



Lucena Research
Predictive Analytics

Application of Alternative Data and Convolutional Neural Networks for Forecasting Stock Prices

Erez Katz

CEO & Co-Founder, Lucena Research

About

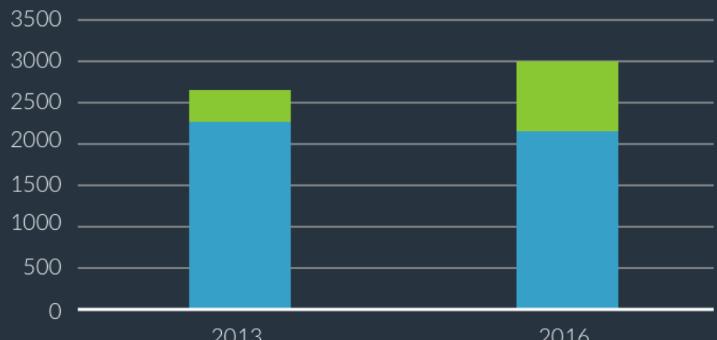
About Lucena

- Leader in Machine Learning and AI
- Connect big data with investment professionals



Erez Katz, Lucena CEO

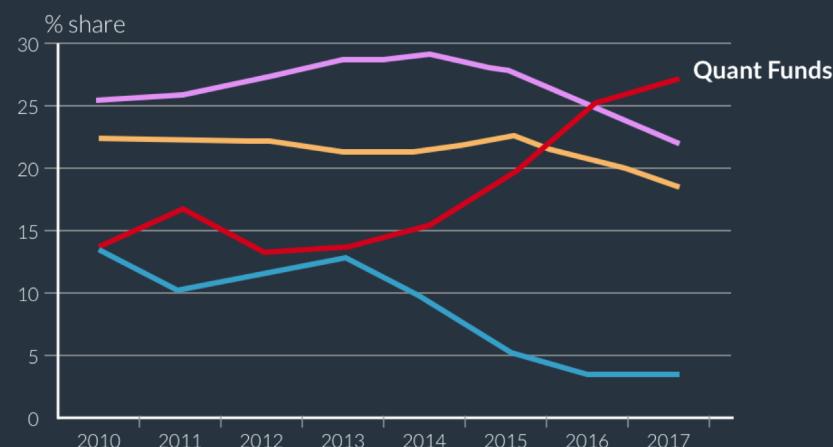
Buy Side Alternative Data Adoption



Quant Hedge Funds

Non-Quant Hedge Funds

Image Source: Jefferies, Hedge Fund Research



Quant Funds

Other Hedge Funds

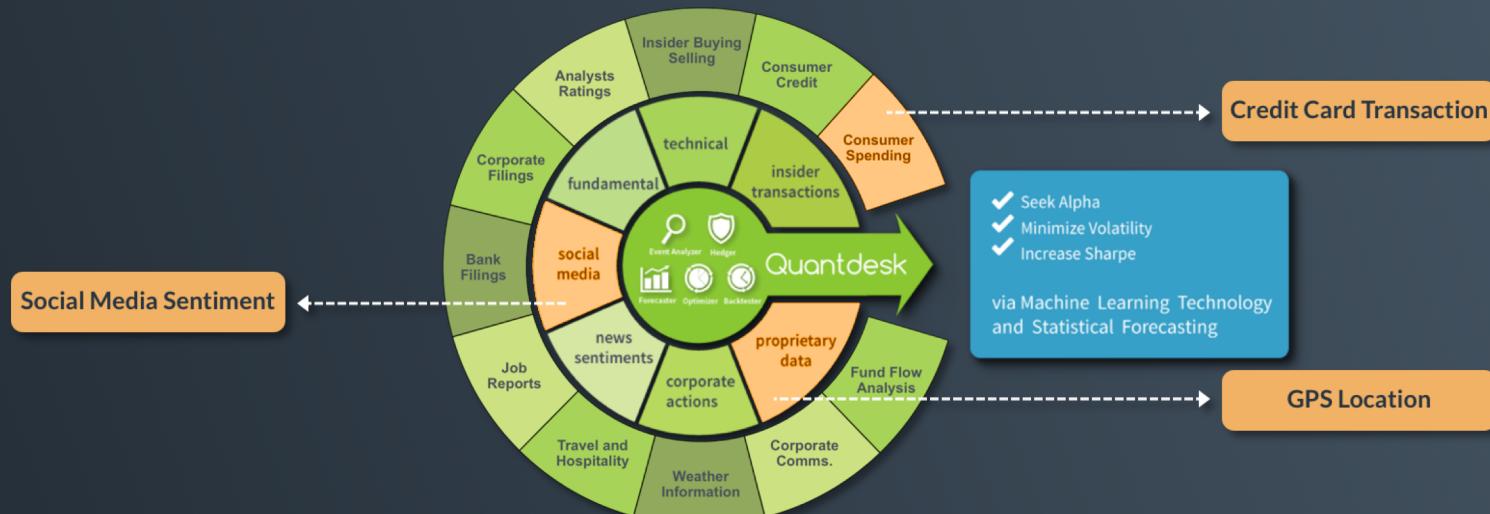
Traditional Asset Managers

Banks & Investment Banks

Image Source: WSJ, Tabb Group

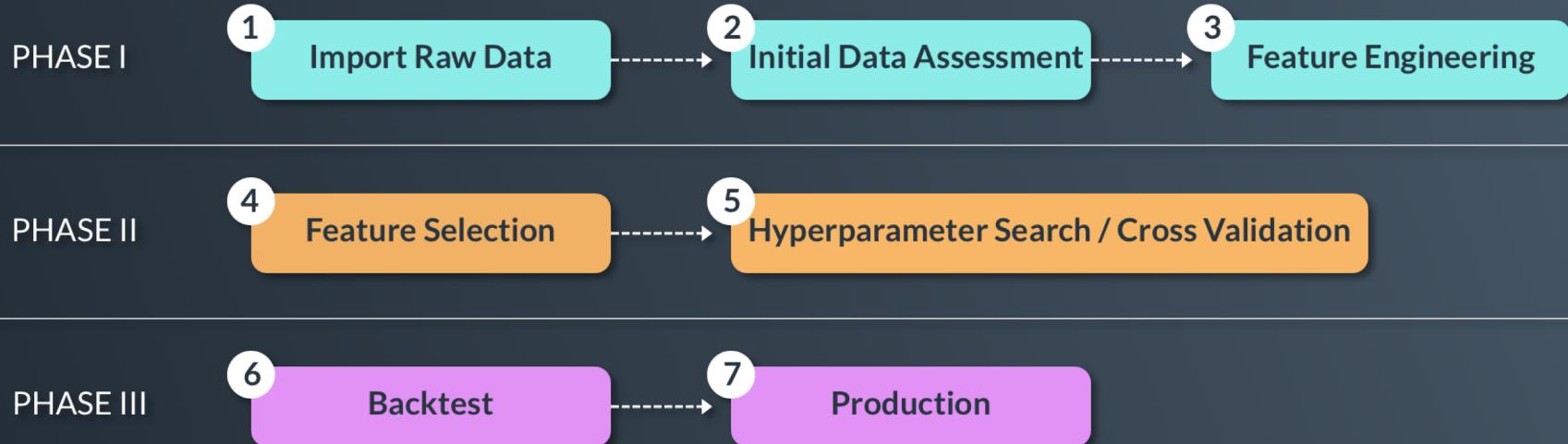
