

A decorative header featuring several overlapping, semi-transparent spheres in green, blue, red, and yellow at the top of the slide.

Large-Scale, Computationally Intensive Forecasting in R

Farzan Rohani, Eric Tassone,
and Murray Stokely



Overview



"Large-Scale" = Lots of time series

- Manual analysis impractical
- Diversity of time series
 - Growth
 - Seasonality
 - Shocks, regime change, emergence
- Automatic and robust methods needed
 - Ensemble method (a.k.a., "Many Models" approach)

"Computationally Intensive"...

- Two sources
 - Thousands of forecasts every day
 - Arising from statistical method
 - Fitting the "many models"
 - Quantifying statistical uncertainty
 - Simulation-based confidence intervals (requiring 1K to 10K forecasts per time series)

"Forecasting in R"

Just what it sounds like! :-)



Data



The Data

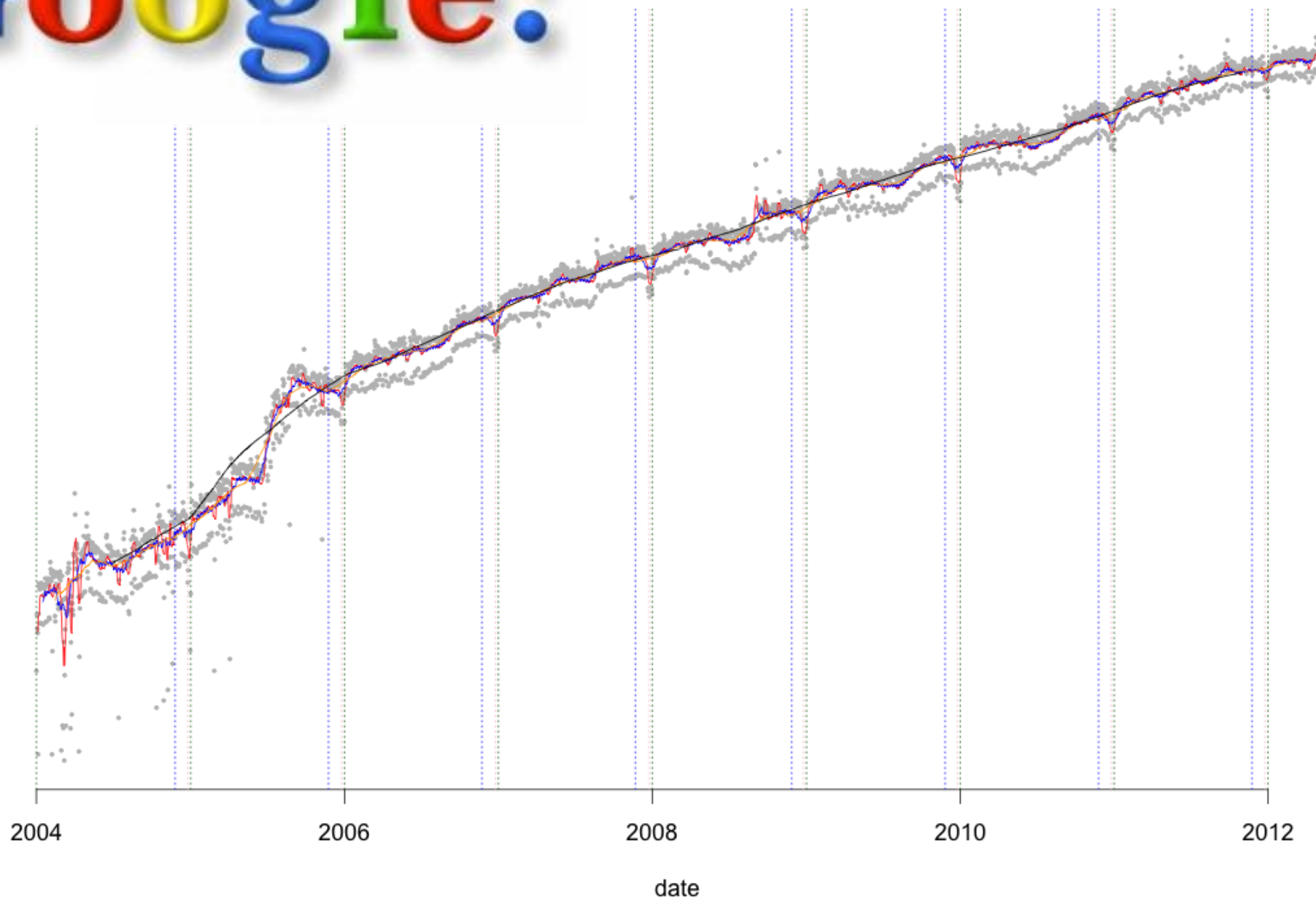
- 1000 time series from Google Trends:
 - Looked at past 5 years of data
 - Hence, persistently popular queries
 - Forecasts based on subset of 5 years
- Let's look at some of them...
 - Data since 2004, not just 5 years
 - Shown on logarithmic scale (y-axis)
 - Grey is data, 7-/31-/91-/365-day moving averages shown too
 - T-giving, New Years marked



Google!

