

Using social media data to enhance predictive models in finance

Nicholas Mitchell

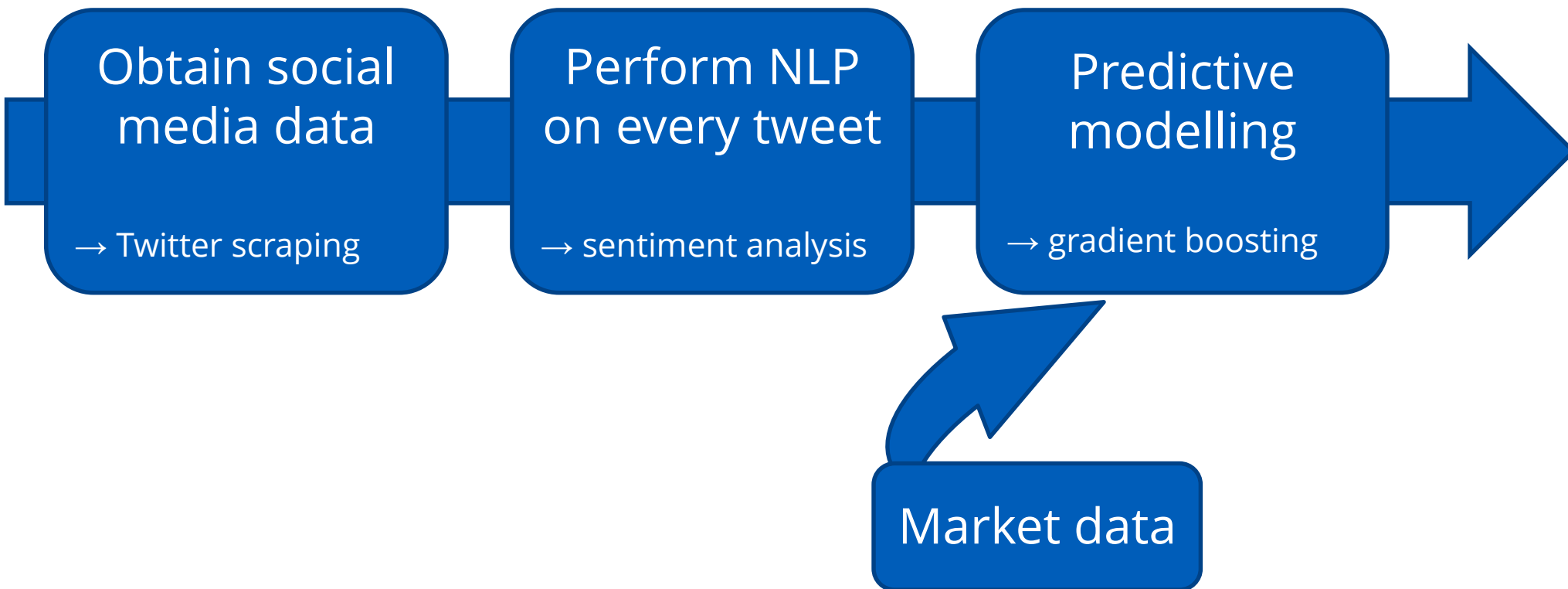
Munich, September 2016

We aimed to answer the question(s):

Can social media data improve predictive models of financial markets?

...by how much?

The Roadmap



Tapping the *Twitterverse*

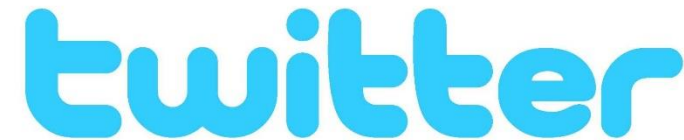
1. Purchase them

- Every tweet since Twitter inception
- Enterprise solution
- Expensive



2. Official API

- Free
- Fast
- Limited to last 10 days of data



3. Unofficial API... Scraping

- Free
- (currently) slow
- (currently) unlimited data*



Scraping

- Twitter's advanced search:
 - Choose a search term
 - Select a date range
 - Tweets are displayed
- Page loads - starting with youngest tweets
- Scroll downwards → backwards in time
- Dynamically loading page
- Save all HTML code from Browser – WYSINAYG!

ONE DOES NOT SIMPLY



GET SOME TWITTER DATA

Search

Parsing

- Take raw HTML → parse using Xpath-Element-trees

- Create a one-tweet-per-line CSV file

Number of search terms: 13

- In addition to tweet text, extract useful meta-data:

Date range: 14.01.2013 – 11.09.2015

– Date

– Number of likes

Total timeline: 982 days (695 weekdays)

– Number of retweets

– Unique tweet ID

Total tweets obtained: 2, 350, 217

– Username & ID

– Geo-coordinates

Cleaning the tweets

- Required to facilitate accurate sentiment analysis

- Tweets had to go from:

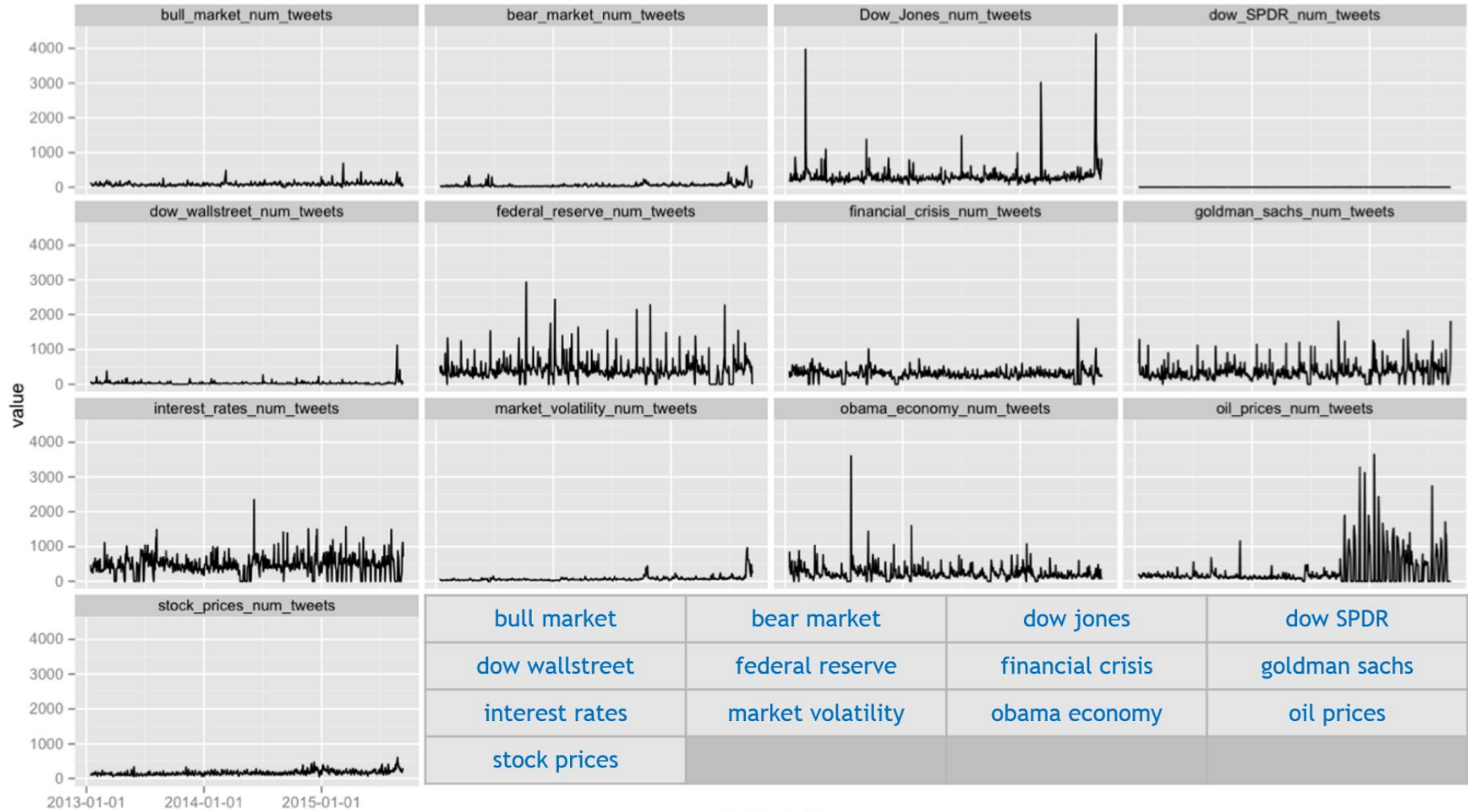
""""I wonder what people think about the
""Death Cross"" now? :)^M #trendfollowing

