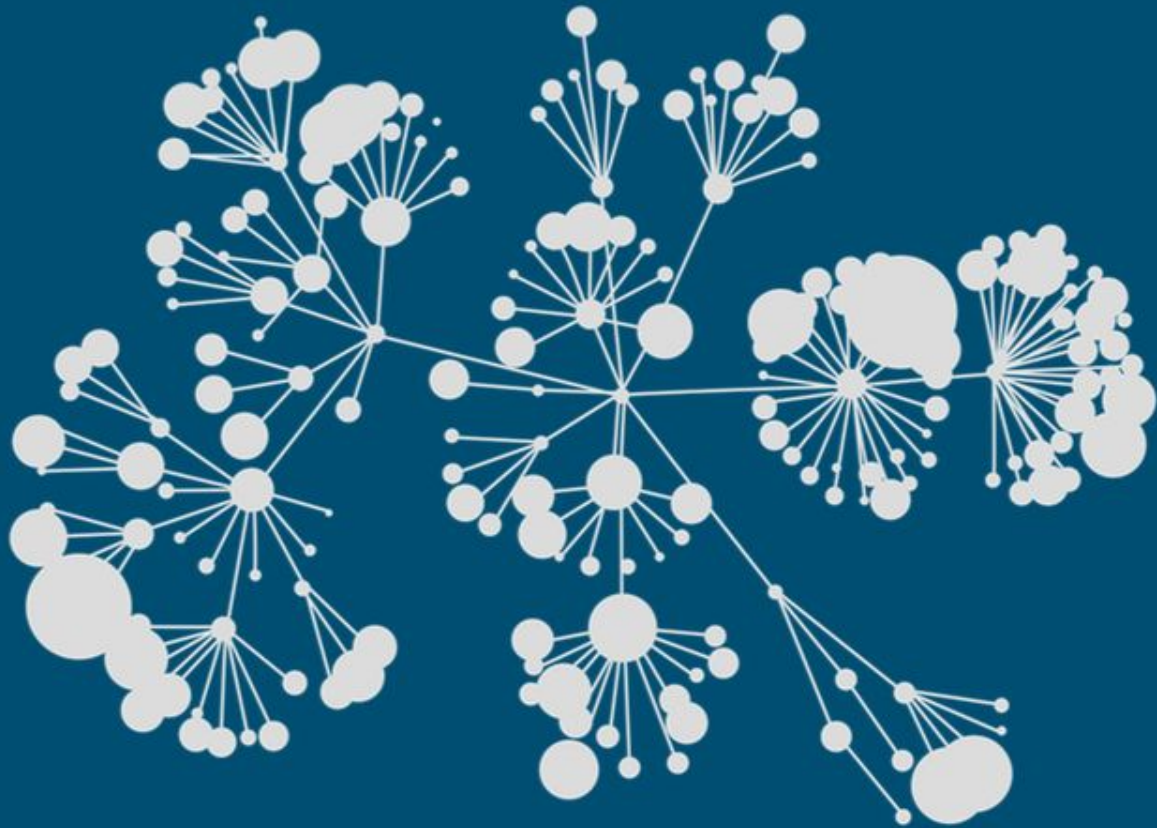# Kaggle

**Two Sigma: Using News to Predict Stock Movements**

**Andrzej Urbanowicz**
**Krzysztof Urbanowicz**

kaggle

Agenda

1. Background

2. Summary

3. Feature selection & Engineering

4. Training methods

5. Important findings

6. Simple model

# Andrzej Urbanowicz

- MSc in Computer Science
- Bioinformatics [RNA Sequencing analysis at the University of California Los Angeles (UCLA)]
- Big Data Architect/Data Engineer/Software Engineer working in Poland (Warsaw), Switzerland, and Los Angeles (USA)
- 15+ years experience in programming [working in big international IT companies such as IBM, Hewlett Packard (HP) and CERN] and startups

# Krzysztof Urbanowicz

- PhD in Physics
- Post-Doc in Max Planck Institute in Dresden
- Many years working as a Quant specialist, especially experienced in derivatives (futures, options)
- Experience in nonlinear time series analysis and econophysics
- 15+ years experience in finance and algo trading

- Model type:
  - ANN
- Most important features:
  - Volume
  - Close price
  - Open price
  - Returns from 1 and 10 days as raw data and residual to market
- The tools that were used include:
  - TensorFlow, Keras, Pandas
- Time needed to train our model:
  - up to 20 minutes

We generated several different approaches to cope with the problem and we found out what was the most important in this quest:

- It is very probable to overfit financial model, so we would like to make some non ML models or very simple ML models.

- We left only given features and did not generate new ones by any transformation.