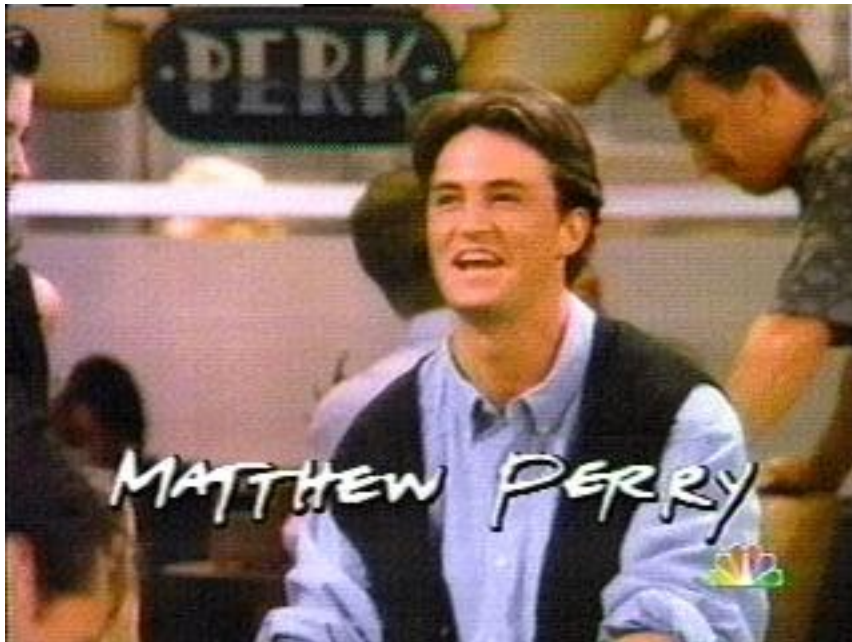
$\log(\text{google})$

value

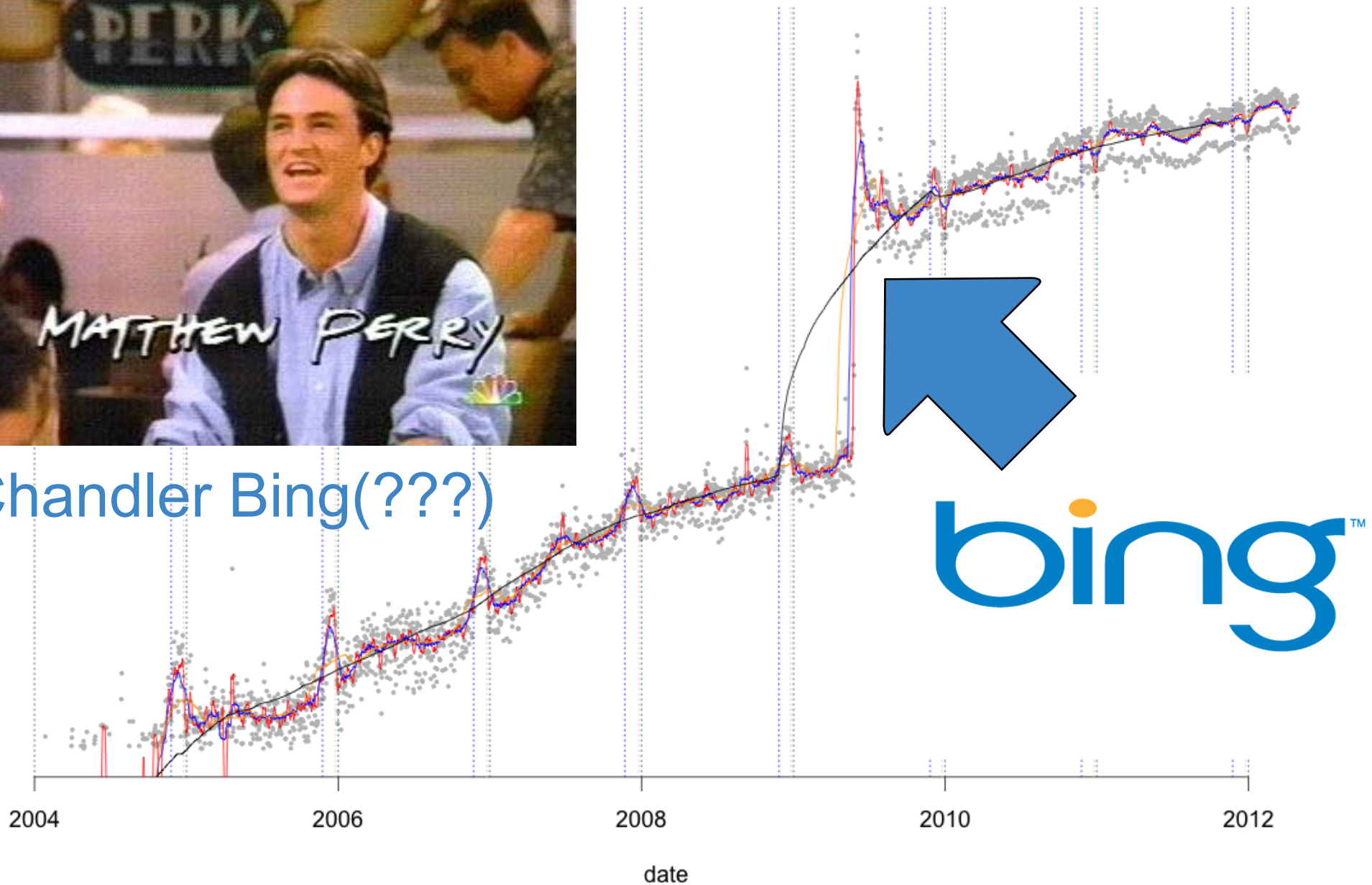


Google™

