# Machine Learning

Lecture 7a

# Let's apply the Adaboost algorithm to our data set

```python
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  %matplotlib inline
5  from sklearn.metrics import accuracy_score, confusion_matrix
6  from sklearn.tree import DecisionTreeClassifier
7  from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
8  pd.set_option('use_inf_as_na', True)
9  from collections import Counter
```

```python
1  raw_data = pd.read_pickle(r'C:\Users\niels\data\dataset.pkl')
2  data = raw_data[raw_data['market_cap'] > 1000.0]
3  data.fillna(0.0,inplace=True)
```

```
C:\Users\niels\Anaconda3\lib\site-packages\pandas\core\frame.py:3790: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy
  downcast=downcast, **kwargs)
```

```python
1  def f(x):
2      if x > 0.01:
3          return 1
4      elif x < -0.01:
5          return -1
6      else:
7
8          return 0
```

An AdaBoost classifier.

An AdaBoost [1] classifier is a meta-estimator that begins by fitting a classifier on the original dataset and then fits additional copies of the classifier on the same dataset but where the weights of incorrectly classified instances are adjusted such that subsequent classifiers focus more on difficult cases.

This class implements the algorithm known as AdaBoost-SAMME [2].

As a base estimator we choose a DecisionTreeClassifier with max_depth = 4

```
1  ada_clf = AdaBoostClassifier(DecisionTreeClassifier(max_depth=4),n_estimators=25)
```

# It takes a bit of time to train the classifier on the full data set

```
1  ada_clf.fit(train_1,y_1)
```

```
AdaBoostClassifier(algorithm='SAMME.R',
          base_estimator=DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=4,
            max_features=None, max_leaf_nodes=None,
            min_impurity_decrease=0.0, min_impurity_split=None,
            min_samples_leaf=1, min_samples_split=2,
            min_weight_fraction_leaf=0.0, presort=False, random_state=None,
            splitter='best'),
          learning_rate=1.0, n_estimators=25, random_state=None)
```

```
1  %timeit ada_clf.fit(train_1,y_1)
```

```
30.8 s ± 1.96 s per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

# Feature Importances

