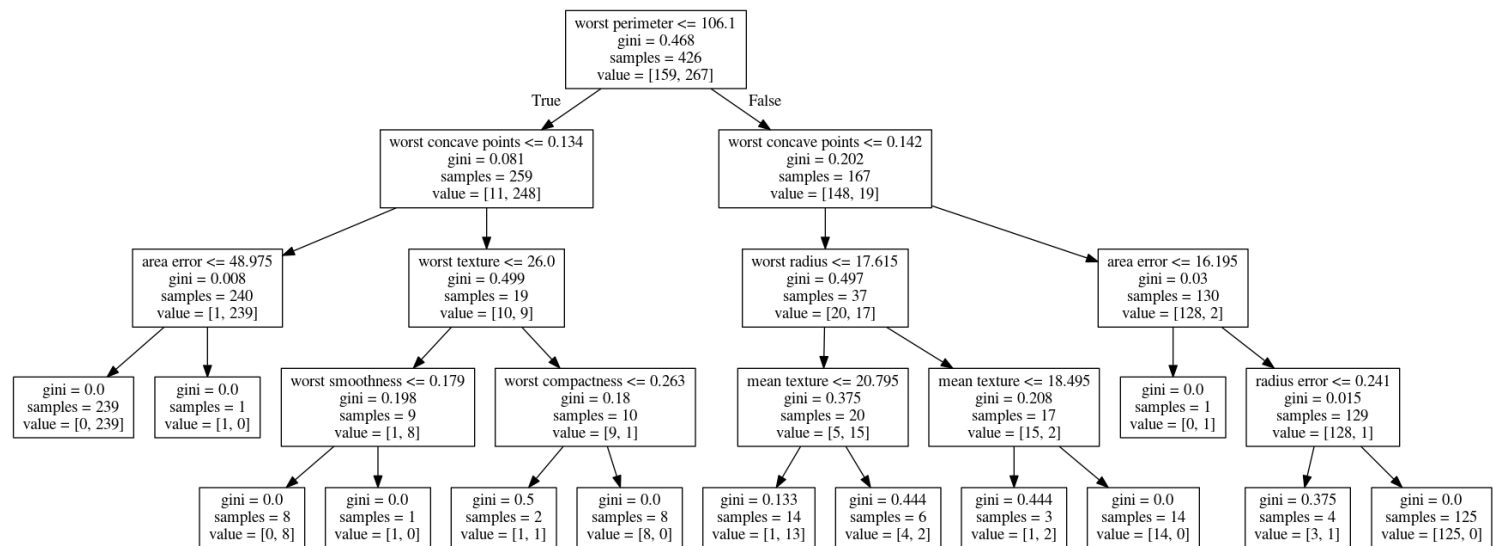
Parameter Tuning

- Pre-pruning and post-pruning (not in sklearn yet)
- Limit tree size (pick one, maybe two):
 - max_depth
 - max_leaf_nodes
 - min_samples_split
 - min_impurity_decrease

