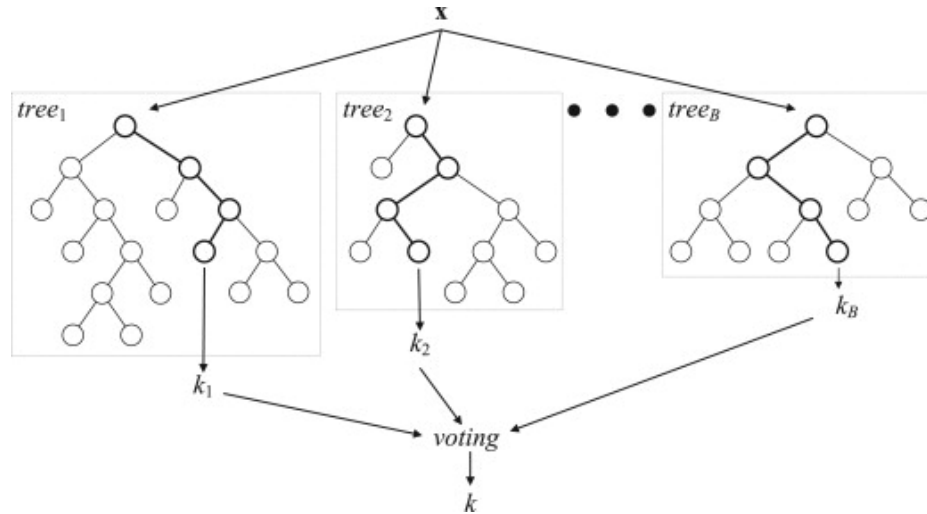


Introduction to Data Science

(Decision Trees and) Random Forests



Random Forests

- Train an "**ensemble**" of decisions trees
- RF decision is
 - the mean of decisions by individual trees (regression)
 - the majority vote of decisions by individual trees (classification)

Pros

- Reduces variance of individual decision tree models
- Can rank the importance of features in a natural way
- Can operate on many (1000s) of features without dimensionality reduction
- Can train well with relatively few samples

Cons

- Loss of interpretability, compared to individual decision trees
- Not as good at regression:
 - Output is not truly continuous valued, but rather, discretized (less problematic than for individual trees, however)

Random Forest trees are different from regular decision trees

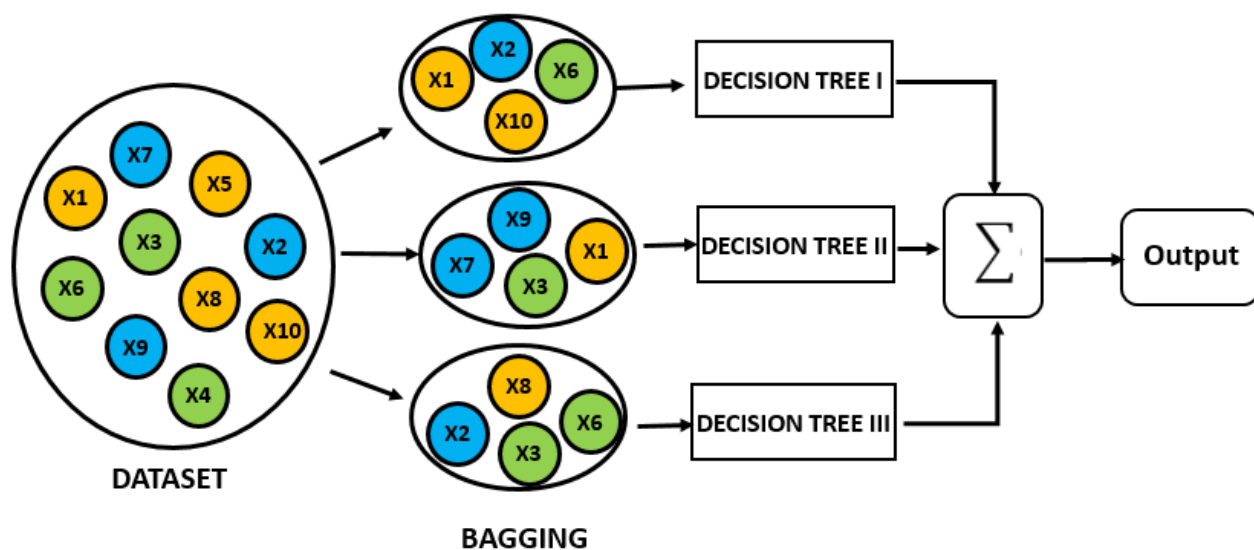
If the same training set and method is used for each tree, all the trees will be correlated (produce nearly identical outputs) and nothing will be gained by having an ensemble.