- To this:

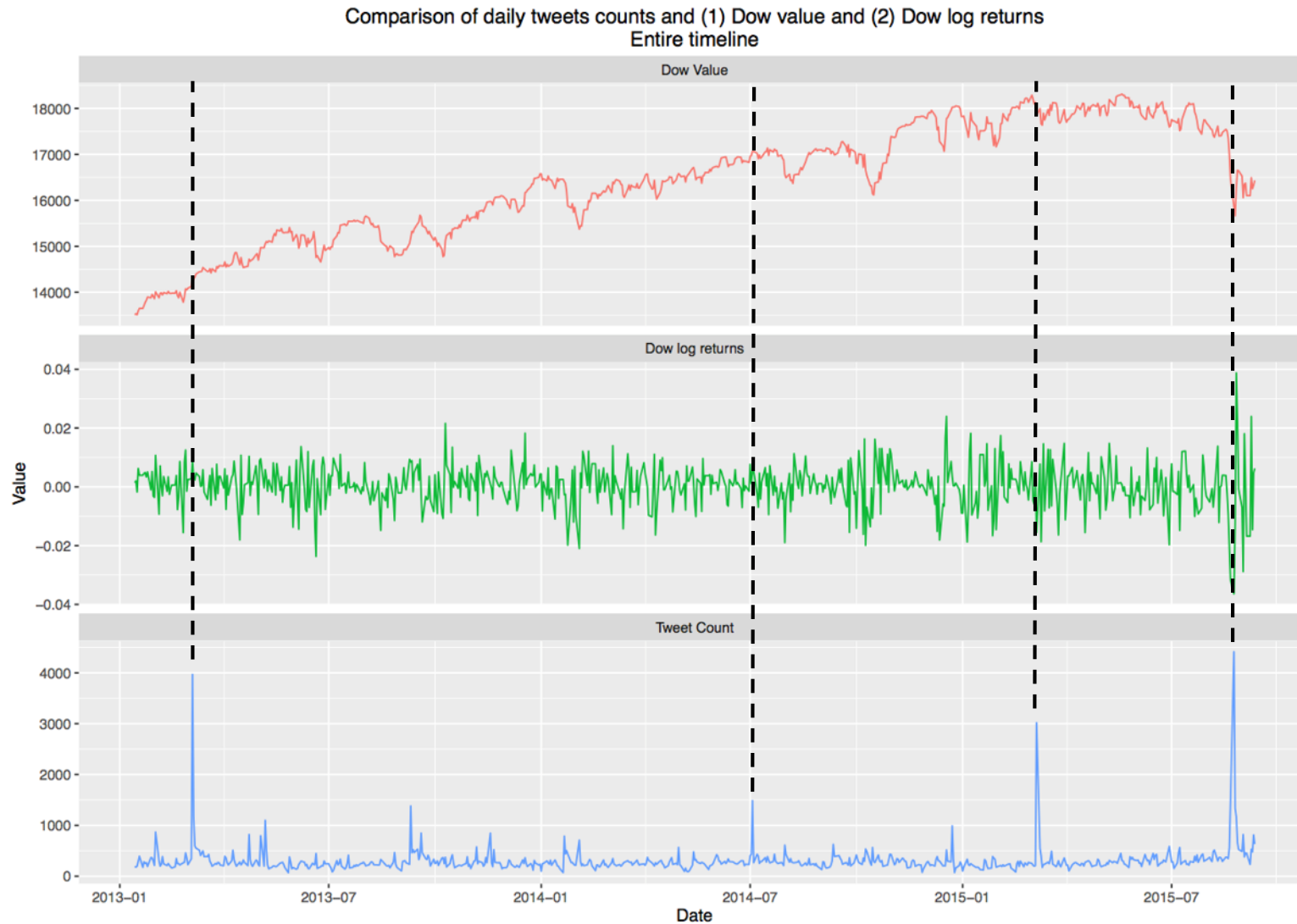
I wonder what people think about the
Cross now? :) trendfollowing

→ RegEx using Perl engine
→ hexadecimal char definitions
→ using an ASCII table

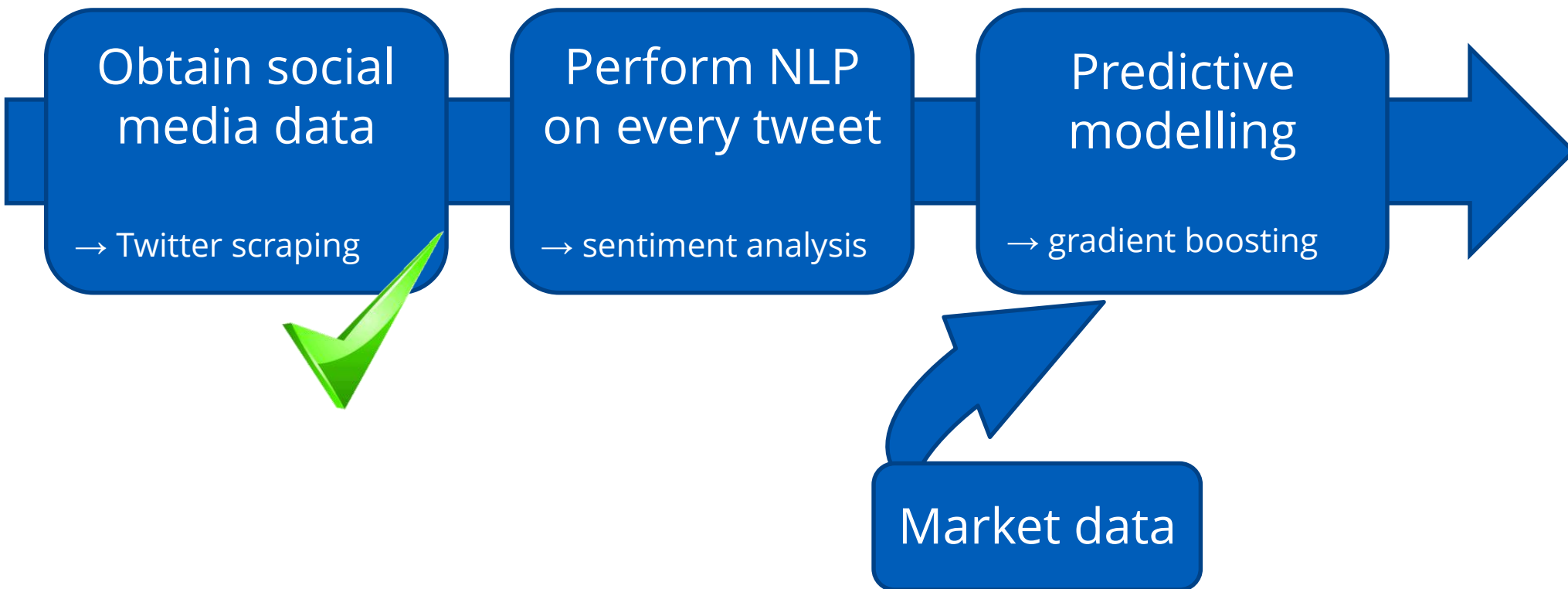
Breakdown of search terms and tweet count



Inspecting the Twitter data (1)