No pruning



max_depth = 4



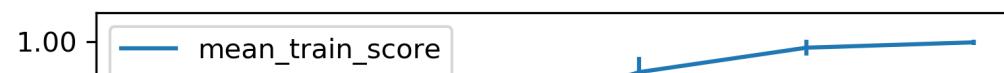
`max_leaf_nodes = 8`

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

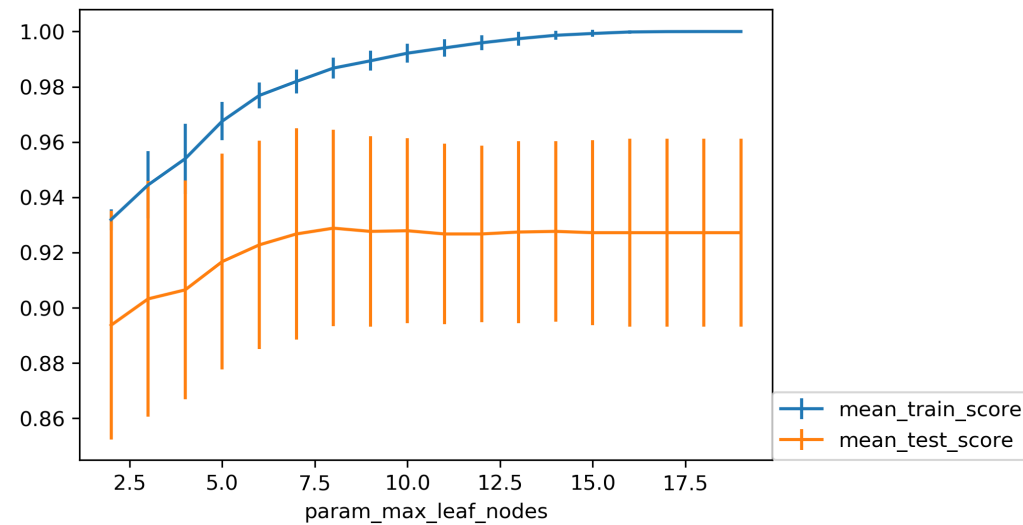
`min_samples_split = 50`

```
worst perimeter <= 106.1  
  gini = 0.468  
  samples = 426  
  value = 0.60 0.671
```

```
from sklearn.model_selection import GridSearchCV
param_grid = {'max_depth': range(1, 7)}
grid = GridSearchCV(DecisionTreeClassifier(random_state=0),
                    param_grid=param_grid, cv=10)
grid.fit(X_train, y_train)
```




```
from sklearn.model_selection import GridSearchCV
param_grid = {'max_leaf_nodes':range(2, 20)}
grid = GridSearchCV(DecisionTreeClassifier(random_state=0),
                    param_grid=param_grid, cv=10)
grid.fit(X_train, y_train)
```



Relation to Nearest Neighbors

- Predict average of neighbors – either by k , by epsilon ball or by leaf.
- Trees are much faster to predict.
- Both can't extrapolate

Extrapolation



Extrapolation

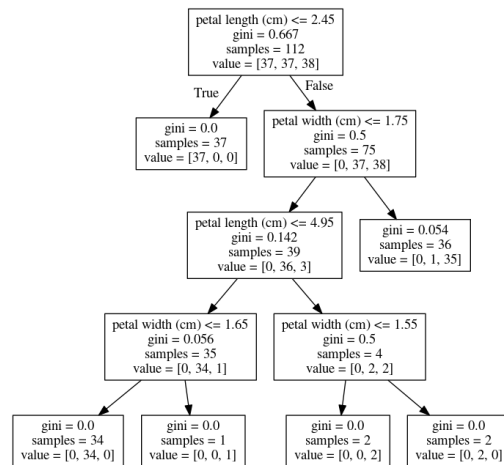


Extrapolation

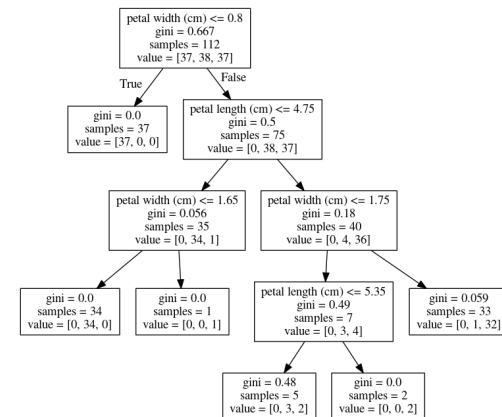


Instability

```
X_train, X_test, y_train, y_test = train_test_split(
    iris.data, iris.target, stratify=iris.target, random_state=0)
tree = DecisionTreeClassifier(max_leaf_nodes=6)
tree.fit(X_train, y_train)
```