```python
1  def adaboost_feat_importances(m, df):
2
3      return pd.DataFrame({'cols':df.columns, 'feat_imp': m.feature_importances_}
4                          ).sort_values('feat_imp', ascending=False)
5
6  def plot_fi(fi): return fi.plot('cols', 'feat_imp', 'barh', figsize=(12,7), legend=False)
```
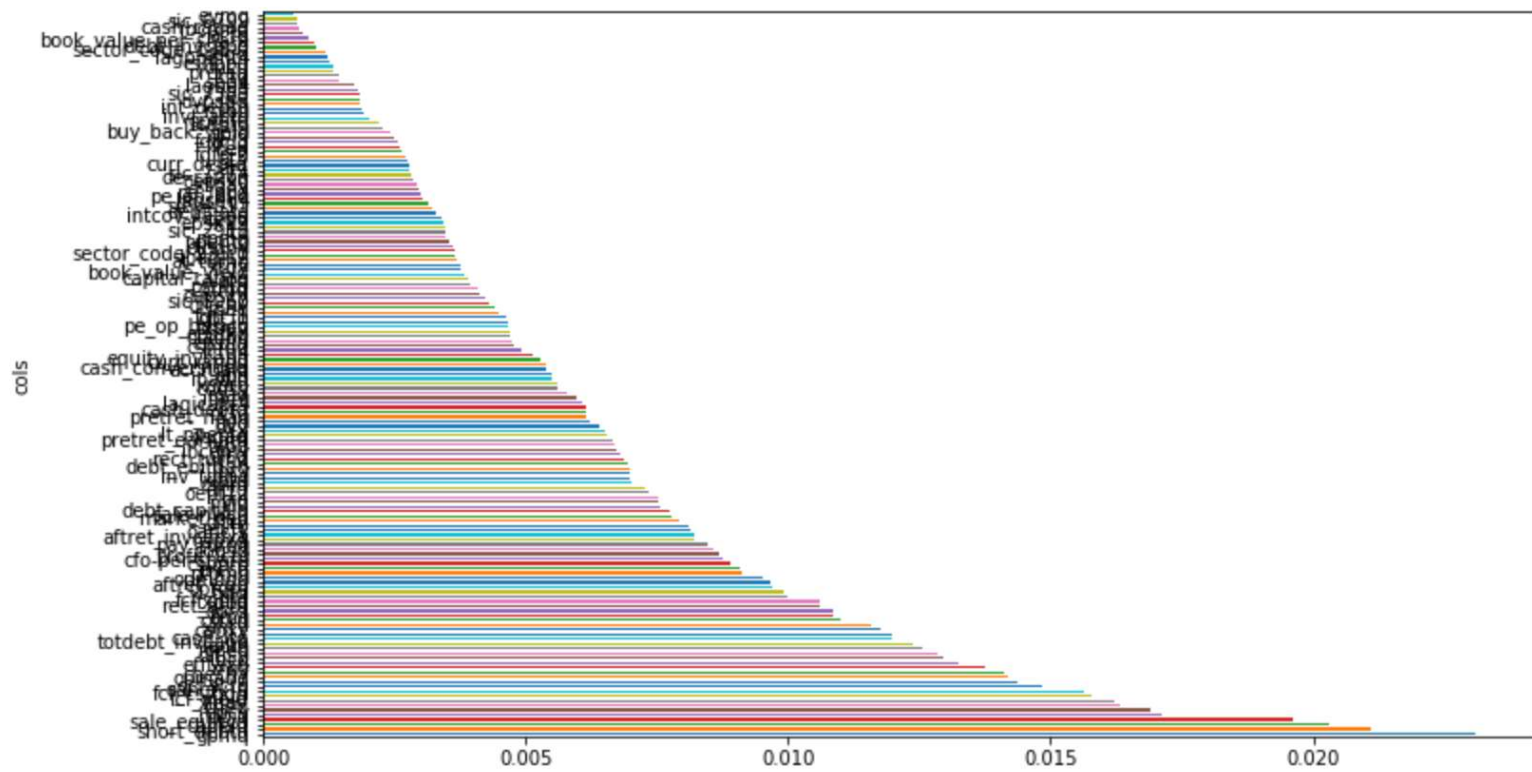
```python
1  fi = adaboost_feat_importances(ada_clf,train_1)
```

```python
1  features = fi[(fi['feat_imp'] > 0.0)]
2  features.shape
```

)]: (122, 2)

```python
1  plot_fi(features);
```

```
1  plot_fi(features);
```

```
1  features
```

| | cols | feat_imp |
|---|---|---|
| 112 | gpmq | 0.024416 |
| 148 | short_debtq | 0.021091 |
| 99 | capeiq | 0.020853 |
| 174 | sale_equityq | 0.019603 |
| 36 | nopiq | 0.017122 |
| 71 | dpcy | 0.016910 |
| 4 | cheq | 0.016338 |
| 155 | fcf_csfhdq | 0.015632 |
| 156 | fcf_yield | 0.015610 |
| 102 | oancfy_q | 0.014830 |
| 152 | ocf_lctq | 0.014370 |
| 111 | opmadq | 0.014169 |
| 66 | chechy | 0.014123 |
| 82 | xidoy | 0.014021 |
| 49 | rectq | 0.014010 |
| 75 | miiy | 0.013920 |
| 120 | offtaxq | 0.013235 |

| | cols | feat_imp |
|---|---|---|
| 138 | debt_invcapq | 0.000983 |
| 90 | book_value_per_share | 0.000875 |
| 0 | actq | 0.000767 |
| 27 | ibcomq | 0.000682 |
| 165 | cash_ratioq | 0.000650 |
| 540 | sic_6799 | 0.000633 |
| 95 | evmq | 0.000572 |

153 rows × 2 columns

# We cut down to the features that actually have an influence on the classifier

```
1  train_1 = train_1[features['cols'].values]
2  valid = valid[features['cols'].values]
```