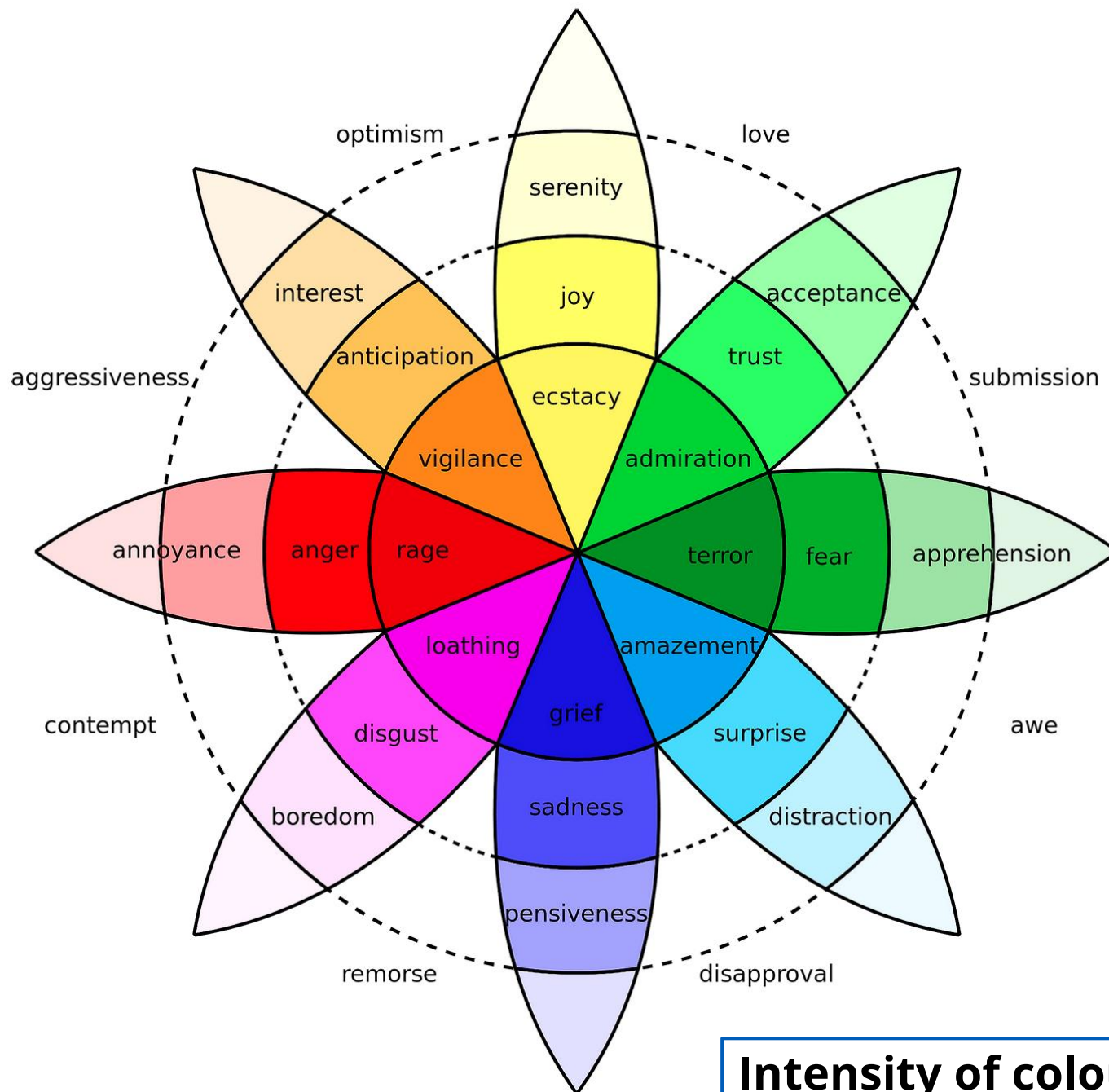


The Roadmap



Sentiment Analysis (1)

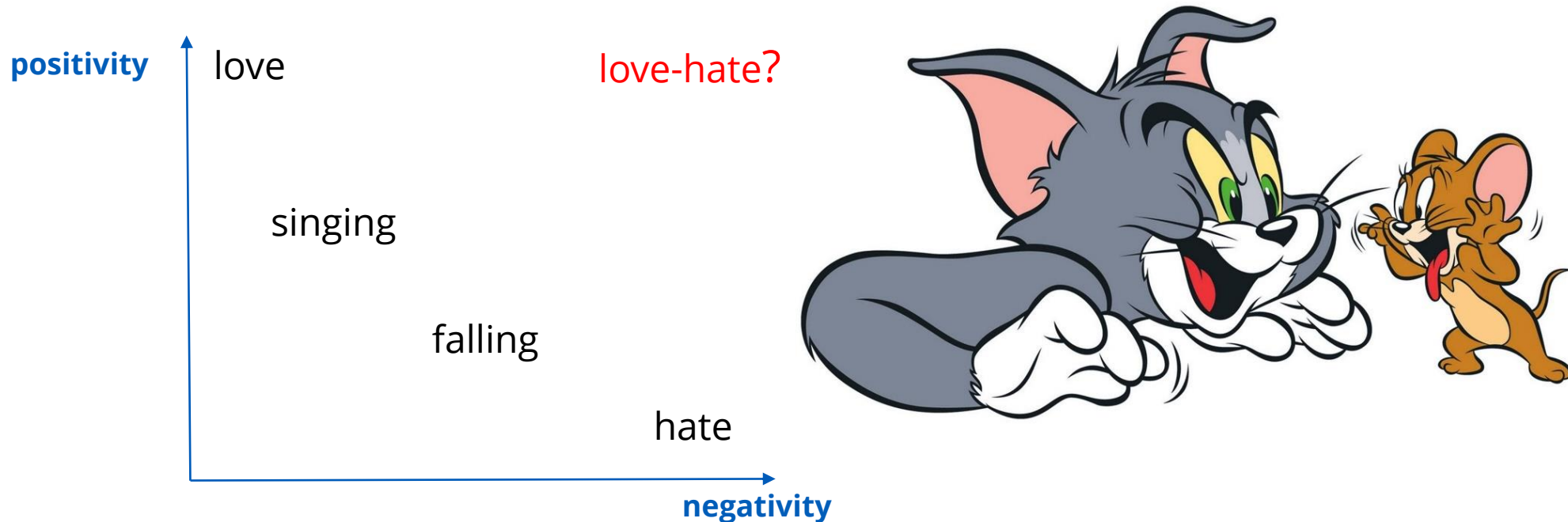
- 5 differing models
 - Each focussed on short texts / social media data
 - Employ diverse linguistic approaches
 - Example approaches:
 - Grammar based – scored word lists for nouns, verbs, adjectives, etc.
 - Informal text – scores for smileys, slang and profanity
 - Pure word list, create by a community/mechanical Turk
- example: *Plutchik's Wheel of Emotion*



Intensity of colour \approx intensity of emotion

Sentiment Analysis (2)

- Basic example: score each word on two dimensions:

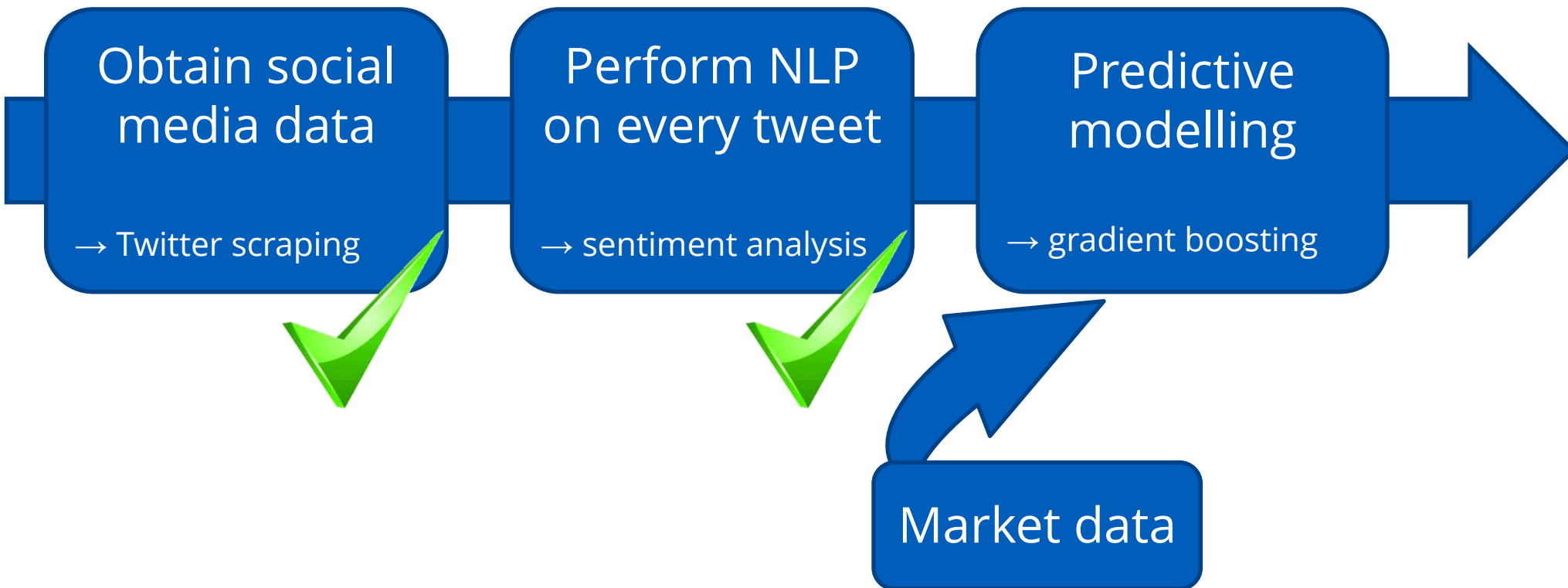


- Example:

Rupert loves safe investments, but Niko loves risky OTM options

$(1,1)$ $(4,1)$ $(3,2)$ $(1,1)$ $(1,1)$ $(1,1)$ $(4,1)$ $(1,3)$ $(1,1)$ $(2,1)$ $\Sigma = (19, 13)$

The Roadmap



Market data

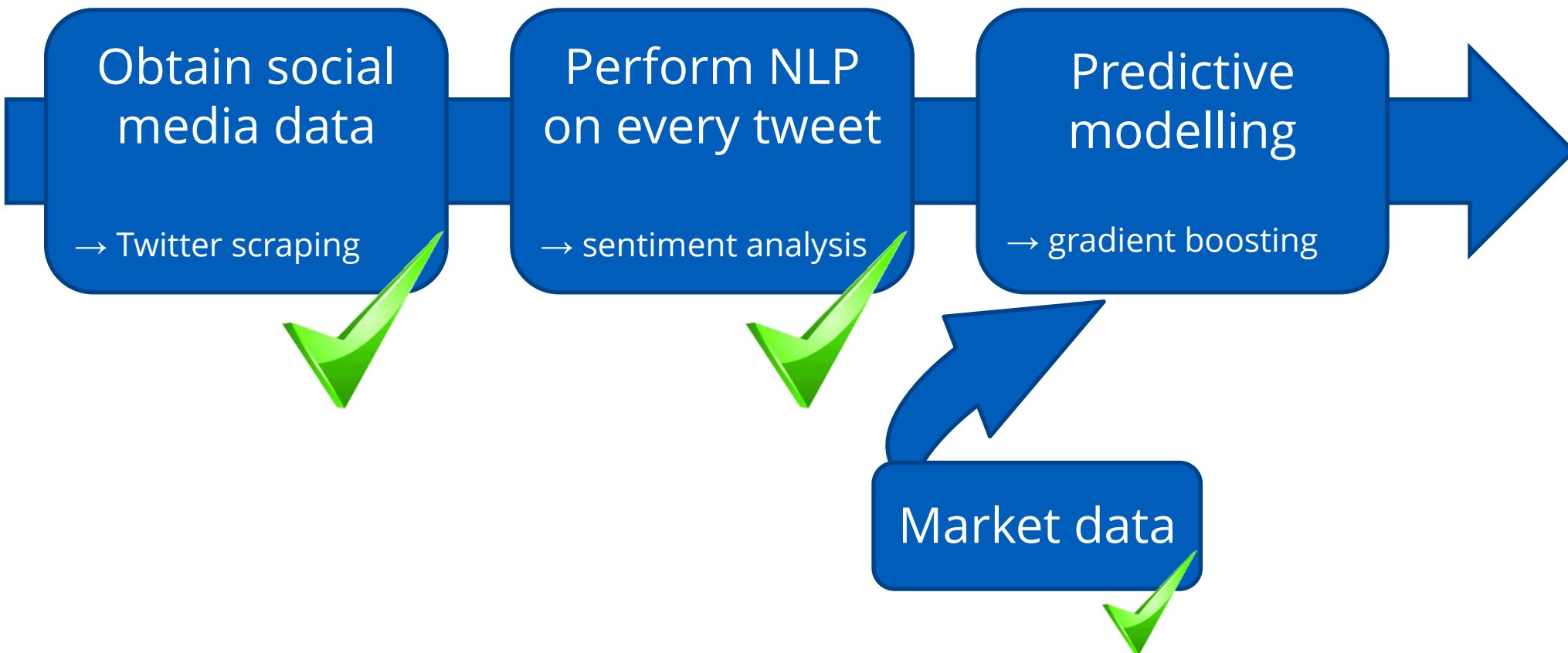
Commodities	Currency pairs	Fixed income
Gold spot	USD-AUD	U.S. Zero-coupon 1Y
Gold 3M	USD-CAD	U.S. Zero-coupon 2Y
Copper spot	USD-EUR	U.S. Zero-coupon 5Y
Copper 3M	USD-GBP	U.S. Zero-coupon 10Y
Oil (WTI)	USD-JPY	U.S. Zero-coupon 15Y
Natural gas		U.S. Zero-coupon 20Y
Indices	Volatility indicators	ETF
DAX	VIX (S&P500)	MSCI Emerging Markets
Dow Jones	Gold spread	
FTSE100	Copper spread	
Nikkei 225		
S&P500		
Shanghai SE		

Datasets for comparison

		Σ
• Traditional market data		
– Small selection	→ <code>traditional_{small}</code>	6
– Everything	→ <code>traditional_{large}</code>	36
• Sentiment data only		
– Aggregated	→ <code>sentiment_{small}</code>	22
– Individual	→ <code>sentiment_{large}</code>	100
• Market + social media data	→ <code>combined</code>	142