- We generally stick to the simple rule of thumb that the natural logarithm from number of data should be much less than number of degrees of freedom.
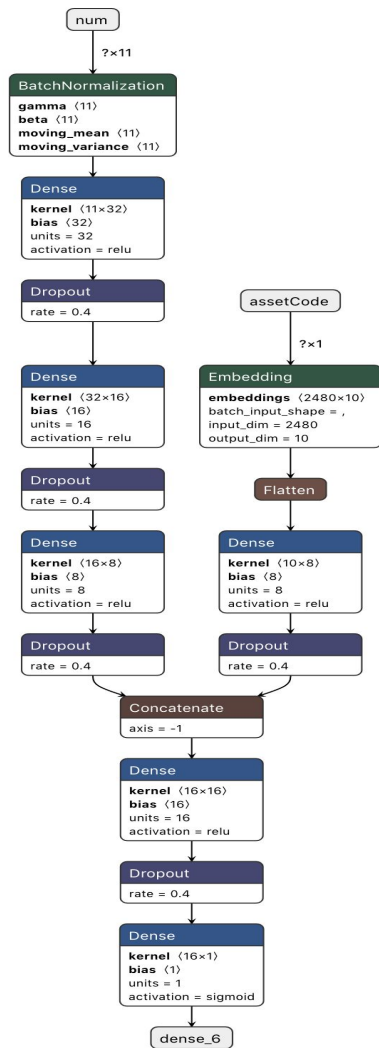
  **ln(#data) >> #features**

  Degrees of freedom can be understood as number of data features as well as of parameters of ANN. That is why, we did not artificially augment number of features by any transformation and we kept the model simple.
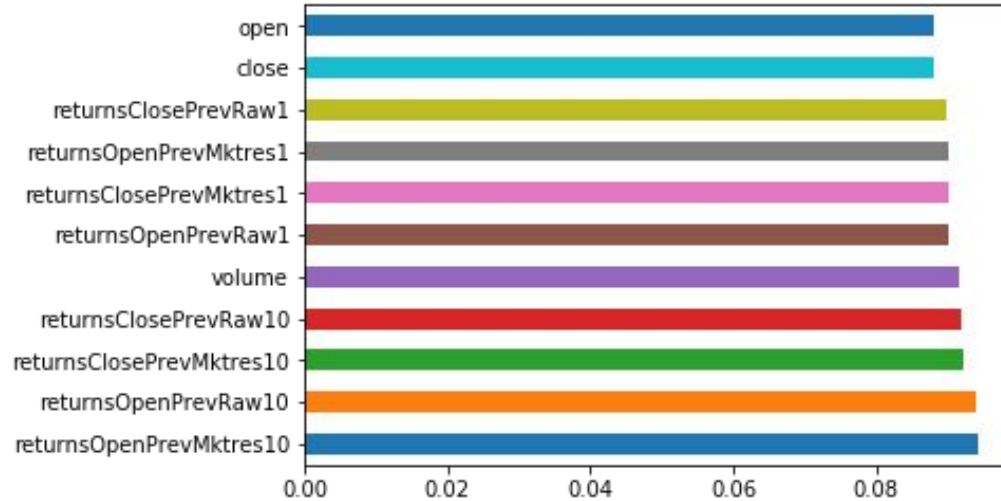
- Having in mind, that ML always tends to heavy fit to the train set, not to an unknown future as in this case, we allowed for this simple ANN a large dropout (0.4).
- For ANN construction, we set relatively small number of neurons with max 32 neurons in hidden layers.
- We did some simple preliminary cleaning of data (replacing with simple mean).
- We did not introduce L1/L2 regularization in ANN kernels.

- We tried to introduce XGBoost to news features training, but we found out that it may worsen the results as a whole.

- We finally removed news from training, because:
  - it introduces too much noise;
  - it is too much subjective and depending on human judgement;
  - timing of news appearance can be misleading and can spoil news' importance;
  - numbers of news features were too big for a reliable simple model.