```
1  %timeit ada_clf.fit(train_1,y_1)
2  ada_clf.score(train_1,y_1)
```

```
24 s ± 714 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

: 0.5565183580514359

```
1  pred_valid = ada_clf.predict(valid)
2  ada_clf.score(valid,y_valid)
```

: 0.44457142857142856

```
1  (pred_valid * valid_stock_returns).sum()
```

: 27.26727199999997

```
1  Counter(pred_valid)
```

: Counter({1: 1110, 0: 56, -1: 584})

# Profit Importance

```python
def profit_importance(m,df,rets):
#     np.random.seed(123)
    profit = []
    for col in df.columns:
        prof = []
        for _ in range(10):
            X = df.copy()
            X[col] = np.random.permutation(X[col].values)
            prediction = m.predict(X)
            prof.append((prediction * rets).sum())
        profit.append(np.mean(prof))
    return profit
```

```python
def adaboost_profit_importance(m, df,rets):
    return pd.DataFrame({'cols':df.columns, 'pi_imp':profit_importance(m,df,rets)}
                        ).sort_values('pi_imp', ascending=True)
```

```python
pi = adaboost_profit_importance(ada_clf,valid,valid_stock_returns)
pi
```

|    | cols | pi_imp |
|----|------|--------|
| 50 | piq | 17.245805 |
| 6 | cheq | 17.895002 |
| 80 | equity_invcapq | 22.199241 |
| 74 | xoprq | 22.533351 |
| 96 | atq | 23.067195 |

# Finding the features that optimize the strategy profit on the validation set

```python
1   profits = []
2   feat=[]
3
4   train = train_1.copy()
5   validation = valid.copy()
6
7   while len(train.columns)>1:
8
9       pred_valid = ada_clf.predict(validation)
10
11      print((pred_valid * df_valid['next_period_return']).sum())
12      profits.append((pred_valid * df_valid['next_period_return']).sum())
13      feat.append(train.columns)
14
15      col_to_drop = pi.iloc[-1]['cols']
16      train.drop(col_to_drop,axis=1,inplace=True)
17      validation.drop(col_to_drop,axis=1,inplace=True)
18
19      ada_clf.fit(train,y_1)
20      pi = adaboost_profit_importance(ada_clf,validation,df_valid['next_period_return'])
21
22
23
```

27.26727199999997

One of the interesting things about AdaBoosting is that we do not have to use just trees as base estimators. Here is an example boosting a RandomForest
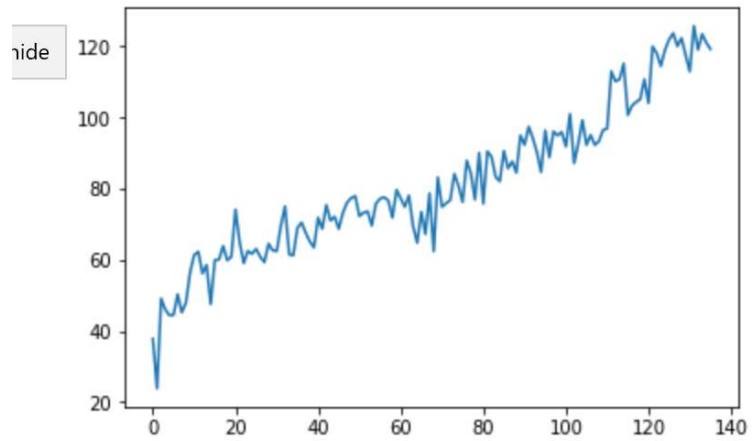
```
1  rf_clf = RandomForestClassifier(n_estimators=15,min_samples_leaf=2400)
```

```
1  ada_clf = AdaBoostClassifier(rf_clf,n_estimators=25)
```

```
1  %timeit ada_clf.fit(train_1,y_1)
```

7.87 s ± 283 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

```
1  plt.plot(profits);
```



```
1  optim_feats = feat[n]
2  optim_feats
3
```