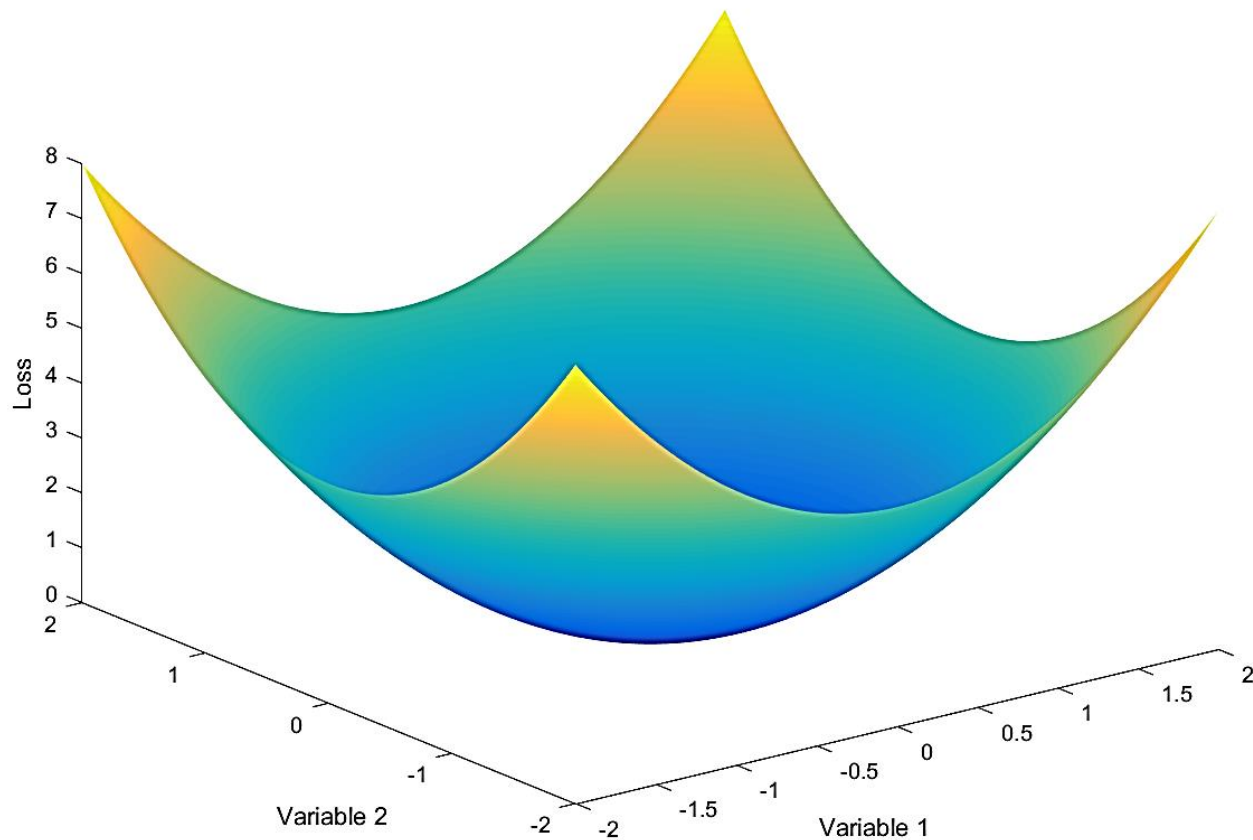
→ **“Correlation cut-off” used, reducing datasets for final model**

The Roadmap



Modelling: Gradient Boosting

- Define loss function
 - convex
 - Convert to weak learners
 - Continuous
 - Boosting
-
- Fit so
 - Measure
 - Update
 - Repeat
 - Until at bottom of the bowl

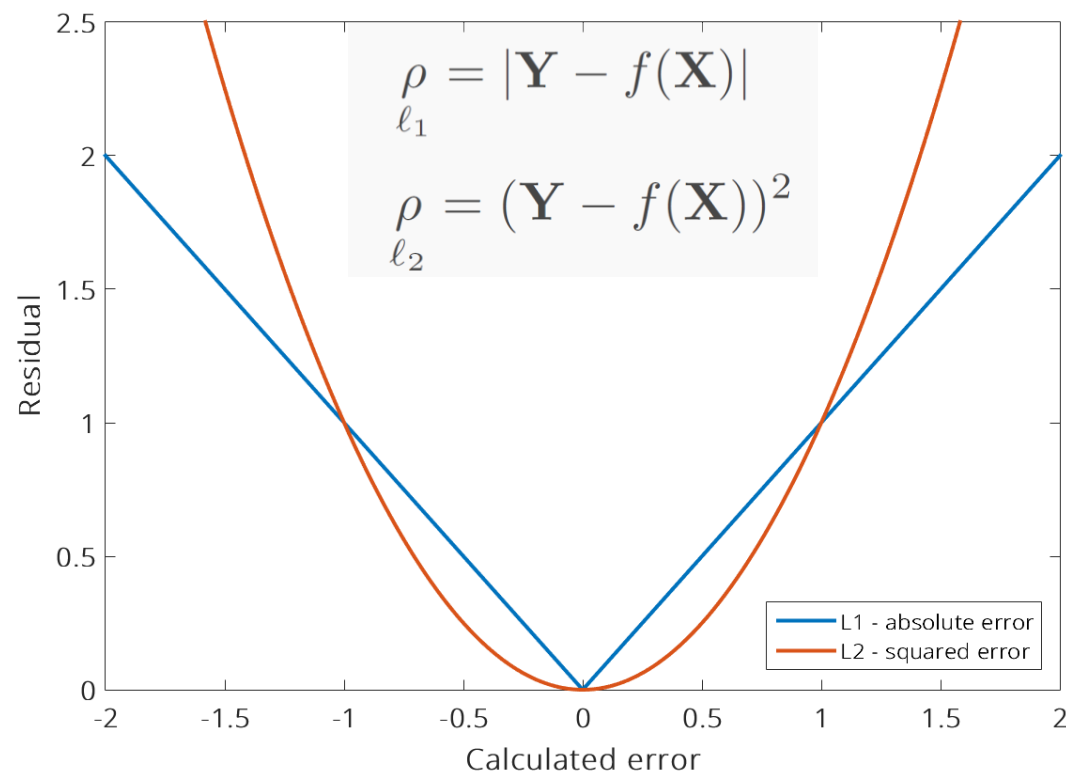


Gradient Boosting: theory in a nutshell

$$f^* := \operatorname{argmin}_{f(\cdot)} \mathbb{E}_{\mathbf{Y}, \mathbf{X}} [\rho(\mathbf{Y}, f(\mathbf{X}))]$$

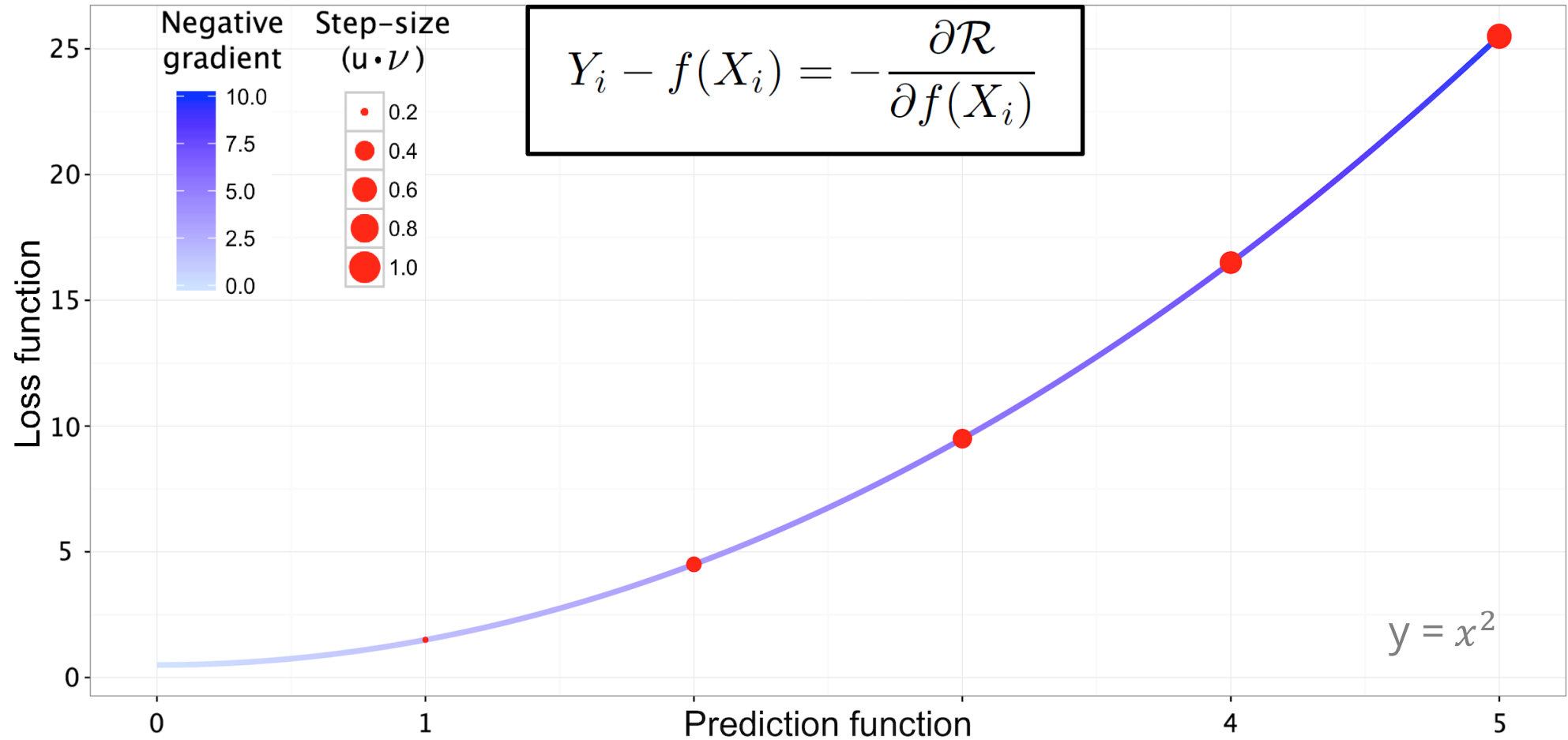
$$\mathcal{R} = \frac{1}{n} \sum_{i=1}^n \rho(Y_i, f(X_i))$$