Use two forms of **randomization** when training each tree, to **reduce tree correlation**.

1. **Bagging**
2. **Feature randomization (feature bagging)**

Bagging (bootstrap aggregating)

The training set for each tree is B samples randomly drawn from the full training set, **with replacement**.



Bonus: We can get an "out-of-bag" (OOB) score from our training set:

For a given sample, we make a prediction using subset of trees that were not trained with that sample. This is a form of cross-validation.

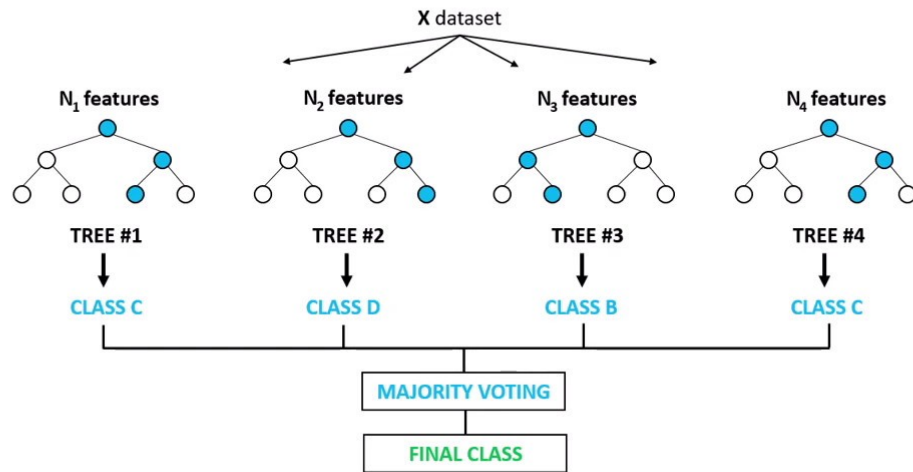
Feature subset randomization

During training of **individual RF trees**, a random subset of **features** is selected at each branch split in the tree. In regular decision trees, all features are considered.

The figure below doesn't quite do this justice. It shows a subset of features being used for each tree.

However, in RF training of a single tree, **a random subset of features is selected at each split, and from those features one is then selected as the feature on which to split the data** (based on Gini impurity or some other metric).

Random Forest Classifier



We'll work with the Titanic data set again

```
In [5]: 1 # Get necessary packages
2 import sklearn
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import graphviz # needed to visualize trained decision tree
7 from sklearn.tree import export_graphviz # needed to visualize trained decision tree
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.model_selection import train_test_split
10 from sklearn import tree
11
12 np.random.seed(1000)
13
14 # Set up for plotting
15 plt.style.use("ggplot")
16 %matplotlib inline
17 %config InlineBackend.figure_format = 'retina'
18
19 # Read in our data and glance at its formatting
20 df = pd.read_csv('titanic_data.csv')
21 df.head()
```

Out[5]:

