# Experiences With XGBoost And The Financial Markets

## (Or How To Avoid An Embarrassing Call From The Bugatti Garage)

**PyData Bristol**
@PyDataBristol
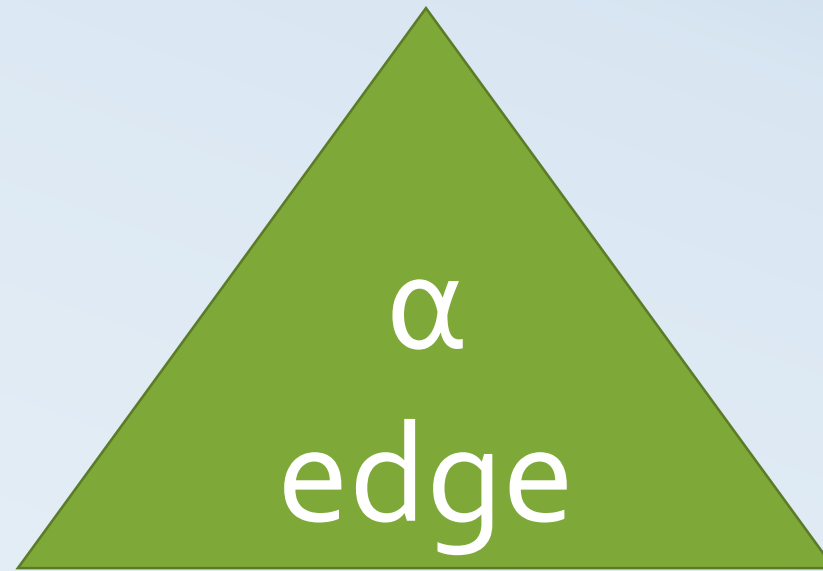
PyData meetups in London for data-loving pythonistas. Powered by @john_sandall, Frank Kelly, Miquel Perello Nieto.

⊙ Bristol, England

🔗 meetup.com/PyData-Bristol/

Sponsored by...

Strategy Paradigms
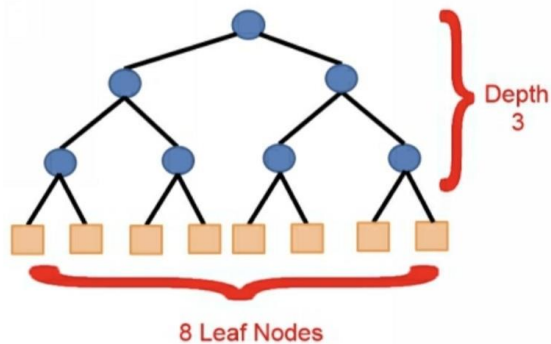
Fast:  < 500 ns?

α
edge

Smart

Dirty /
Naughty /
Black edge

# XGBoost & Decision Trees

- Non-linear

- Feature range insensitive
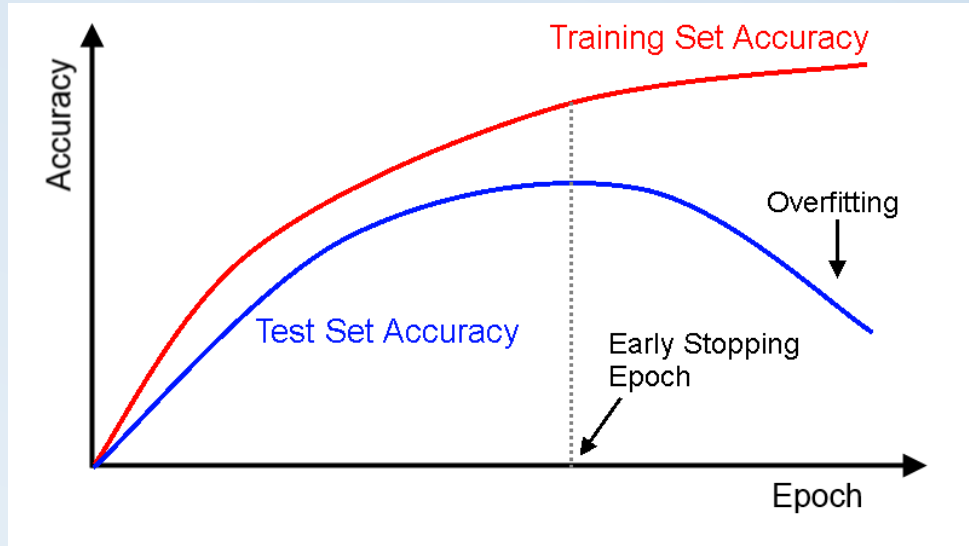
- Interpretable