$$\mathcal{R} = \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2$$



$$\frac{\partial \mathcal{R}}{\partial f(X_i)} = \frac{\partial}{\partial f(X_i)} \left(\sum_{i=1}^n \rho(Y_i, f(X_i)) \right) = \frac{\partial}{\partial f(X_i)} (\rho(Y_i, f(X_i))) = f(X_i) - Y_i$$

Adaptive descent



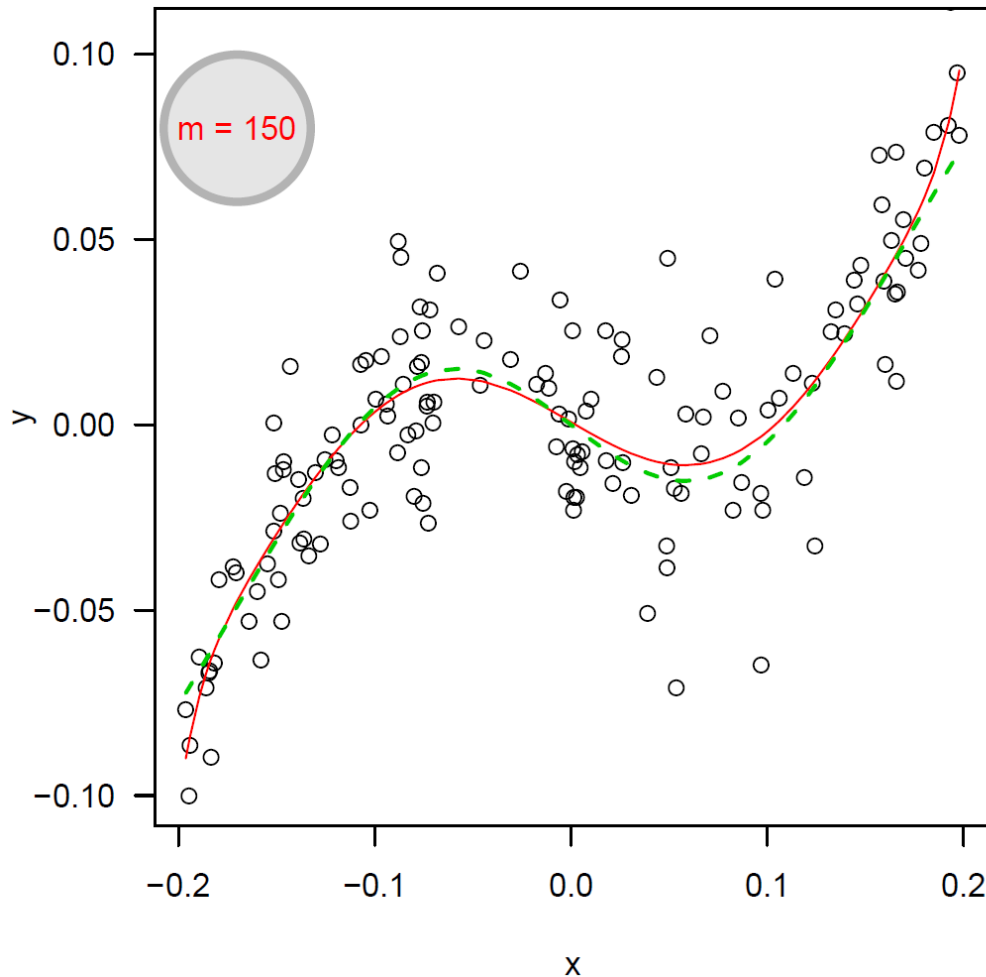
Iteratively improving our estimation

- Iteration counter: m
- Learning rate: v
- Total base-learners: n

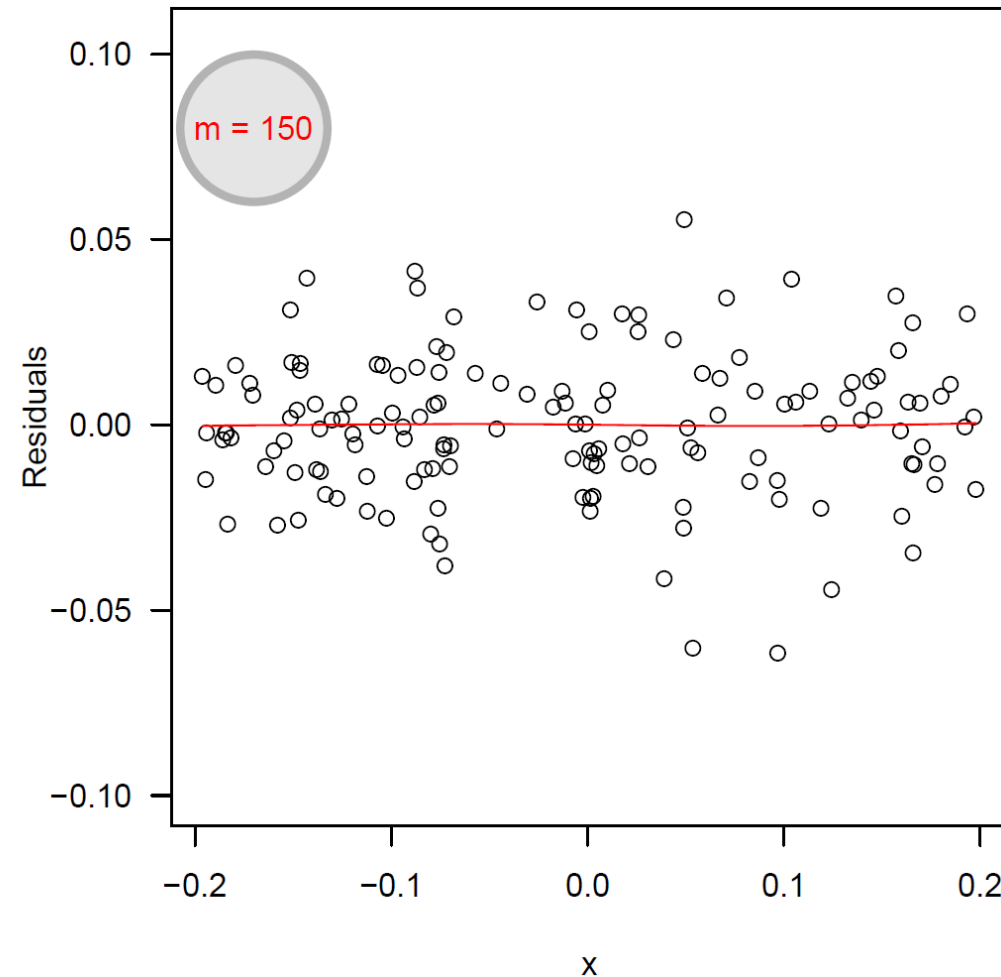
$$\hat{f}_{(1)}^{[m]} = \hat{f}_{(1)}^{[m-1]} + v \cdot -\frac{\partial}{\partial f_{(1)}} \left(\hat{f}_{(1)}^{[m-1]} \right)$$