# ANN Architecture

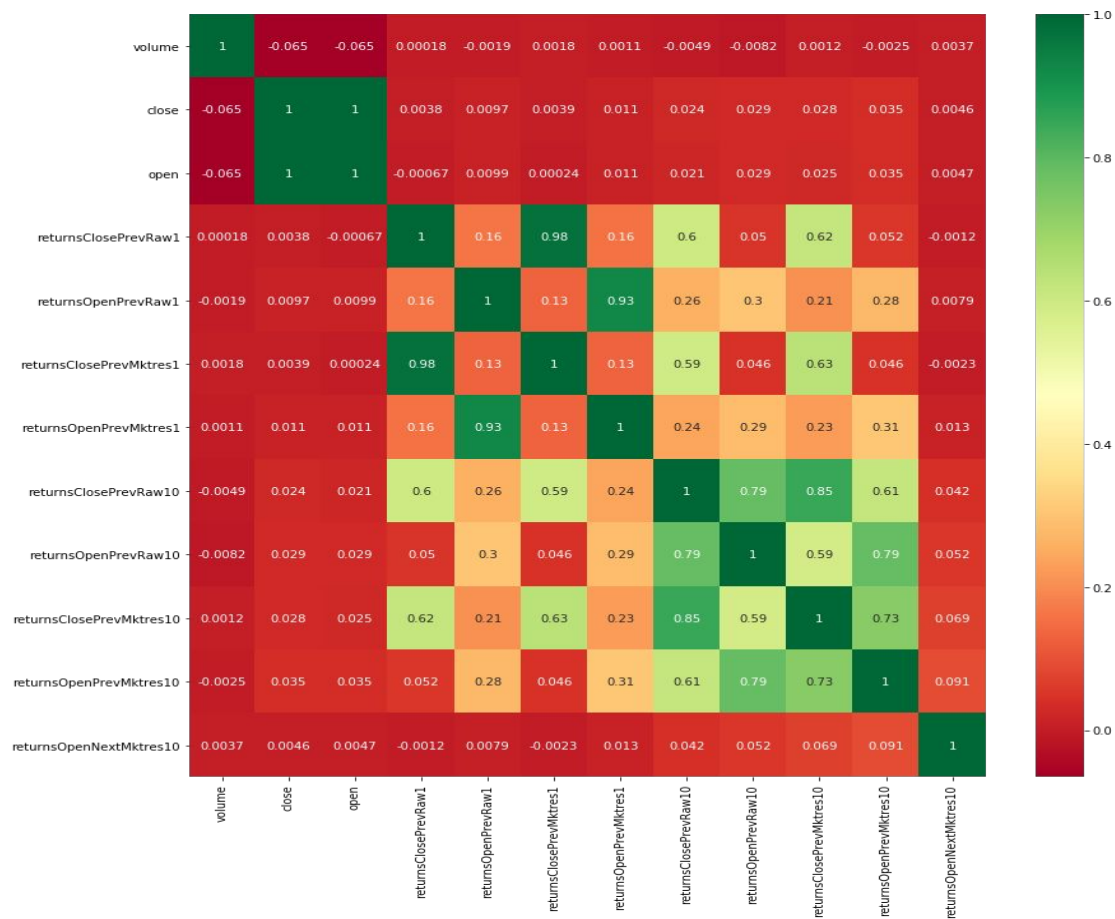# Features Selection / Engineering



Variable Importance Plot

Variable importance plot to be used together with a correlation heatmap (see next slide) in order to exclude variable with a very high cross correlation (dependency to each other). This can be used for choosing data features for a simplified model.

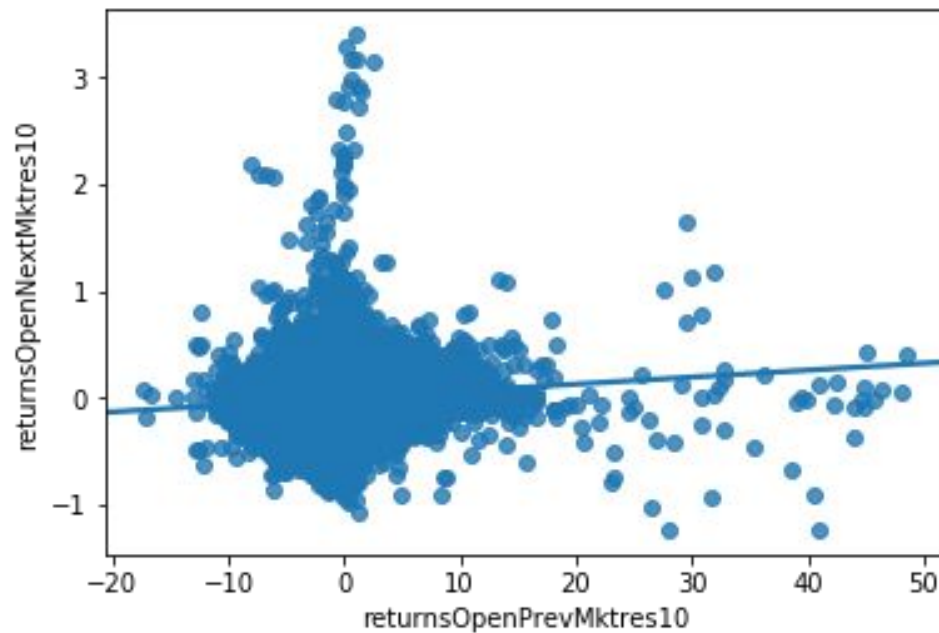# Features Selection / Engineering



Variable Cross-Correlation Heatmap

Most important features:

- **'returnsOpenPrevMktres10'** - positive correlation to target (9.1%) with largest importance;
- **'volume'** - low correlation to feature #1 and other features, but relatively high importance;
- **'returnsOpenPrevRaw1'** - in the same manner selection was performed;
- **'close'** - reason the same as above-mentioned.