log(bing)



Chandler Bing(???)

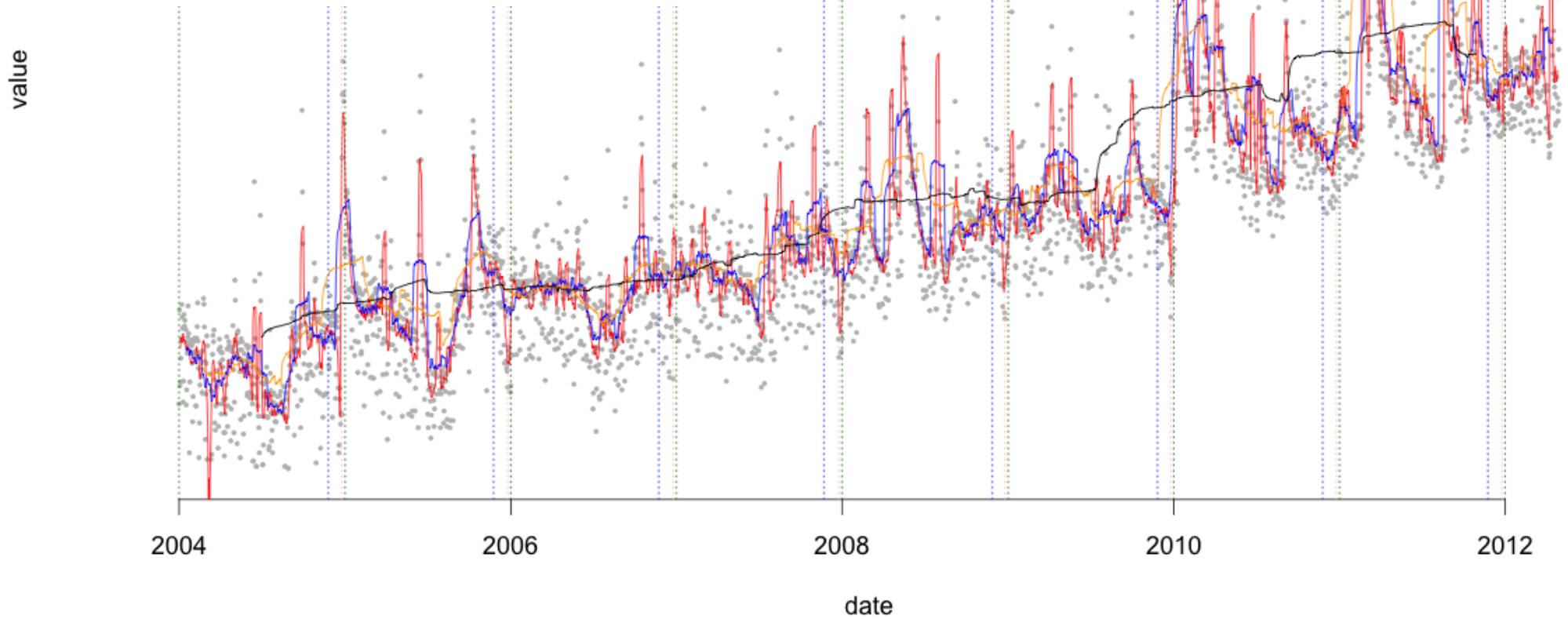
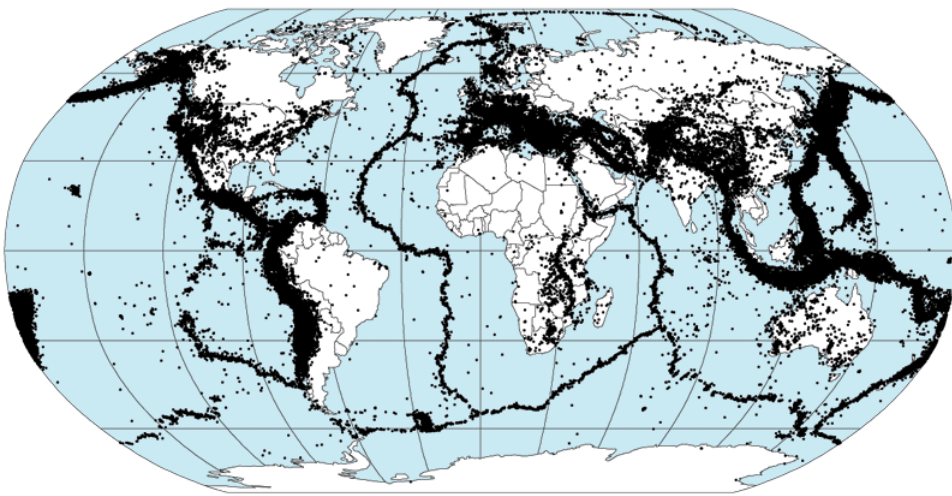