Big Data & Alternative Data

- Traditional Data Types: Mainly Technical and Fundamentals
- Equity Specific Alternative Data: Corporate Action, Insider Buying Selling
- Non Financial Data Repurposed for the financial markets: Credit Card, Weather Data
- Unstructured Data: News Feed, Social Media Feed, Speeches, Earnings Calls



Data Validation Process

- Critical to follow a regimented process by which you can reliably consume, validate, and enhance data in order to generate predictive signals.

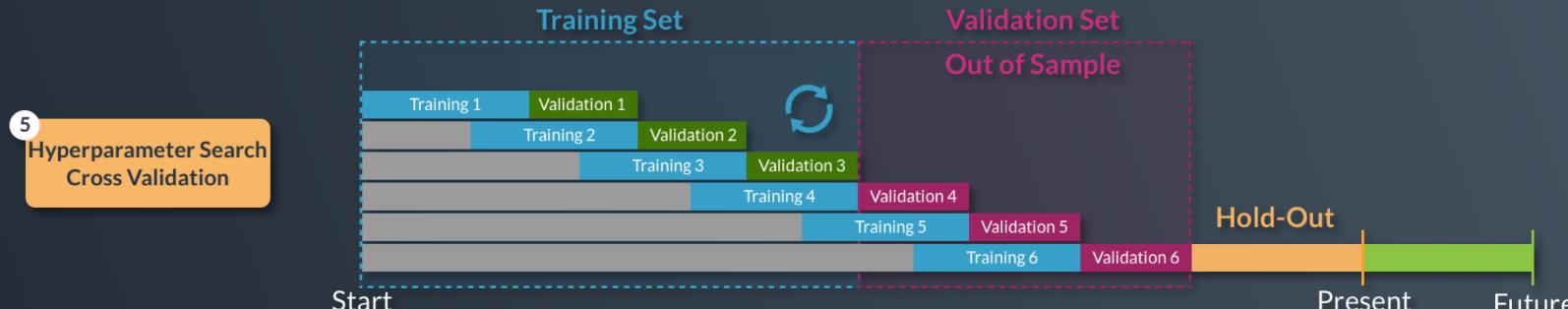


Data Validation Process

PHASE I



PHASE II



Data Validation Process



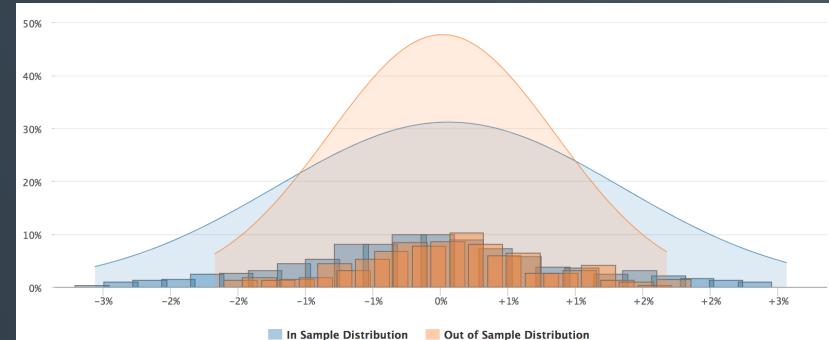
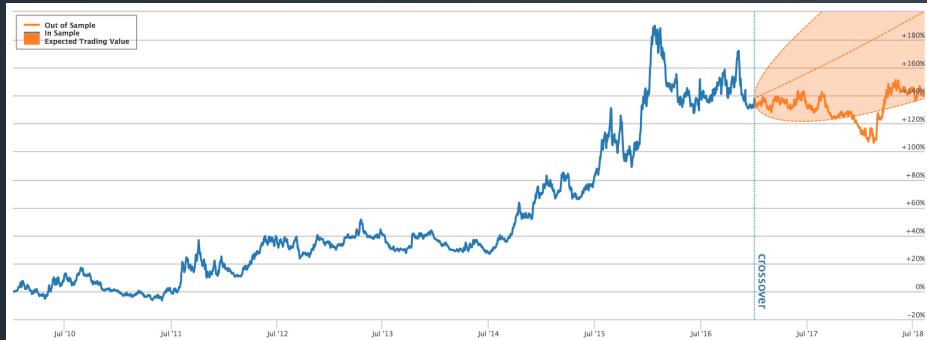
Automating Data Qualification? Why?