	Survived	Pclass	Name	Sex	Age	Siblings/Spouses Aboard	Parents/Children Aboard	Fare
0	0	3	Mr. Owen Harris Braund	male	22.0	1	0	7.2500
1	1	1	Mrs. John Bradley (Florence Briggs Thayer) Cum...	female	38.0	1	0	71.2833
2	1	3	Miss. Laina Heikkinen	female	26.0	0	0	7.9250
3	1	1	Mrs. Jacques Heath (Lily May Peel) Futrelle	female	35.0	1	0	53.1000
4	0	3	Mr. William Henry Allen	male	35.0	0	0	8.0500

```
In [6]: 1 ## Prepare our data for viable input to the sklearn decision tree model
2
3 # Convert Sex feature into a pair of Boolean/binary dummy features (male and female)
4 df_dummy = pd.get_dummies(df, columns=['Sex'])
5
6 # Remove redundant male or female feature, and unneeded Name feature
7 df_dummy = df_dummy.drop(columns=['Sex_male', 'Name'])
8 df_dummy.head()
9
10 # Create train/test data sets
11 X = df_dummy.drop(columns=['Survived'])
12 y = df_dummy['Survived']
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

Let's repeat our work on decision trees, for eventual comparison with random forest models

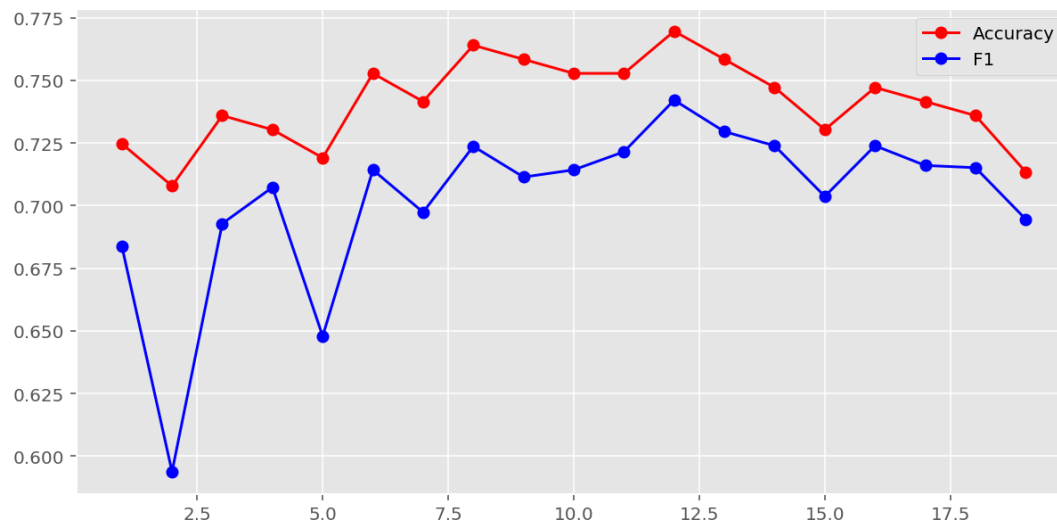
```