- Feature selection & ranking



"Machine learning with boosting, a beginner's guide", Scott Hartshorn

```python
from xgboost import XGBRegressor
import operator


def SelectXGBoostBest(self, features, target):

    model = XGBRegressor(
        n_estimators=100,
        learning_rate=0.15,
        objective='reg:linear')
    model.fit(features, target)
    scores = model.booster().get_score(importance_type='weight')

    sorted_features = sorted(
        scores.items(),
        key=operator.itemgetter(1),
        reverse=True)

    top_n = 10
    top_features = []
    for f in sorted_features[:top_n]:
        top_features += [f[0]]

    return top_features
```

# Gradient Boosted Decision Tree Algorithm

- Iterative algorithm, incremental addition of weak learners focused on areas of poor prediction

- XGBoost has both native and scikit learn style python APIs, very quick to get going

- Multi-platform, for instance linux for simulation and windows for live trading



"Machine learning with boosting, a beginner's guide",  Scott Hartshorn

# Robustness & Overfitting
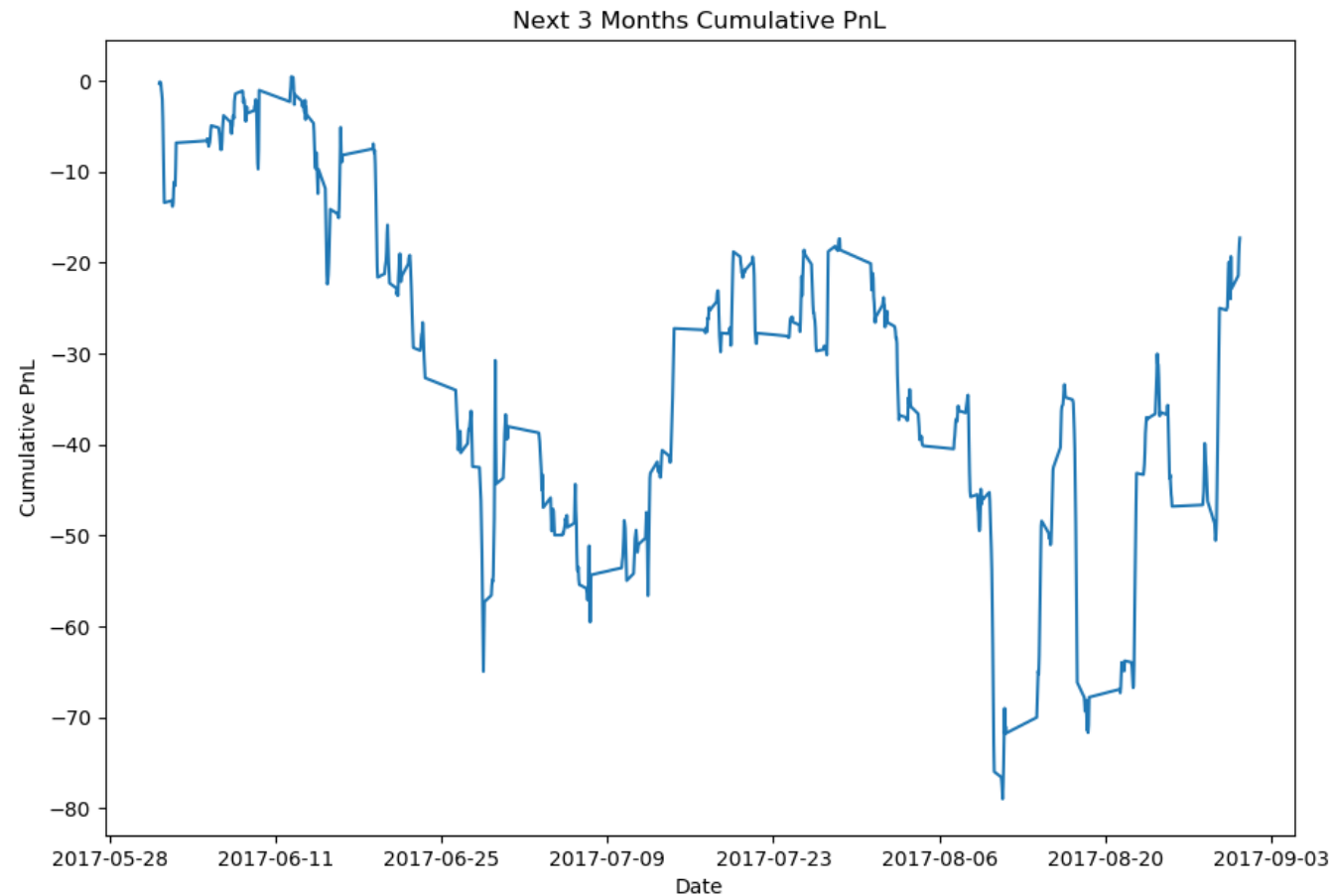


https://deeplearning4j.org/earlystopping

- Loss of robustness when training set improves but test does not

- XGBoost has a whole host of robustness settings e.g. early stopping

# Early Stopping & Evaluation Data Selection

```python
import numpy as np
from sklearn.model_selection import train_test_split
import xgboost as xgb


def _train_xgb(hyper, estimators, rounds, X, Y, test_ratio):

    xtrain, xtest, ytrain, ytest = train_test_split(