Index(['cstkq', 'ibadj12', 'fcf_ocfq', 'cshfdq', 'cash_conversionq', 'cshopq'], dtype='object')

```
1  train_1_optim = train_1[optim_feats]
2  valid_optim = valid[optim_feats]
3
4
5  ada_clf.fit(train_1_optim,y_1)
6  print(ada_clf.score(train_1_optim,y_1))
7  pred_valid_tree = ada_clf.predict(valid_optim)
8  print(ada_clf.score(valid_optim,y_valid))
9  (pred_valid_tree * valid_stock_returns).sum()
```

```
0.5075820696550575
0.532
```

]: 119.12096800000003

```
1  train_2_tree = train_2[optim_feats]
2  test_tree = test[optim_feats]
3  ada_clf.fit(train_2_tree,y_2)
4  pred_test_tree = ada_clf.predict(test_tree)
5  (pred_test_tree * test_stock_returns).sum()
```

]: 61.964654999999965

```
1  train_1_optim = train_1[optim_feats]
2  valid_optim = valid[optim_feats]
3
4
5  ada_clf.fit(train_1_optim,y_1)
6  print(ada_clf.score(train_1_optim,y_1))
7  pred_valid_tree = ada_clf.predict(valid_optim)
8  print(ada_clf.score(valid_optim,y_valid))
9  (pred_valid_tree * valid_stock_returns).sum()
```

```
0.5075820696550575
0.532
```

]: 119.12096800000003

```
1  train_2_tree = train_2[optim_feats]
2  test_tree = test[optim_feats]
3  ada_clf.fit(train_2_tree,y_2)
4  pred_test_tree = ada_clf.predict(test_tree)
5  (pred_test_tree * test_stock_returns).sum()
```