Example iteration – function approximation

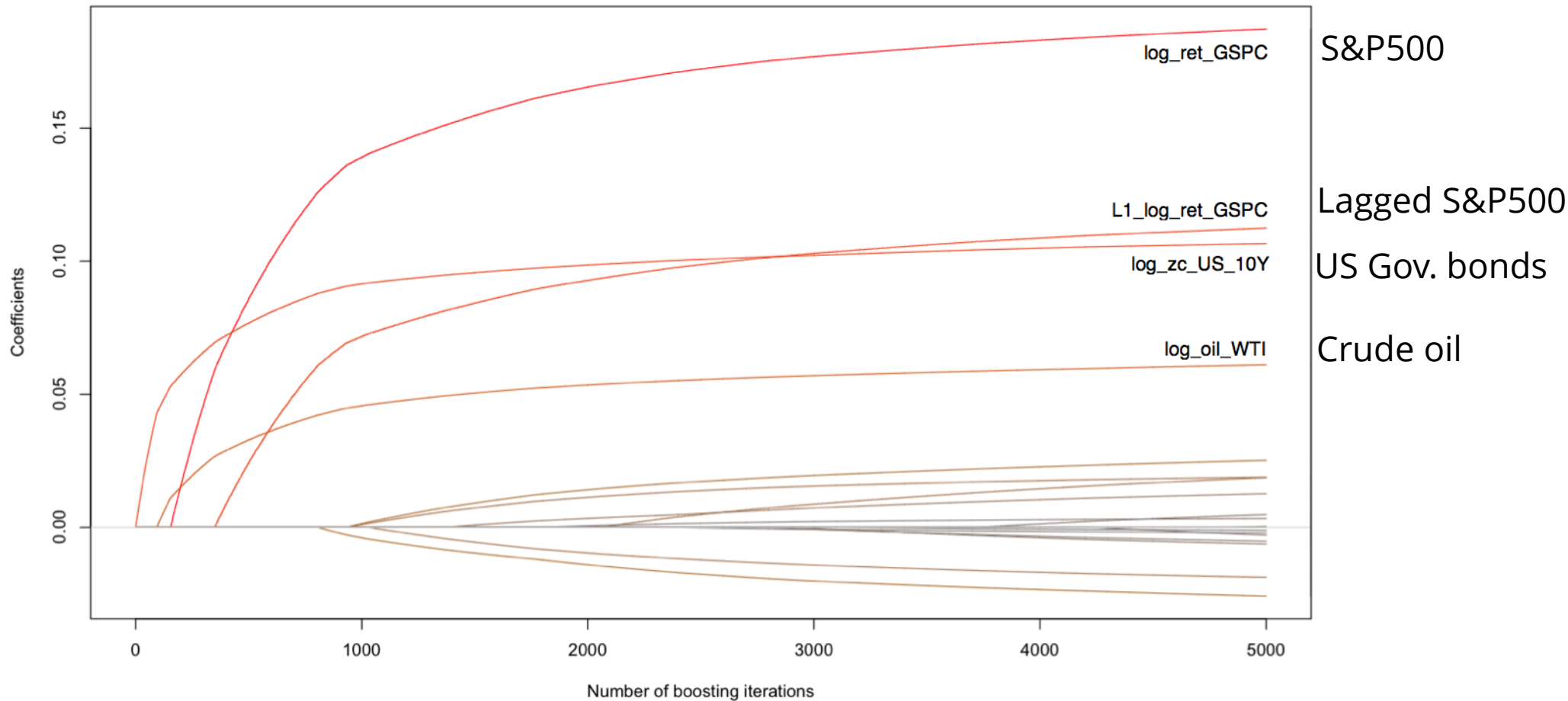
$$y = (0.5 - 0.9 e^{-50 x^2}) x + 0.02 \varepsilon$$