Relation and linear regression of the most important feature with the target
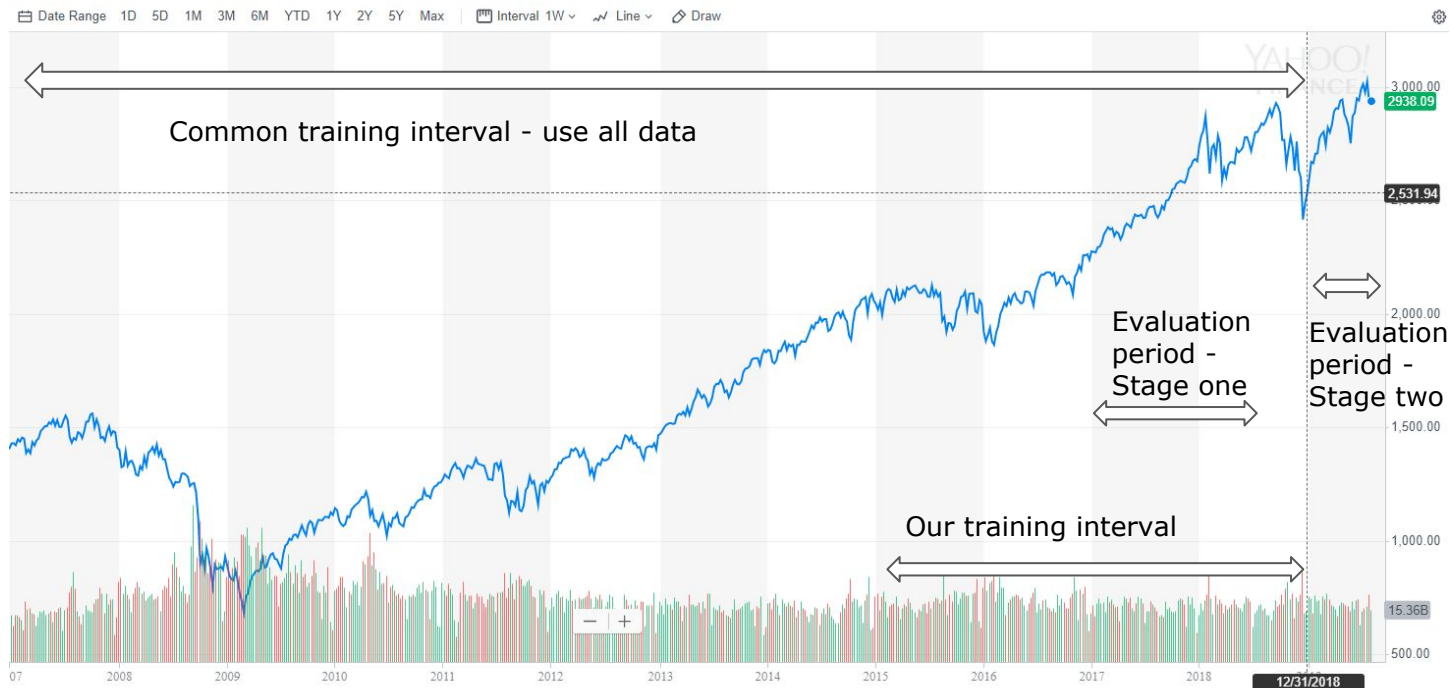
- Training methods:
    - used only data from 2015-01-01
    - 75% data used for training
    - 10 epochs
    - EarlyStopping callback
    - not used hyperparameter tuning
    - Adam optimizer with custom loss function
- Did you ensemble?
    - No. Try to keep model as simple as possible in order not to overfit.
    - Ensembling current model with news model generated always the worst results.

- What set you apart from others in the competition?

  - Simplicity;
  - Training was started from beginning 2015 year. Financial system is structurally unstable, so one should operate only on the most recent data;
  - Custom loss function.

  After many years of experience in trading, we observed that the past performance of algo has almost nothing to do with live evaluation and start trading, so we decided to keep our model simple and test on data in short time interval.

# Important and Interesting Findings

## Train and test period versus evaluation period based on S&P 500 index



Source: Yahoo Finance

Looking in the picture above, it can be seen that starting from 2018 price patterns seem to be sharper than ever before.

- A subset of features that could contribute to reach 90-95% of your final performance:

  - Please, forget about 3+ years past data. Train your model on the last 2-3 years data.
  - In custom objective function maximizes the Sharpe ratio.
  - In order not to overfit, keep model simple. Please, rather do not use ensemble of models.
  - Reduce number of features to 4 (see slide 13th).
  - Score for simplified model should be around: 0.6 - 0.8.

# Thank you very much!