bing™

Google™

Preliminary Determination of Epicenters
358,214 Events, 1963 - 1998

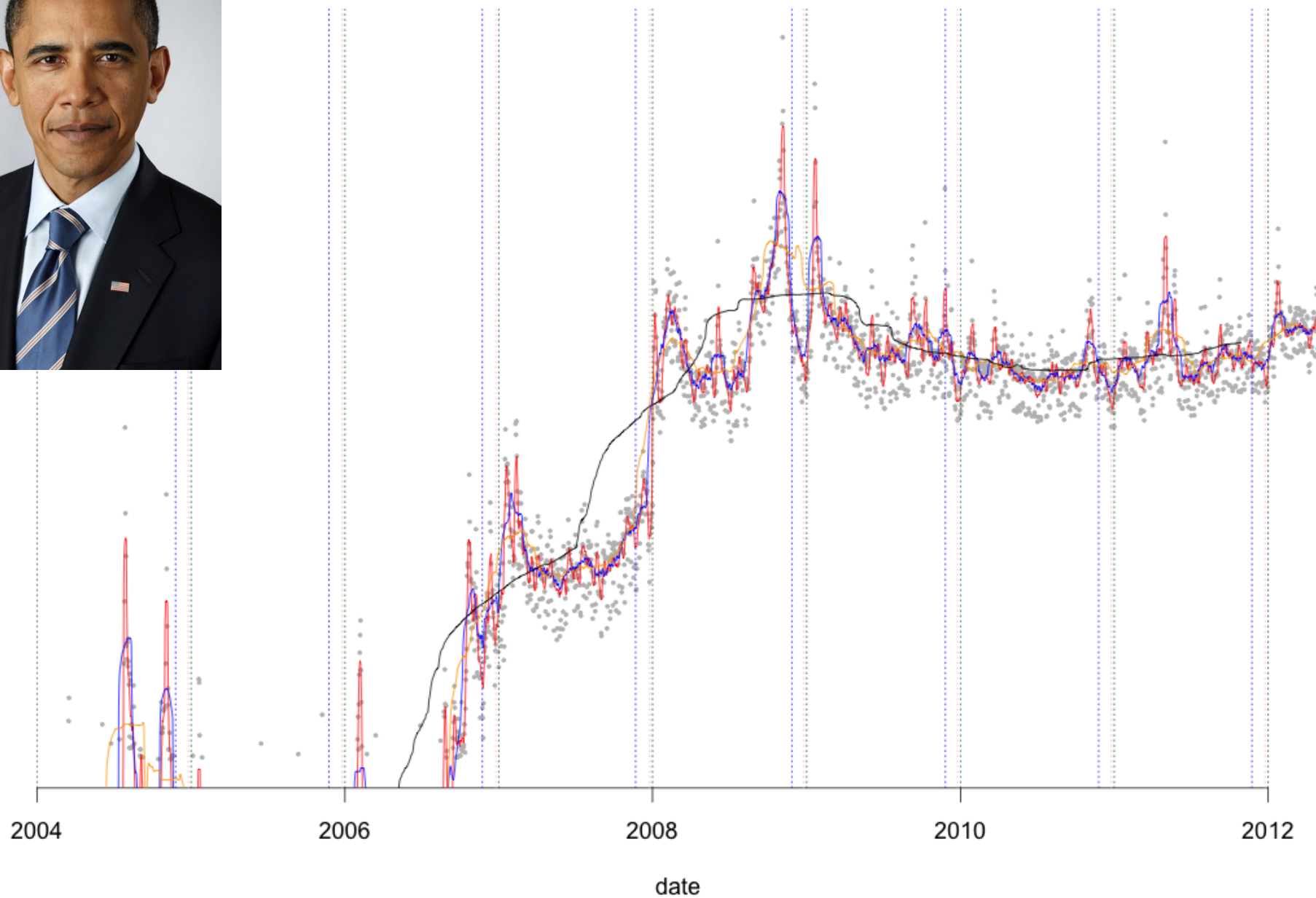
$\log(\text{earthquake})$





log(obama)