In [7]: 1  ## Let's train models over a range of depths, and score them with the test set
2  from sklearn.metrics import f1_score
3
4  acc_scores_dt = []
5  f1_scores_dt = []
6  max_max = 20
7  max_depth = np.arange(1, max_max)
8
9  for depth in max_depth:
10     # Build model and train
11     titanic_tree = DecisionTreeClassifier(max_depth=depth, random_state=0)
12     titanic_tree.fit(X_train, y_train)
13
14     # Test
15     acc_scores_dt.append(titanic_tree.score(X_test, y_test))
16     y_test_hat = titanic_tree.predict(X_test)
17     f1_scores_dt.append(f1_score(y_test, y_test_hat))
18
19 # Plot results
20 plt.figure(figsize=(10, 5))
21 plt.plot(max_depth, acc_scores_dt, 'ro-', label='Accuracy')
22 plt.plot(max_depth, f1_scores_dt, 'bo-', label='F1')
23 plt.legend()
24
25 # Print values for best test score
26 ix_best = np.argmax(acc_scores_dt)
27 print('Best accuracy score is %0.3f, for max_depth=%d' % (acc_scores_dt[ix_best], max_depth[ix_best]))
28 ix_best = np.argmax(f1_scores_dt)
29 print('Best F1-score is %0.3f, for max_depth=%d' % (f1_scores_dt[ix_best], max_depth[ix_best]))

```

Best accuracy score is 0.770, for max_depth=12

Best F1-score is 0.742, for max_depth=12



Now we'll use sklearn to train random forests on the same data set

In addition to scores for the test set, we'll measure the *out-of-bag* error from the training set

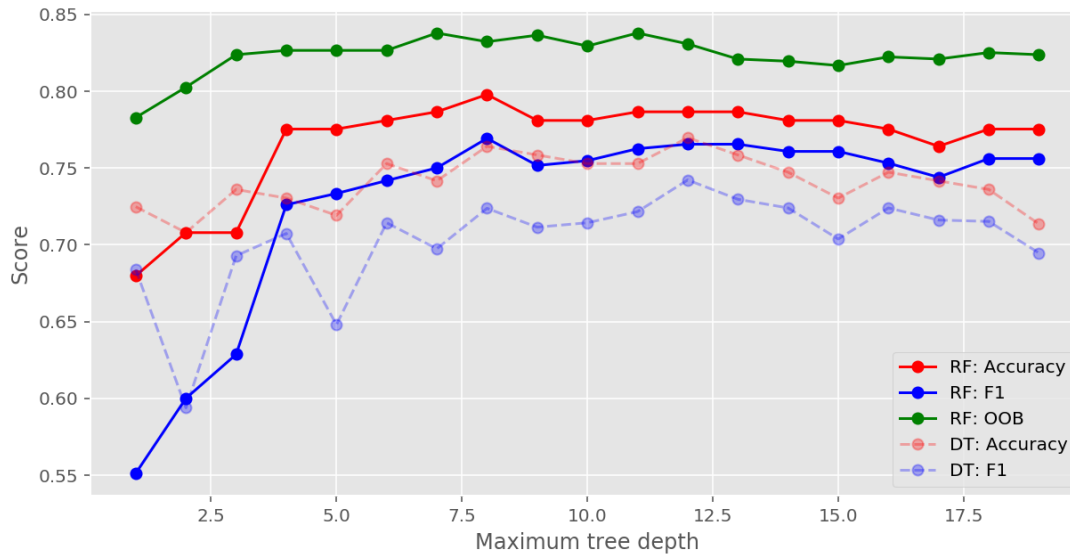
```
In [8]: 1  ## We'll train models over a range of tree depths
2
3  from sklearn.ensemble import RandomForestClassifier
4
5  n_estimators = 100
6  max_depth = np.arange(1, max_max)
7  acc_scores_rf = []
8  fl_scores_rf = []
9  oob_scores_rf = []
10 for depth in max_depth:
11     # Build model and train
12     titanic_forest = RandomForestClassifier(n_estimators=n_estimators,
13                                           max_depth=depth,
14                                           oob_score=True,
15                                           random_state=0) # call the RandomForest
16     titanic_forest.fit(X_train, y_train)
17
18     # Save out-of-bag (OOB) score
19     oob_scores_rf.append(titanic_forest.oob_score_)
20
21     # Test
22     acc_scores_rf.append(titanic_forest.score(X_test, y_test))
23     y_test_hat = titanic_forest.predict(X_test)
24     fl_scores_rf.append(f1_score(y_test, y_test_hat))
25
26 # Print values for best test score
27 ix_best = np.argmax(acc_scores_rf)
28 print('Best accuracy score is %0.3f, for max_depth=%d' % (acc_scores_rf[ix_best], max_depth))
29 ix_best = np.argmax(fl_scores_rf)
30 print('Best F1-score is %0.3f, for max_depth=%d' % (fl_scores_rf[ix_best], max_depth))
31 ix_best = np.argmax(oob_scores_rf)
32 print('Best OOB score is %0.3f, for max_depth=%d' % (oob_scores_rf[ix_best], max_depth))
```

Best accuracy score is 0.798, for max_depth=8

Best F1-score is 0.769, for max_depth=8

Best OOB score is 0.838, for max_depth=7

```
In [9]: 1 # Plot results
2 plt.figure(figsize=(10, 5))
3 plt.plot(max_depth, acc_scores_rf, 'ro-', label='RF: Accuracy')
4 plt.plot(max_depth, f1_scores_rf, 'bo-', label='RF: F1')
5 plt.plot(max_depth, oob_scores_rf, 'go-', label='RF: OOB')
6 plt.plot(max_depth, acc_scores_dt, 'ro--', alpha=0.3, label='DT: Accuracy')
7 plt.plot(max_depth, f1_scores_dt, 'bo--', alpha=0.3, label='DT: F1')
8 plt.xlabel('Maximum tree depth')
9 plt.ylabel('Score')
10 _ = plt.legend()
```



Let's see which features have the highest "importance" for an RF with the best max_depth.

Importance, in sklearn, is calculated as the **Mean Decrease in Impurity (MDI)**.

That is, the **average decrease in Gini impurity** across nodes that use the feature for splitting, weighted by the number of samples that reach that node.

```

In [10]: 1 ix_best = np.argmax(acc_scores_rf)
2 max_depth_best = max_depth[ix_best]
3
4 # Build model and train
5 titanic_forest = RandomForestClassifier(n_estimators=n_estimators, random_state=0,
6 titanic_forest.fit(X_train, y_train)
7
8 plt.figure(figsize=(10,5))
9 column_names = X_train.columns
10 x_tick = np.array(range(len(column_names)))
11 plt.bar(x_tick, titanic_forest.feature_importances_)
12 plt.xticks(x_tick, column_names, rotation=90, fontsize=18)
13 _ = plt.ylabel('Importance', fontsize=18)

```