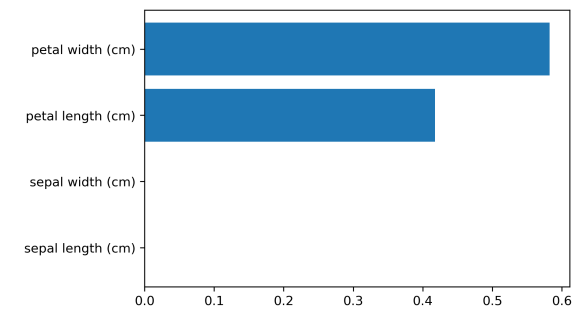
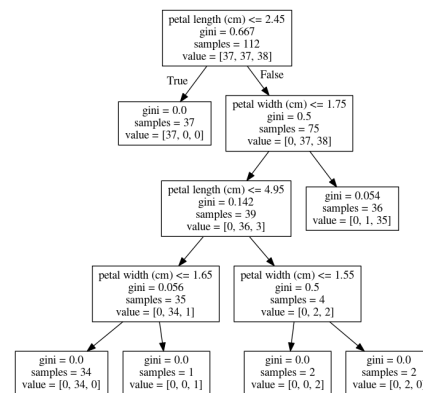
```
X_train, X_test, y_train, y_test = train_test_split(
    iris.data, iris.target, stratify=iris.target, random_state=1)
tree = DecisionTreeClassifier(max_leaf_nodes=6)
tree.fit(X_train, y_train)
```



Feature importance

```
X_train, X_test, y_train, y_test = train_test_split(
    iris.data, iris.target, stratify=iris.target, random_state=0)
tree = DecisionTreeClassifier(max_leaf_nodes=6)
tree.fit(X_train, y_train)
```

```
tree.feature_importances_
array([0.0, 0.0, 0.414, 0.586])
```



Categorical Data

- Can split on categorical data directly
- Intuitive way to split: split in two subsets
- 2^n many possibilities
- Possible to do in linear time exactly for gini index and binary classification.
- Heuristics done in practice for multi-class.
- Not in sklearn :(

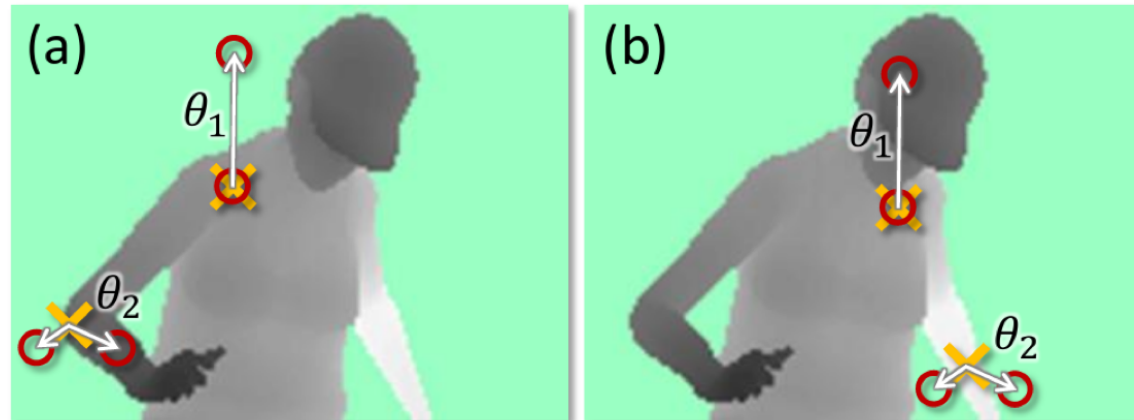
Predicting probabilities

- Fraction of class in leaf.
- Without pruning: Always 100% certain!
- Even with pruning might be too certain.

Conditional Inference Trees

- Select “best” split with correcting for multiple-hypothesis testing.
- More “fair” to categorical variables.
- Only in R so far (party)

Different splitting methods



(taken from Shotton et. al. Real-Time Human Pose Recognition ..)