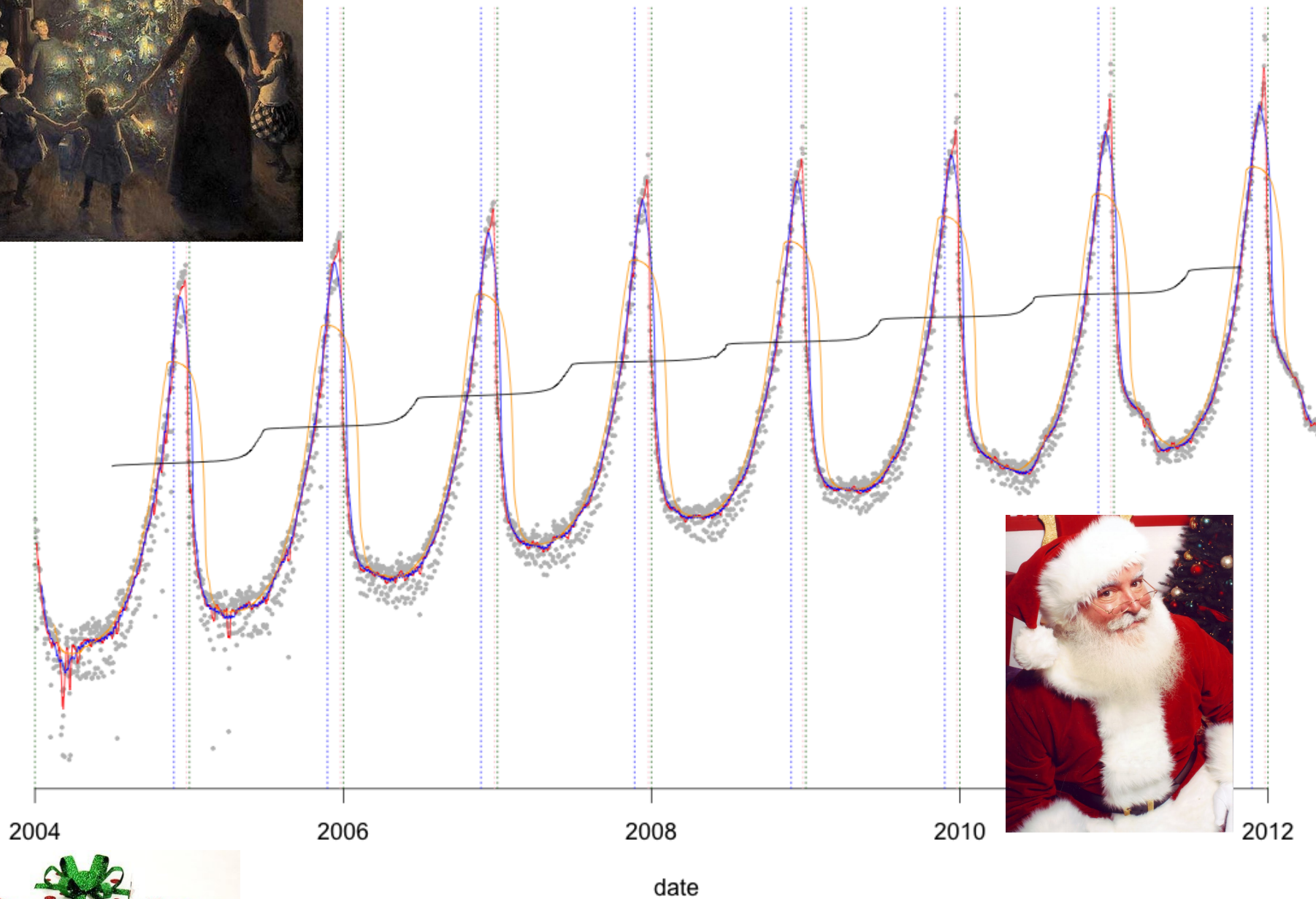
value





log(christmas)