Qualification Statistics

	Value	Best	Grade
Percentage of NaNs	0.0799	Lower	✗
Percentage of Symbols Mapped	0.84	Higher	✓
Likelihood of Survivorship Bias	0.83	Lower	✗
Best Decile Sharpe	0.7353	Higher	✓
Signal Distribution Normality	0.91	Higher	✓
Percentage of Missing Trade Days	0.16	Lower	✗
Percentage of Extra Trade Days	0.01	Lower	✓
Percentage of Days with Multiple Signals per Asset	0.2037	Lower	✗
Average Data Frequency (Days)	65.8316	Closer to 1	✗
Data Type Prevalence	1.00	Higher	✗
Likelihood of Underlying Model Change	0.64	Lower	✗

- Completeness
- Authenticity
- Consistency
- Survivorship Bias

Below is an example of an inconsistent model:



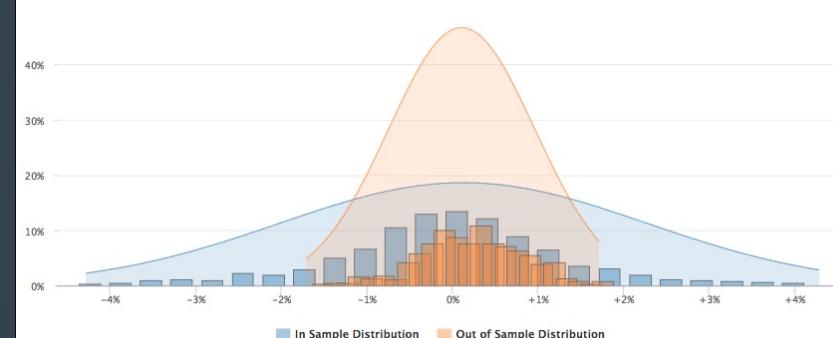
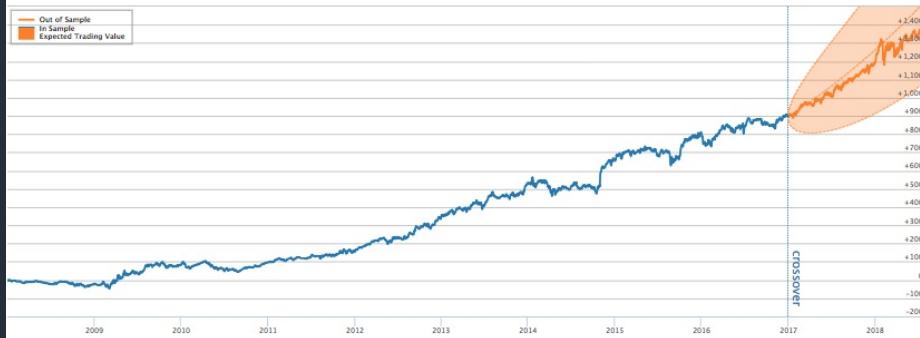
Data Qualification What To Look For?

Qualification Statistics

	Value	Best	Grade
Percentage of NaNs	0.0799	Lower	✗
Percentage of Symbols Mapped	0.84	Higher	✓
Likelihood of Survivorship Bias	0.83	Lower	✗
Best Decile Sharpe	0.7353	Higher	✓
Signal Distribution Normality	0.91	Higher	✓
Percentage of Missing Trade Days	0.16	Lower	✗
Percentage of Extra Trade Days	0.01	Lower	✓
Percentage of Days with Multiple Signals per Asset	0.2037	Lower	✗
Average Data Frequency (Days)	65.8316	Closer to 1	✗
Data Type Prevalence	1.00	Higher	✗
Likelihood of Underlying Model Change	0.64	Lower	✗

- Completeness
- Authenticity
- Consistency
- Survivorship Bias

Below is an example of a consistent model:



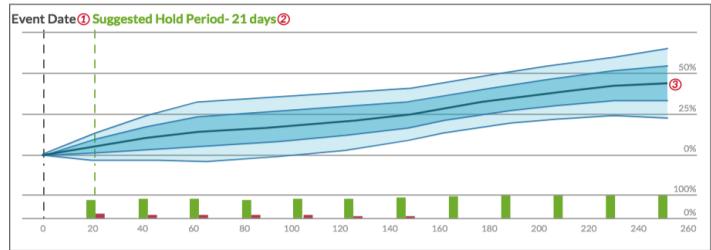
Data Matching Why?

The Following Data Sources Are Recommended

Company	About	Relevancy Score
Credit Card Transaction	Company A provides sentiment scores that can be utilized for market/index timing transactions across time horizons.	
Social Media Sentiment	Since 2000, Company B has been the world leader in capturing sentiment information from retail investors, owning the world's largest database of retail investor sentiment poll opinions.	
GPS Location	Company C is a registered investment advisory firm. Over the last thirty years, their research platform, SMA/UMA strategies, indexes, and team of analysts have created more informed, conversations for advisors, investment managers and their clients.	