Ensemble Models

Poor man's ensembles

- Build different models
- Average the result
- Owen Zhang (long time kaggle 1st): build XGBoosting models with different random seeds.
- 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 hard and soft voting

VotingClassifier

```
voting = VotingClassifier(
    [('logreg', LogisticRegression(C=100)),
     ('tree', DecisionTreeClassifier(max_depth=3, random_state=0))],
    voting='soft')
voting.fit(X_train, y_train)
lr, tree = voting.estimators_
voting.score(X_test, y_test), lr.score(X_test, y_test), tree.score(X_test, y_test)
```

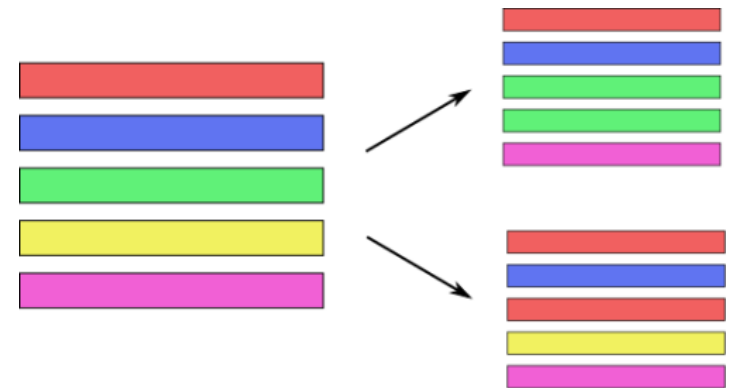
0.88 0.84 0.80



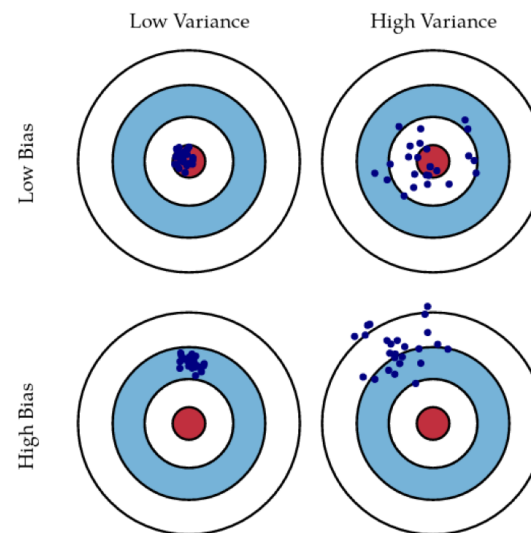
30 / 43

Bagging (Bootstrap AGGregation)

- Generic way to build “slightly different” models
- BaggingClassifier, BaggingRegressor



Bias and Variance

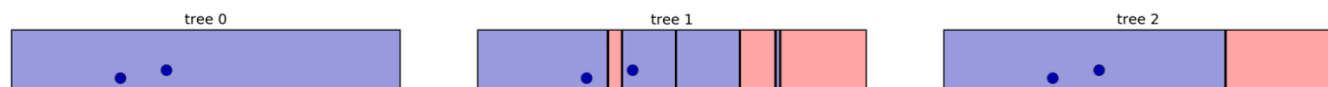


<http://scott.fortmann-roe.com/docs/BiasVariance.html>

Bias and Variance in Ensembles

- Breiman showed that generalization depends on strength of the individual classifiers and (inversely) on their correlation
- Uncorrelating them might help, even at the expense of strength