        X, Y, test_size=test_ratio)

    xgtrain = xgb.DMatrix(
        xtrain,
        label=ytrain,
        weight=np.ones(
            len(xtrain)))

    xgtest = xgb.DMatrix(xtest, label=ytest, weight=np.ones(len(xtest)))

    watchlist = [(xgtrain, 'train'), (xgtest, 'eval')]
    results = {}

    regressor = xgb.train(
        hyper,
        xgtrain,
        estimators,
        watchlist,
        early_stopping_rounds=rounds,
        verbose_eval=True,
        evals_result=results)

    return (regressor, results)
```

# The Next 3 Months…

# Beware Random Cross Validation!

```python
1  import numpy as np
2  #from sklearn.model_selection import train_test_split
3  import xgboost as xgb
4
5
6  def _train_xgb(hyper, estimators, rounds, X, Y, test_ratio):
7
8  #    xtrain, xtest, ytrain, ytest = train_test_split(
9  #        X, Y, test_size=test_ratio)
10
11     split = int(len(X) * (1 - test_ratio))
12     xtrain = X[:split]
13     xtest = X[split + 1:]
14     ytrain = Y[:split]
15     ytest = Y[split + 1:]
16
17     xgtrain = xgb.DMatrix(
18         xtrain,
19         label=ytrain,
20         weight=np.ones(
21             len(xtrain)))
22
23     xgtest = xgb.DMatrix(xtest, label=ytest, weight=np.ones(len(xtest)))
24
25     watchlist = [(xgtrain, 'train'), (xgtest, 'eval')]
26     results = {}
27
28     regressor = xgb.train(
29         hyper,
30         xgtrain,
31         estimators,
32         watchlist,
33         early_stopping_rounds=rounds,
34         verbose_eval=True,
35         evals_result=results)
36
37     return (regressor, results)
38
```