- Is the data effective for my constituents list?
- Is the signal suitable for my investment style
 - Hold Time
 - Long/Short risk tolerance
 - Concentration and turnover requirements

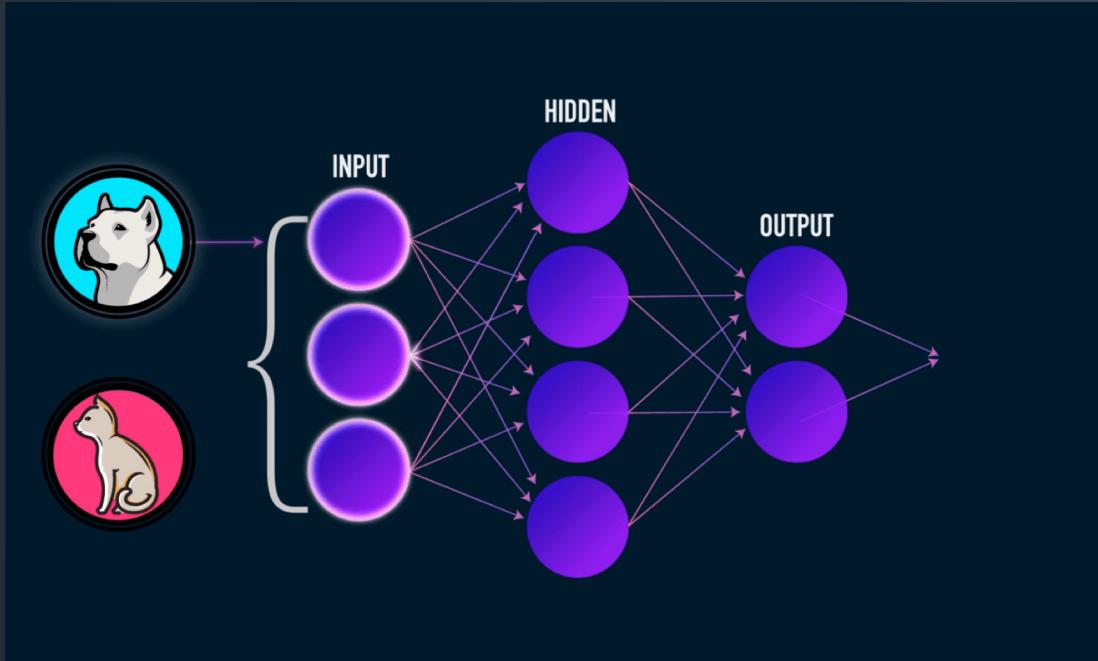
Event Study Analysis



Backtest Charts



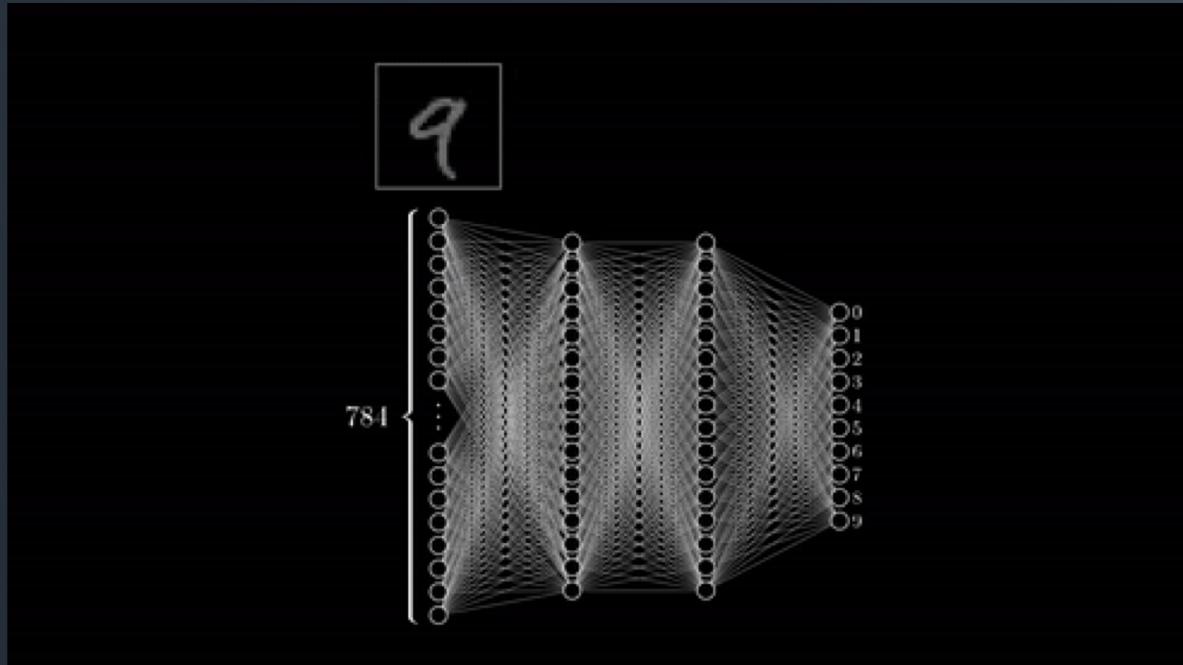
What is a Deep Neural Network



Animated gif by James Loy. Source: <https://towardsdatascience.com/how-to-build-your-own-neural-network-from-scratch-in-python-68998a08e4f6>

- Input Layer of Neurons.
- Multiple Hidden Layers.
- Feeds forward information.
- Output layer.