value

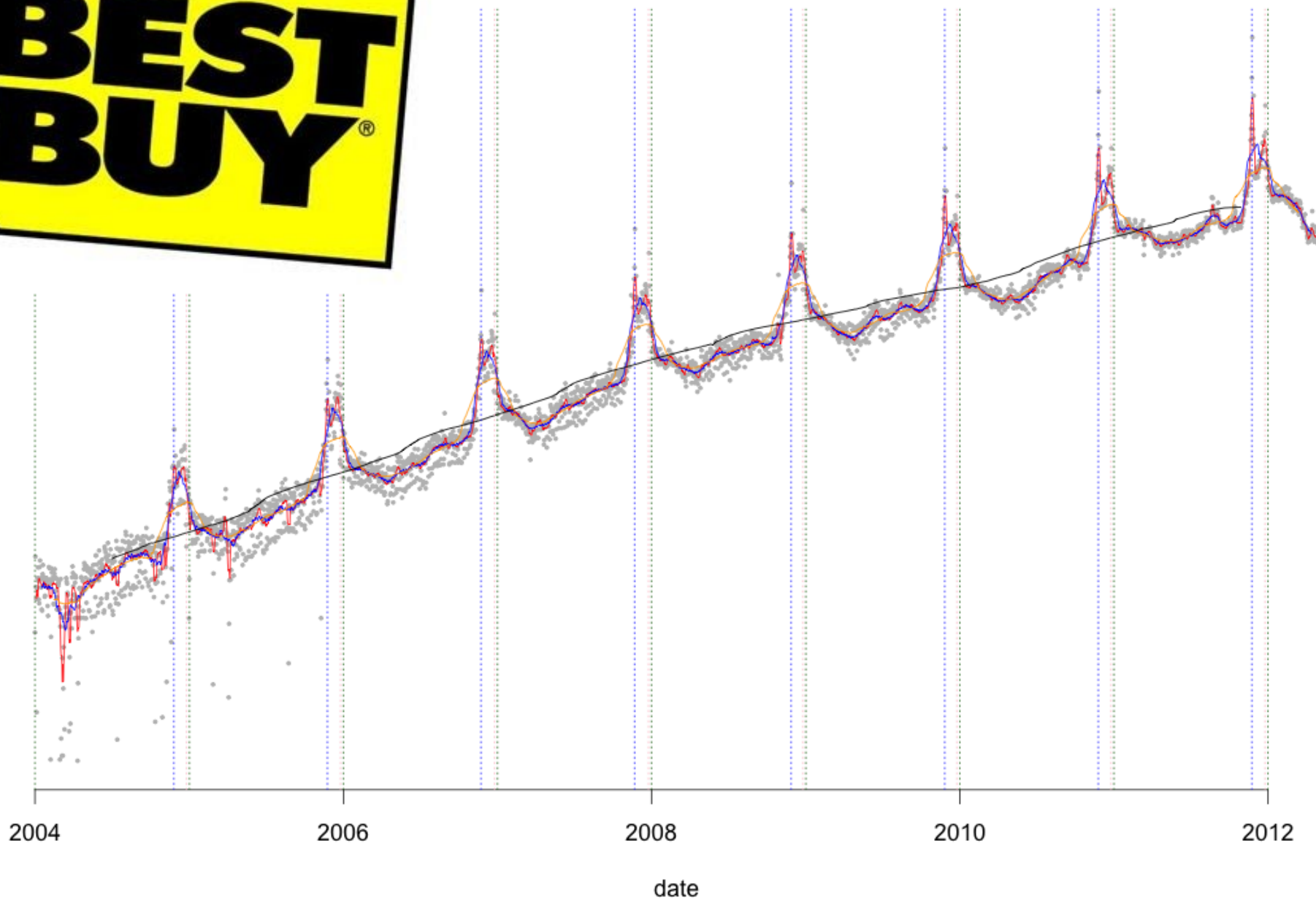


Google™



log(best buy)