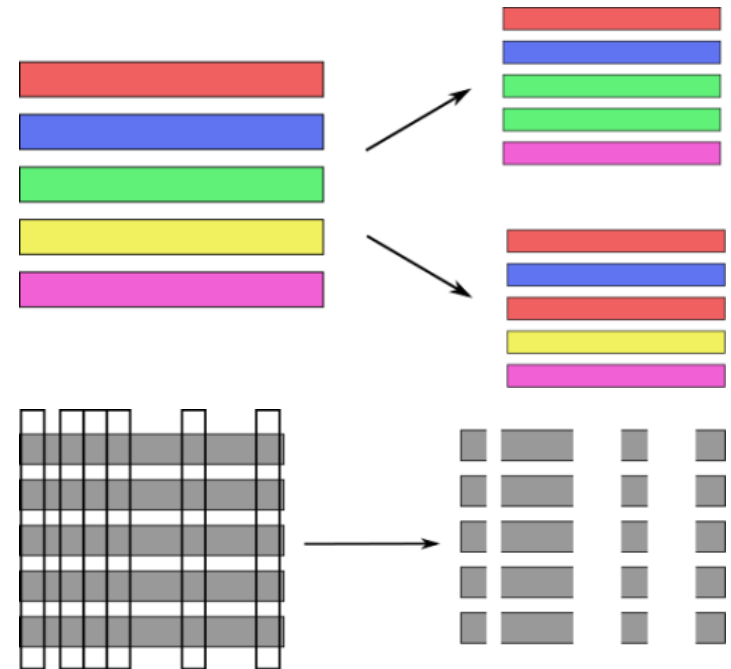
Random Forests



34 / 43

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

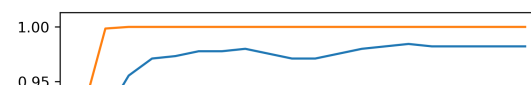
- Main parameter: `max_features`
 - around $\sqrt{n_{\text{features}}}$ for classification
 - Around `n_features` for regression
- `n_estimators > 100`
- Prepruning might help, definitely helps with model size!
- `max_depth`, `max_leaf_nodes`, `min_samples_split` again

Extremely Randomized Trees

- More randomness!
- Randomly draw threshold for each feature!
- Doesn't use bootstrap
- Faster because no sorting / searching
- Can have smoother boundaries

Warm-Starts

```
train_scores = []  
test_scores = []  
  
rf = RandomForestClassifier(warm_start=True)  
estimator_range = range(1, 100, 5)  
for n_estimators in estimator_range:  
    rf.n_estimators = n_estimators  
    rf.fit(X_train, y_train)  
    train_scores.append(rf.score(X_train, y_train))  
    test_scores.append(rf.score(X_test, y_test))
```



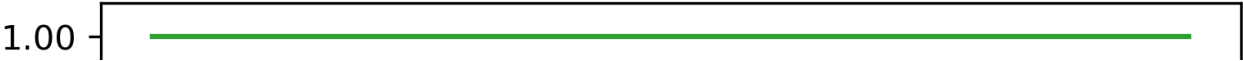
Out-of-bag estimates

- Each tree only uses ~66% of data
- Can evaluate it on the rest!
- Make predictions for out-of-bag, average, score.
- Each prediction is an average over different subset of trees

```
train_scores = []
test_scores = []
oob_scores = []

feature_range = range(1, 64, 5)
for max_features in feature_range:
    rf = RandomForestClassifier(max_features=max_features, oob_score=True,
                               n_estimators=200, random_state=0)
    rf.fit(X_train, y_train)
    train_scores.append(rf.score(X_train, y_train))
    test_scores.append(rf.score(X_test, y_test))
    oob_scores.append(rf.oob_score_)
```

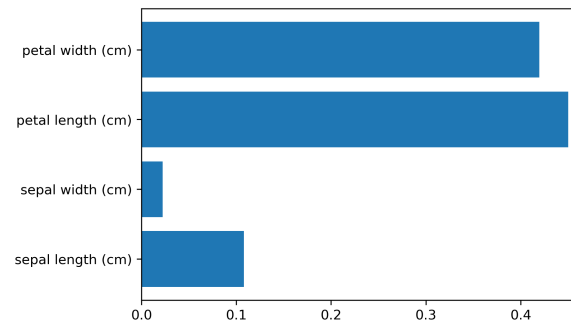
39 / 43



Variable Importance

```
X_train, X_test, y_train, y_test = train_test_split(
    iris.data, iris.target, stratify=iris.target, random_state=1)
rf = RandomForestClassifier(n_estimators=100).fit(X_train, y_train)
rf.feature_importances_
plt.barh(range(4), rf.feature_importances_)
plt.yticks(range(4), iris.feature_names);
```

```
array([ 0.126, 0.033, 0.445, 0.396])
```



Questions ?