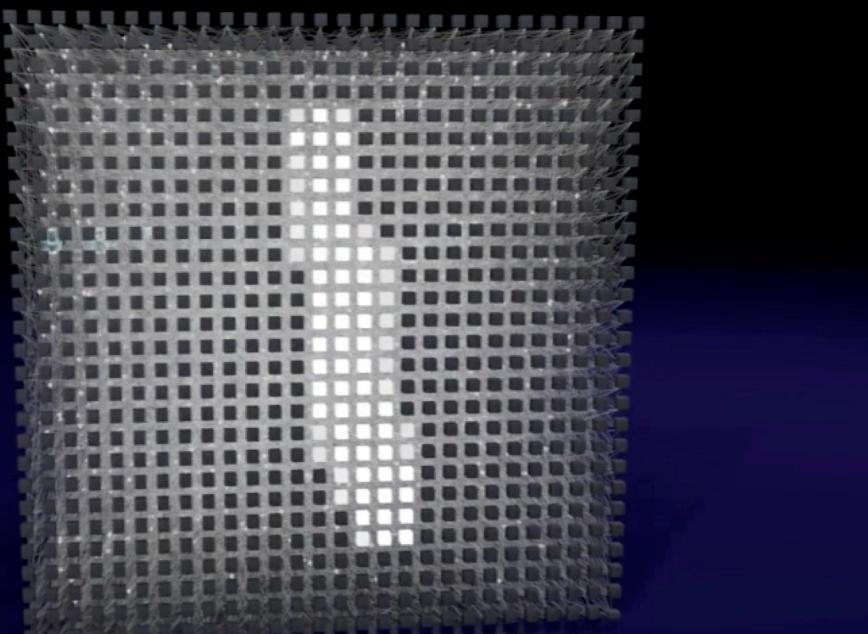
How Does The Network Learns?



- An image of hand written digit One layer identifies edges.
- Subsequent layers identify patterns.
- Output layer identifies the full object.

Video by 3Blue1Brown. Source: <https://www.youtube.com/watch?v=aircAruvnKk>

Realistic Visualization of image classification



- Real world network is complex
- Feed forward Neuron Activation
- Output layer identifies the object
- Backpropagation and gradient descent

Video Source: <https://www.cybercontrols.org/neuralnetworks>

Example of a Modern CNN

- CNN forms additional tasks by filtering special patterns.
- Layers are added to down sample a complex image.

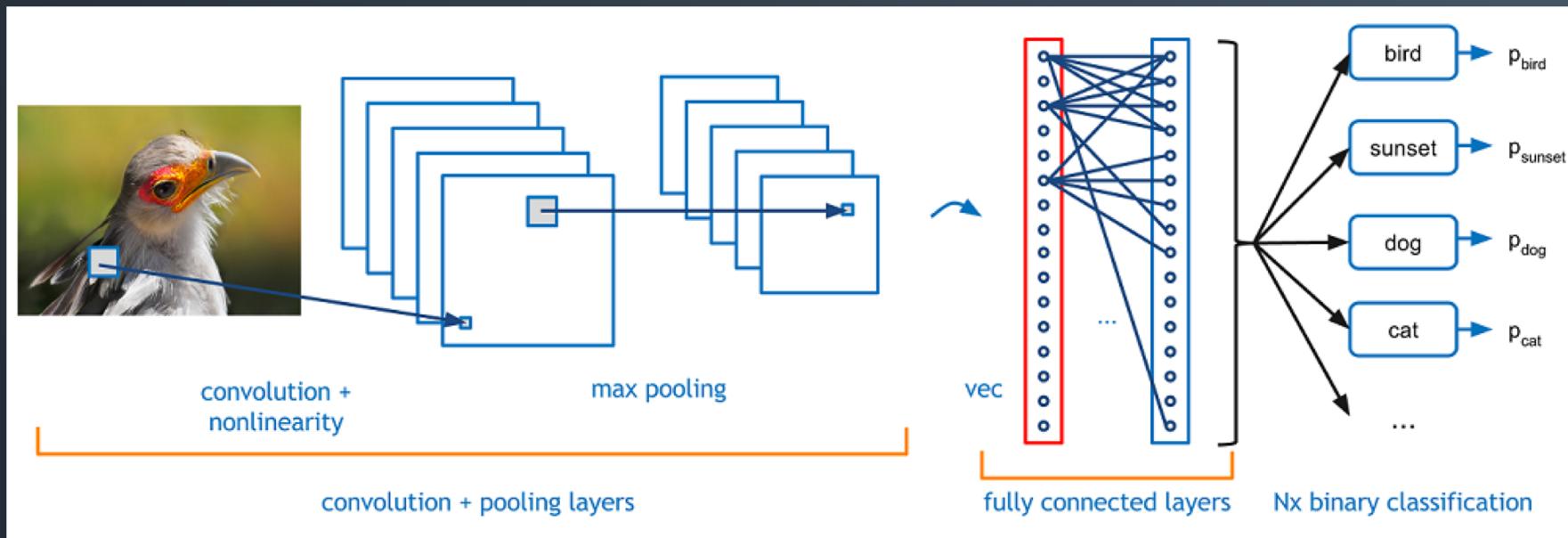
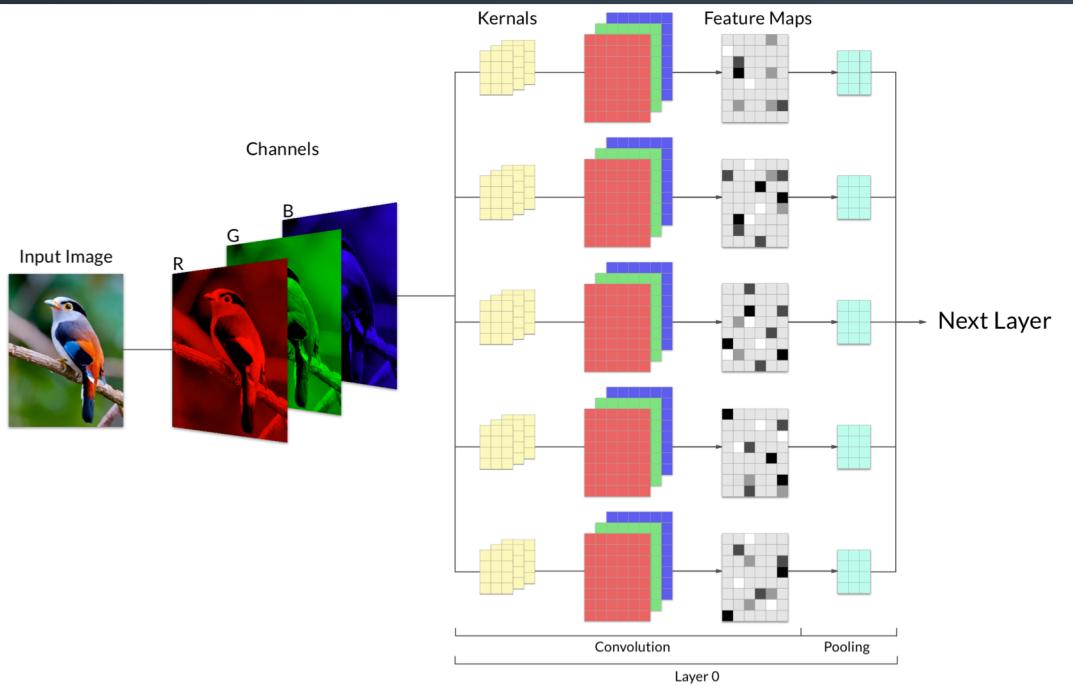