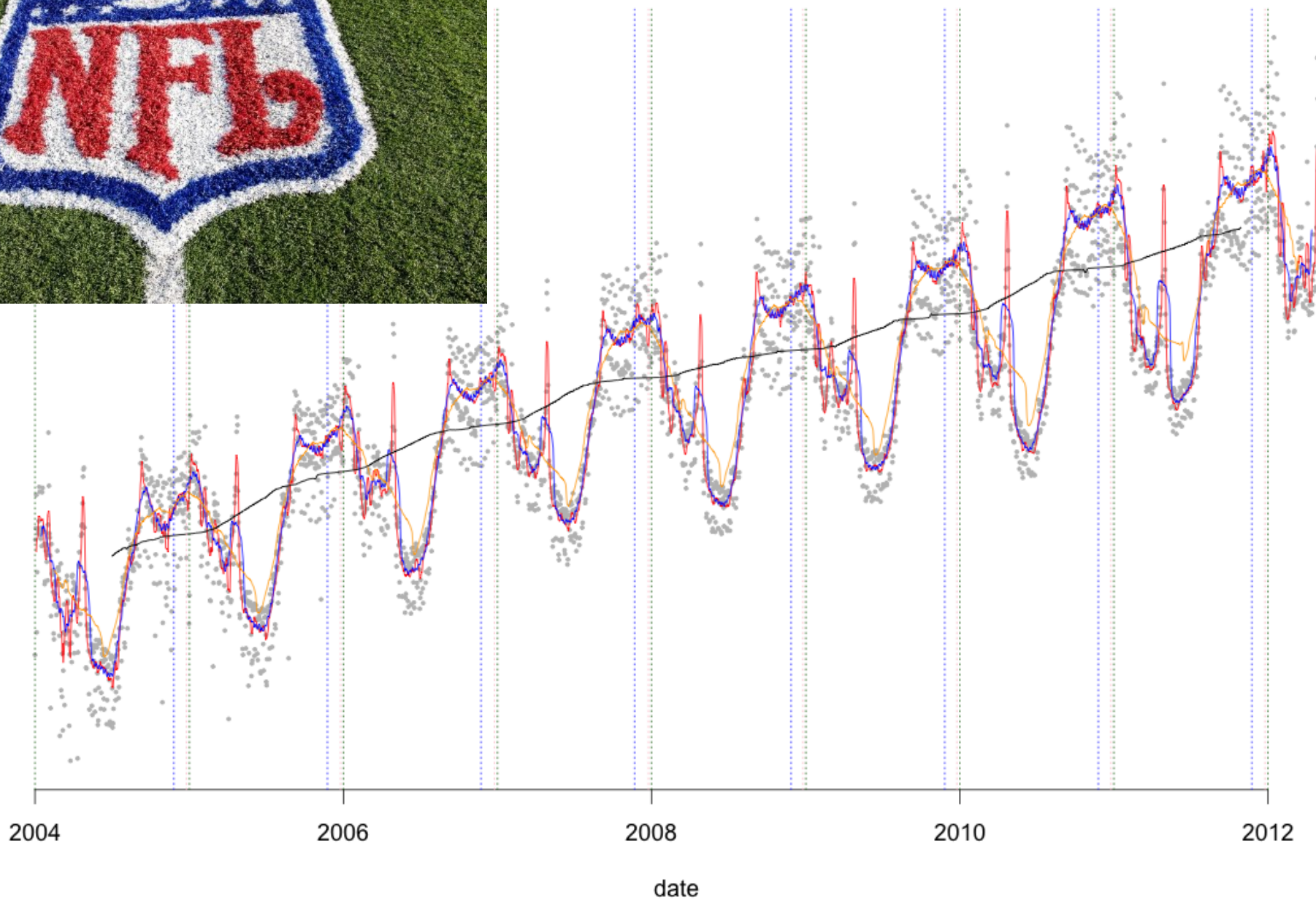
value





$\log(\text{nfl})$

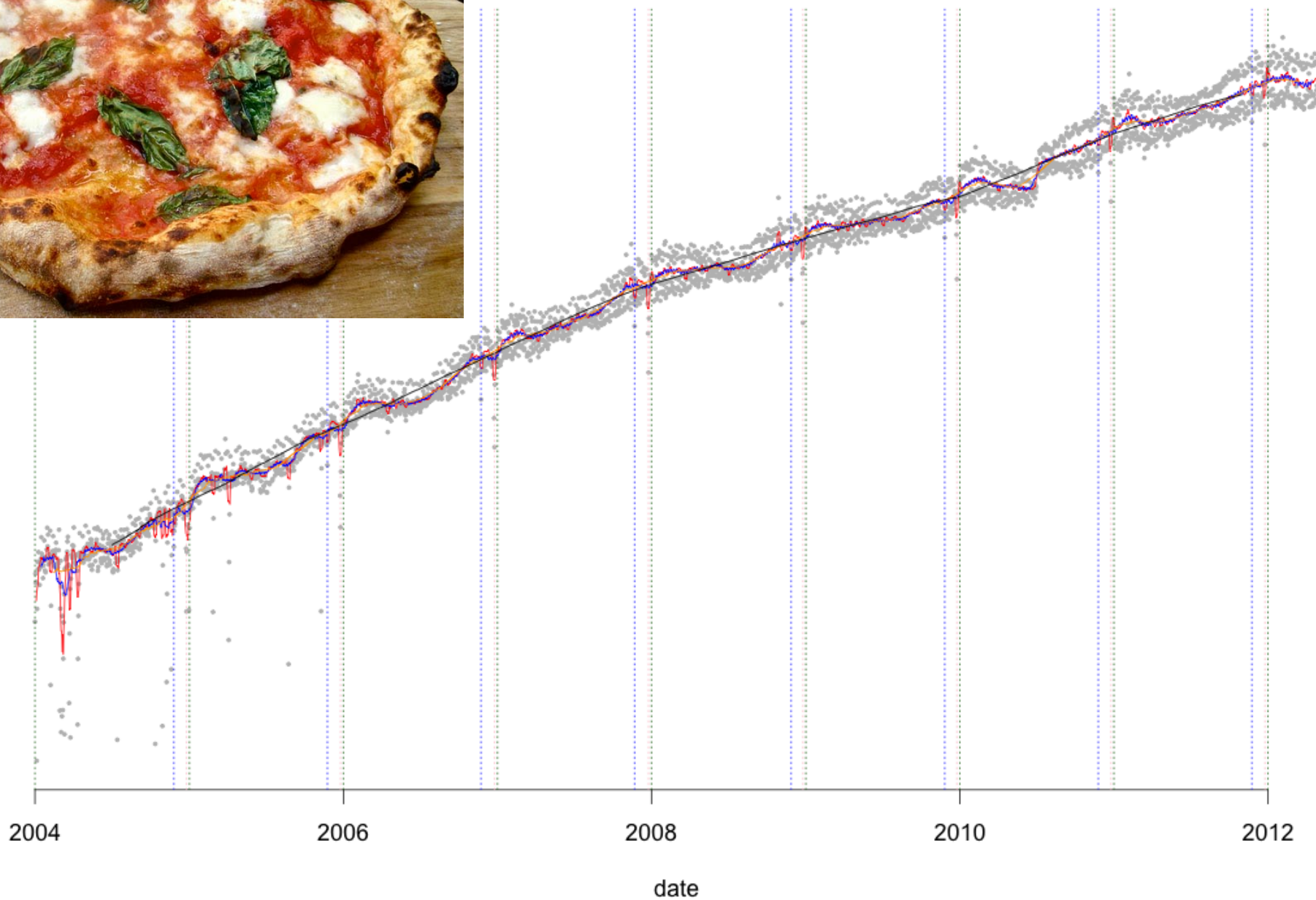