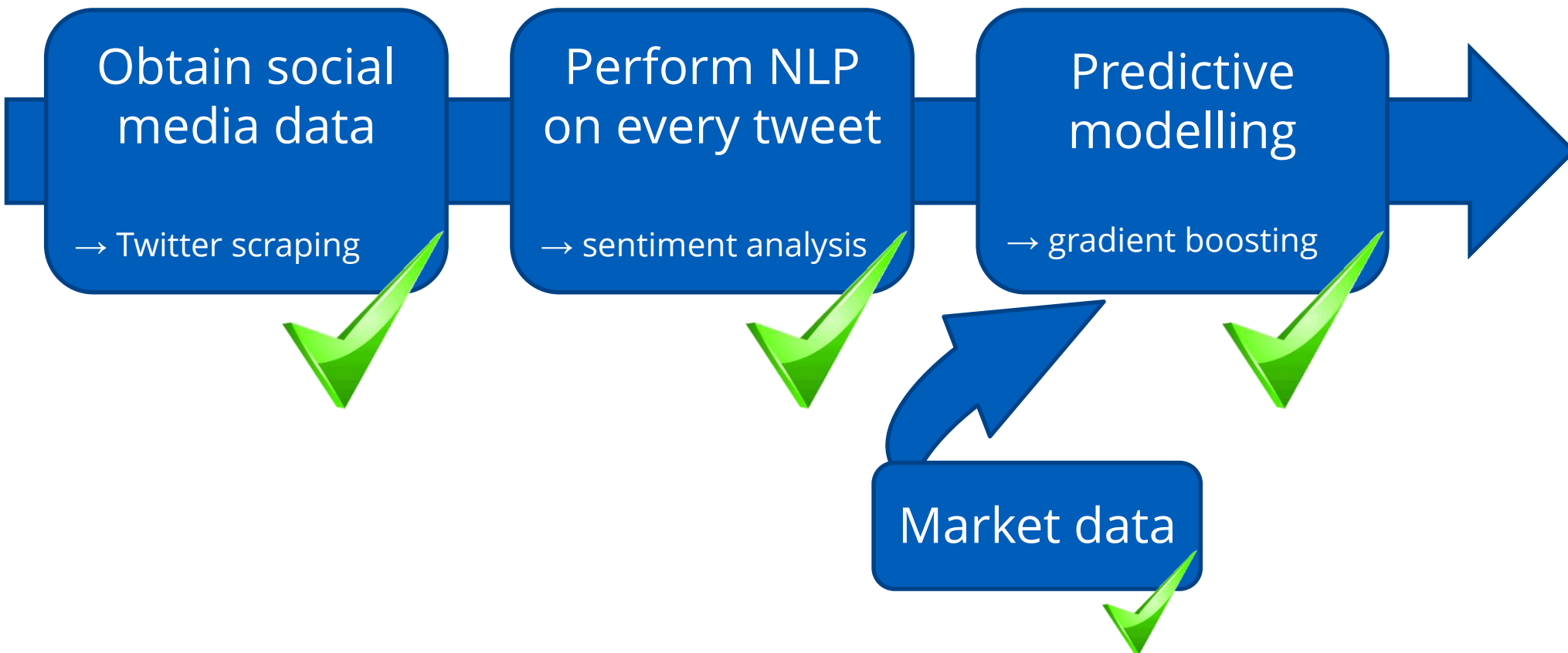
Residuals



Evolution of coefficients



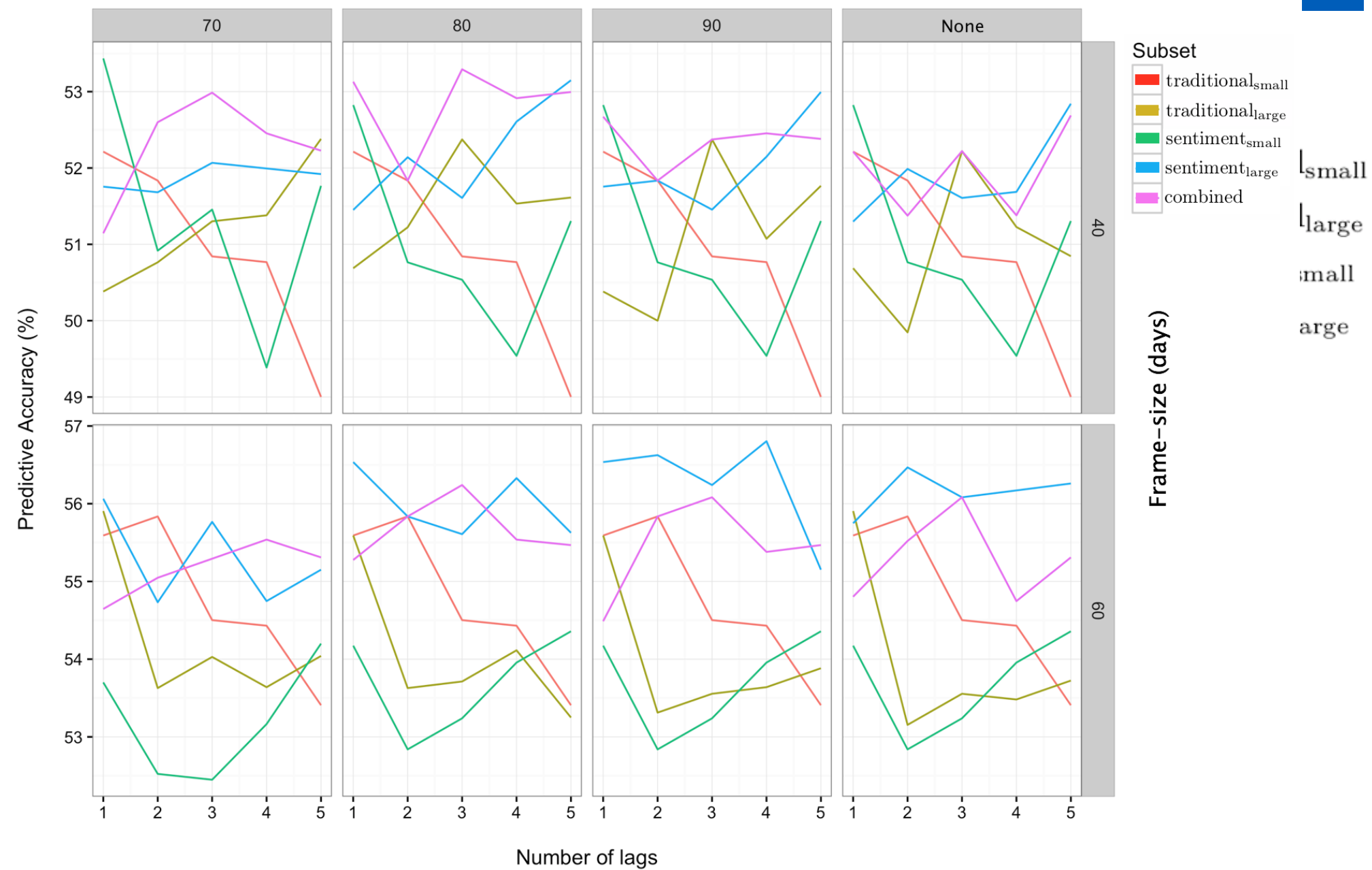
The Roadmap - master



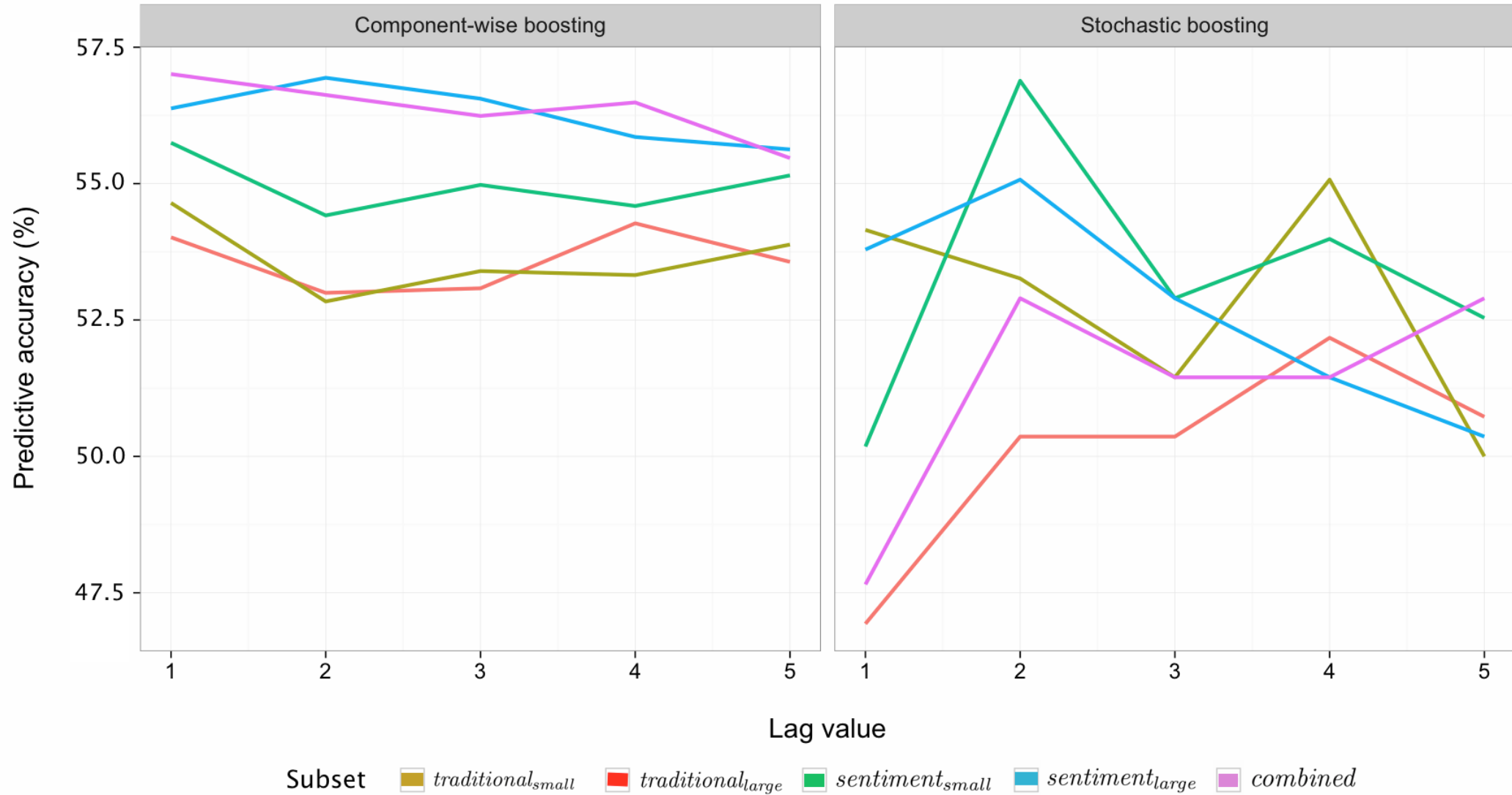
Reminder of the datasets

		Σ
• Traditional market data		
• small selection	→ traditional _{small}	6
• everything	→ traditional _{large}	36
• Sentiment data only		
• Aggregated	→ sentiment _{small}	22
• Individual	→ sentiment _{large}	100
• Market + social media data	→ combined	142

→ Correlation measures taken reduce datasets for final model

Correlation cutoff (κ)

Stochastic Gradient Boosting - comparison



The Roadmap

Questions?

Ideas for future work?

Obtain social media data

→ Twitter scraping

Perform NLP on every tweet

→ sentiment analysis

Predictive modelling

→ gradient boosting

Market data