Image by Rob Hess. Source: <https://code.flickr.net/2014/10/20/introducing-flickr-park-or-bird/>

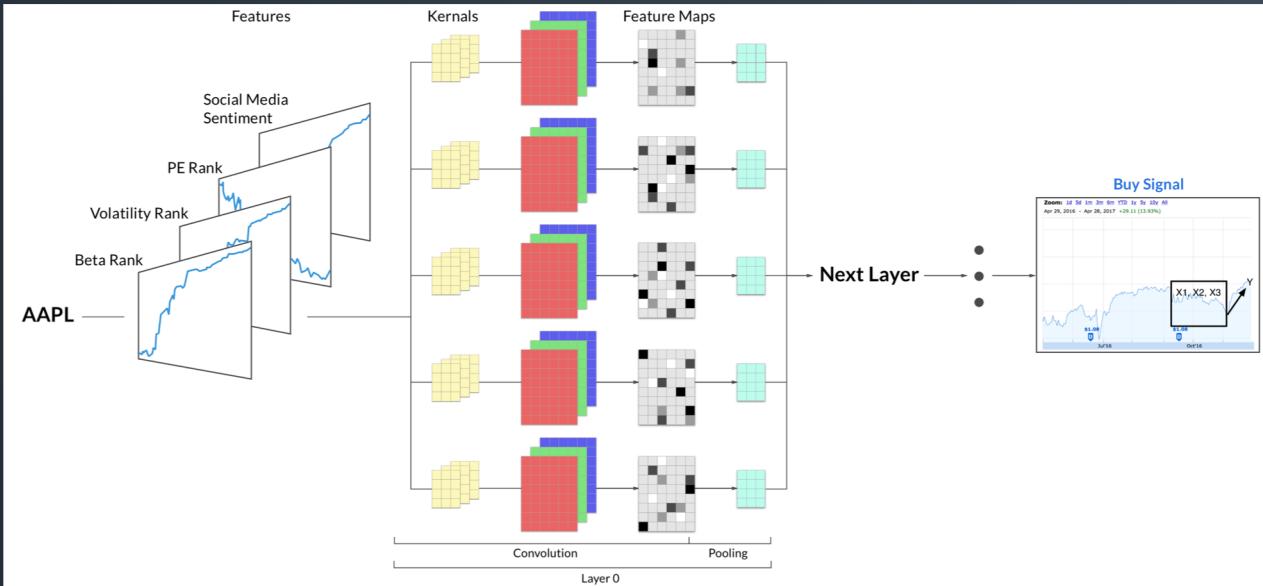
Multichannel Breakdown of an Image



- High degree of accuracy in modern Deep Learning and CNN for image classification or object detection.
 - Autonomous Vehicles
 - Facebook Tagging
 - Amazon purchase habits
 - Netflix
 - Autocomplete Text

Can we apply the same principles to stocks?

We Can Apply the Same Concept to Stocks



Multiple timeseries images of feature values over time.