]: 61.964654999999965

```python
1  train_1_optim = train_1[optim_feats]
2  valid_optim = valid[optim_feats]
3
4
5  ada_clf.fit(train_1_optim,y_1)
6  print(ada_clf.score(train_1_optim,y_1))
7  pred_valid_tree = ada_clf.predict(valid_optim)
8  print(ada_clf.score(valid_optim,y_valid))
9  (pred_valid_tree * valid_stock_returns).sum()
```
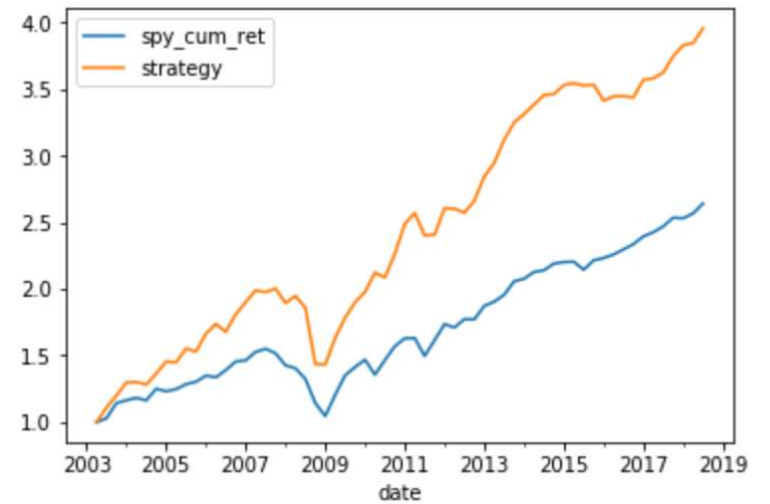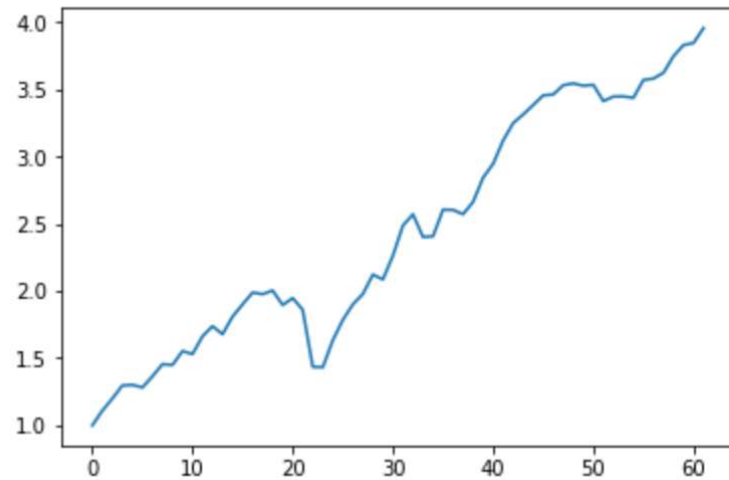
```
0.5075820696550575
0.532
```
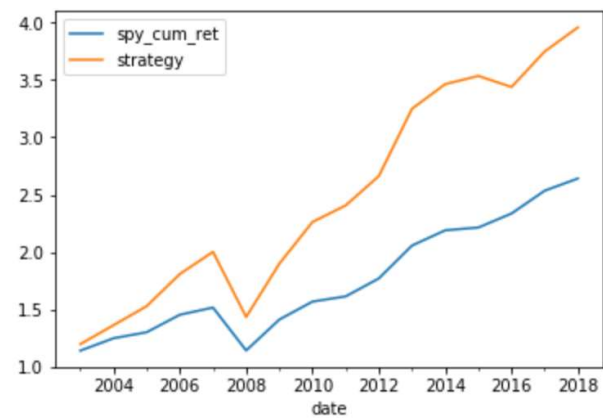
]: 119.12096800000003

```python
1  train_2_tree = train_2[optim_feats]
2  test_tree = test[optim_feats]
3  ada_clf.fit(train_2_tree,y_2)
4  pred_test_tree = ada_clf.predict(test_tree)
5  (pred_test_tree * test_stock_returns).sum()
```
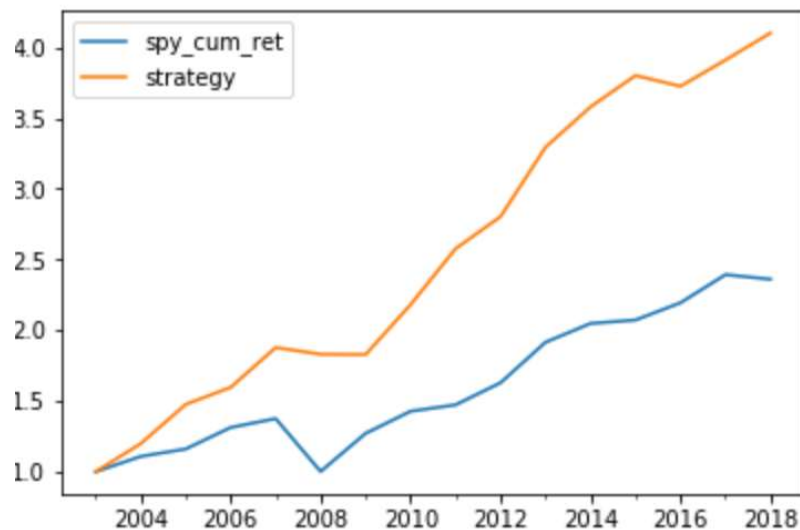
]: 61.964654999999965

```
1  plt.plot(x);
```





```
1  SPY = SPY.resample('Y').ffill()
2  SPY.plot();
```

# Variable optimal features



```
1  strategy_mean_ret = (SPY['strategy'] - 1).diff().mean()
2  strategy_std = (SPY['strategy'] - 1).diff().std()
3  strategy_sr = strategy_mean_ret/strategy_std
4  print('Strategy Sharpe Ratio: ',strategy_sr)
```

Strategy Sharpe Ratio:  1.3047801841589615

```
1  strategy_ret = (SPY['strategy'] - 1).diff().values[1:]
2  spy_ret = (SPY['spy_cum_ret'] - 1).diff().values[1:]
3
4  beta = (np.cov(spy_ret,strategy_ret)/np.var(spy_ret))[1,0]
5  beta
```

|: 0.403773768479095

```
1  residual_ret = strategy_ret - beta * spy_ret
2  IR = np.mean(residual_ret)/np.std(residual_ret)
3  IR
```

|: 1.200388421384077

```
1  alpha = np.mean(residual_ret)
2  alpha
```

|: 0.16120032434776224