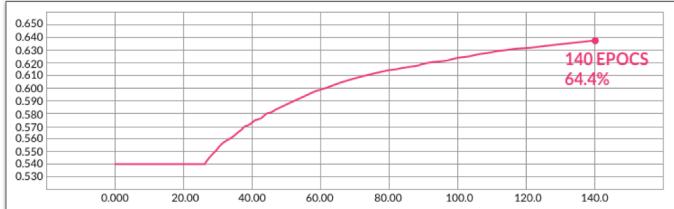
- Technical Features
- Fundamental Features
- Macro Econ Features
- Alternative Data Features
- Social Media Sentiment
- Corporate Action
- Analysts Recommendations
- News Feed Sentiment

Google's Tensor Board

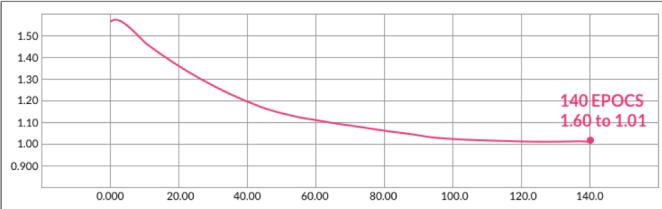


Lucena Research
Predictive Analytics

Validation Accuracy

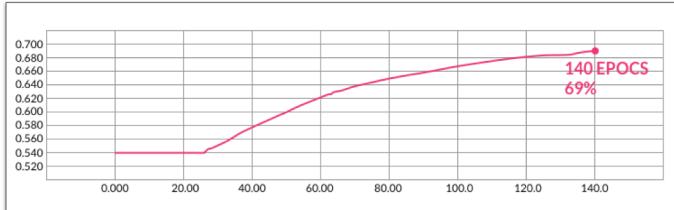


Validation Loss

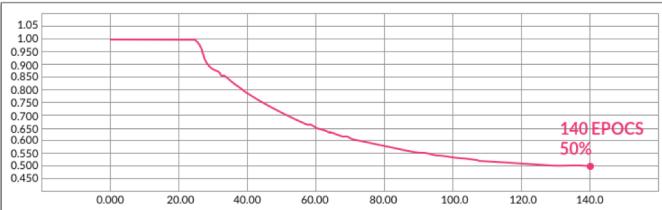


Out of Sample

Validation Precision



Validation Percent 1 Label

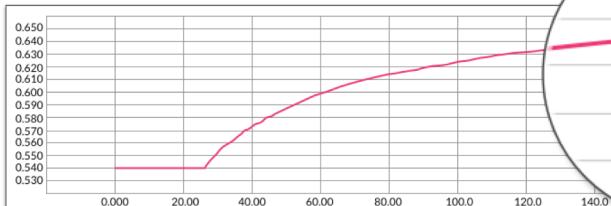


Google's Tensor Board

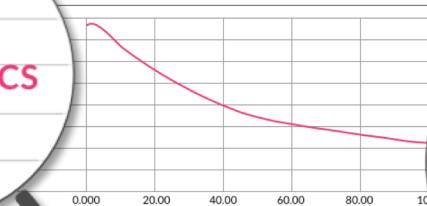


Lucena Research
Predictive Analytics

Validation Accuracy

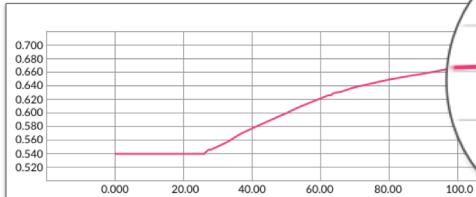


Validation Loss

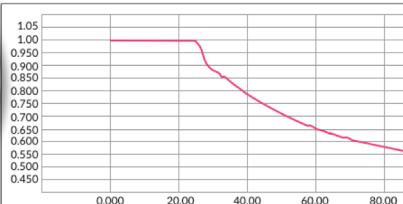


Out of Sample
(Zoom)

Validation Precision



Validation Percent 1 Label



Lucena Research
Predictive Analytics

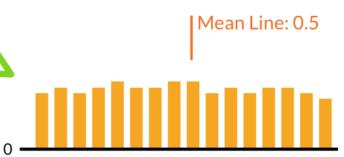
Making It Easier for the Machine to Learn

Uniformly Distributed Features

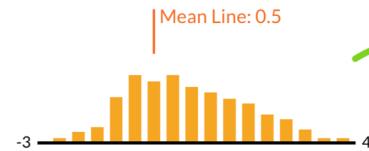
Alpha Value:



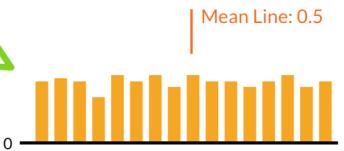
Alpha Rank:



Sharpe Value:

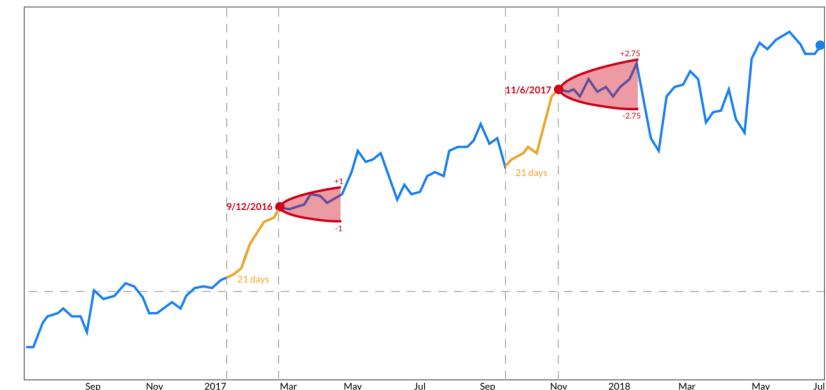


Sharpe Rank:



Dynamic Label Data using ATR

APPL



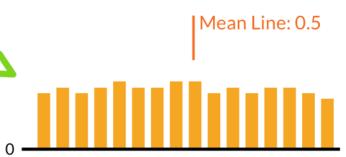
Making It Easier for the Machine to Learn

Uniformly Distributed Features

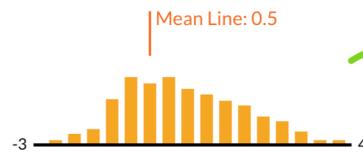
Alpha Value:



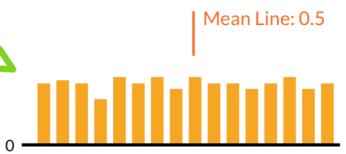
Alpha Rank:



Sharpe Value:

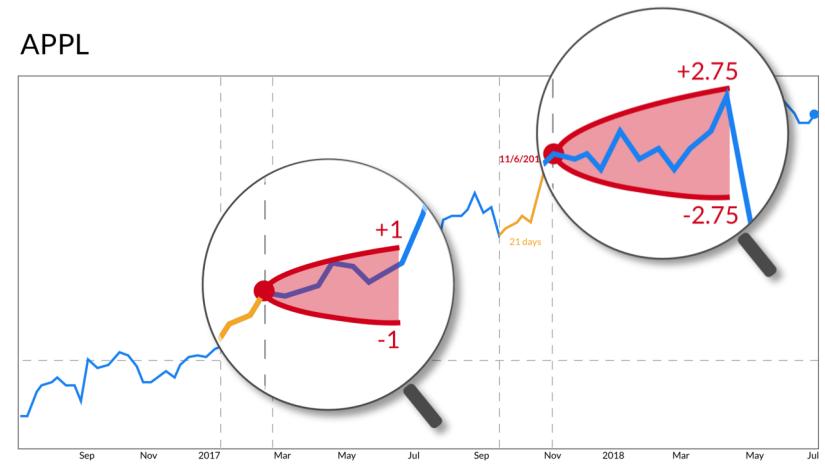


Sharpe Rank:



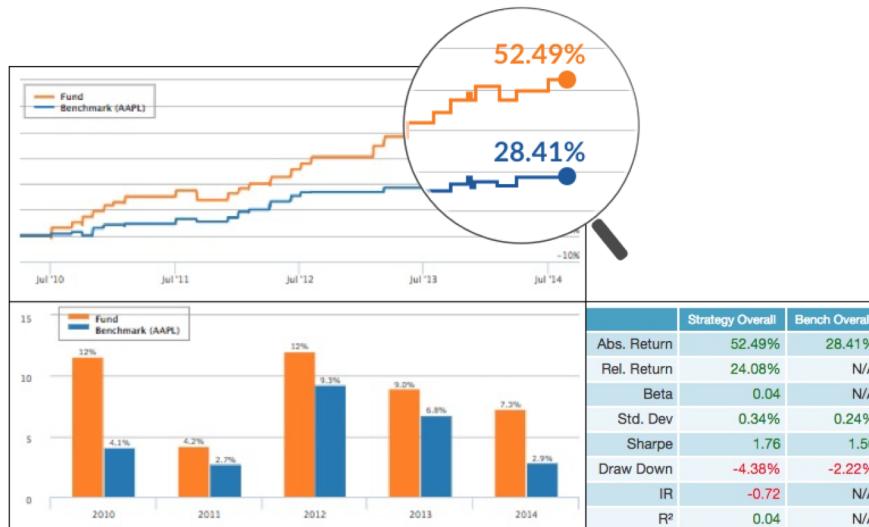
Dynamic Label Data using ATR

APPL

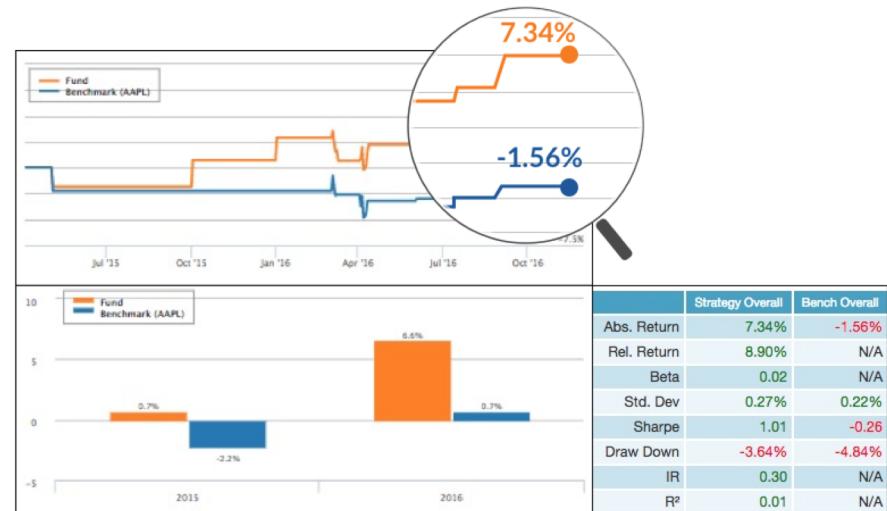


APPL Backtest

Results are Exposure Adjusted
 In Sample: 4/6/2010 to 9/30/2014



Results are Exposure Adjusted
 Out of Sample: 4/6/2015 to 12/1/2016



Thank you!

erez@lucenaresearch.com