value





$\log(\text{pizza})$

value

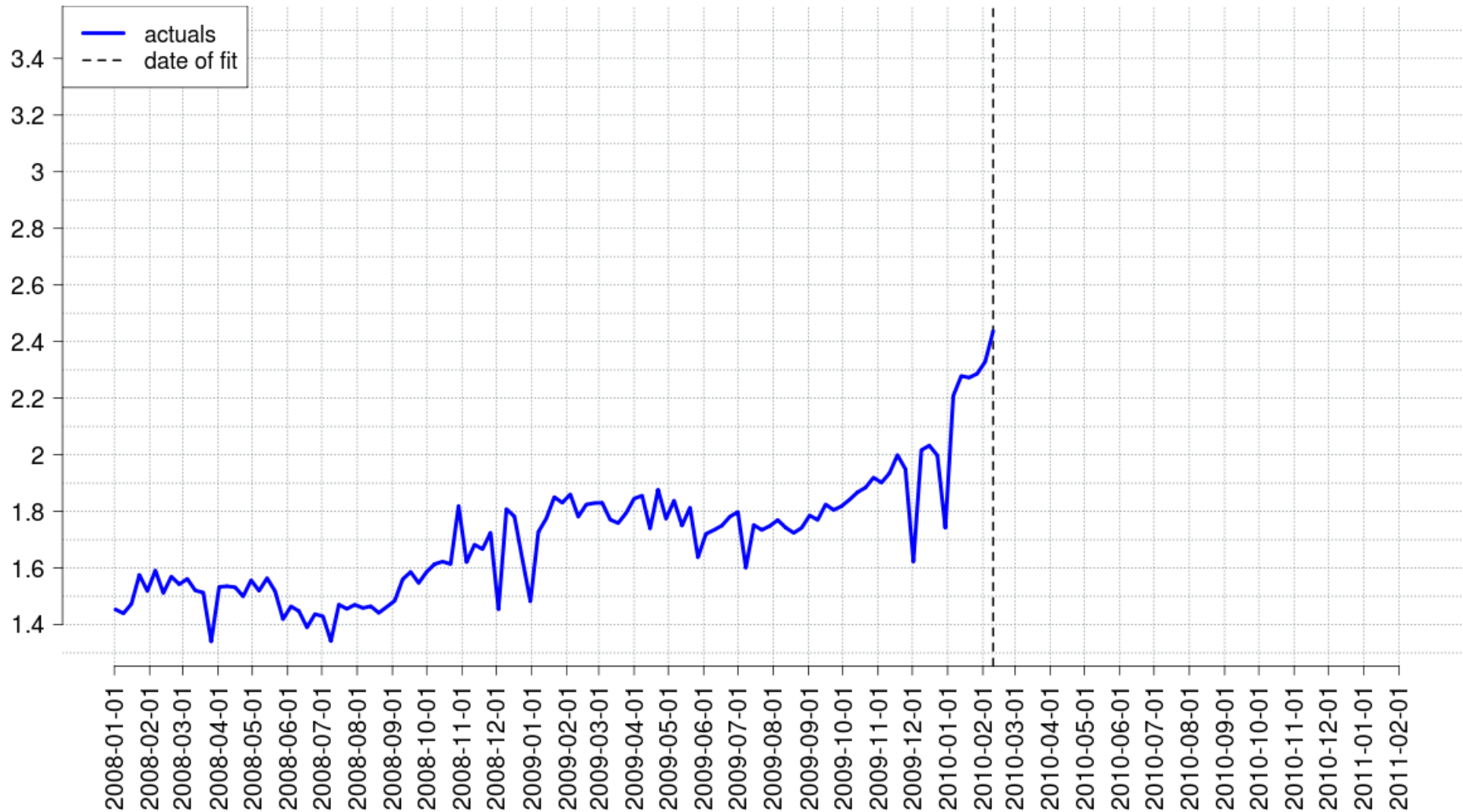


Methods



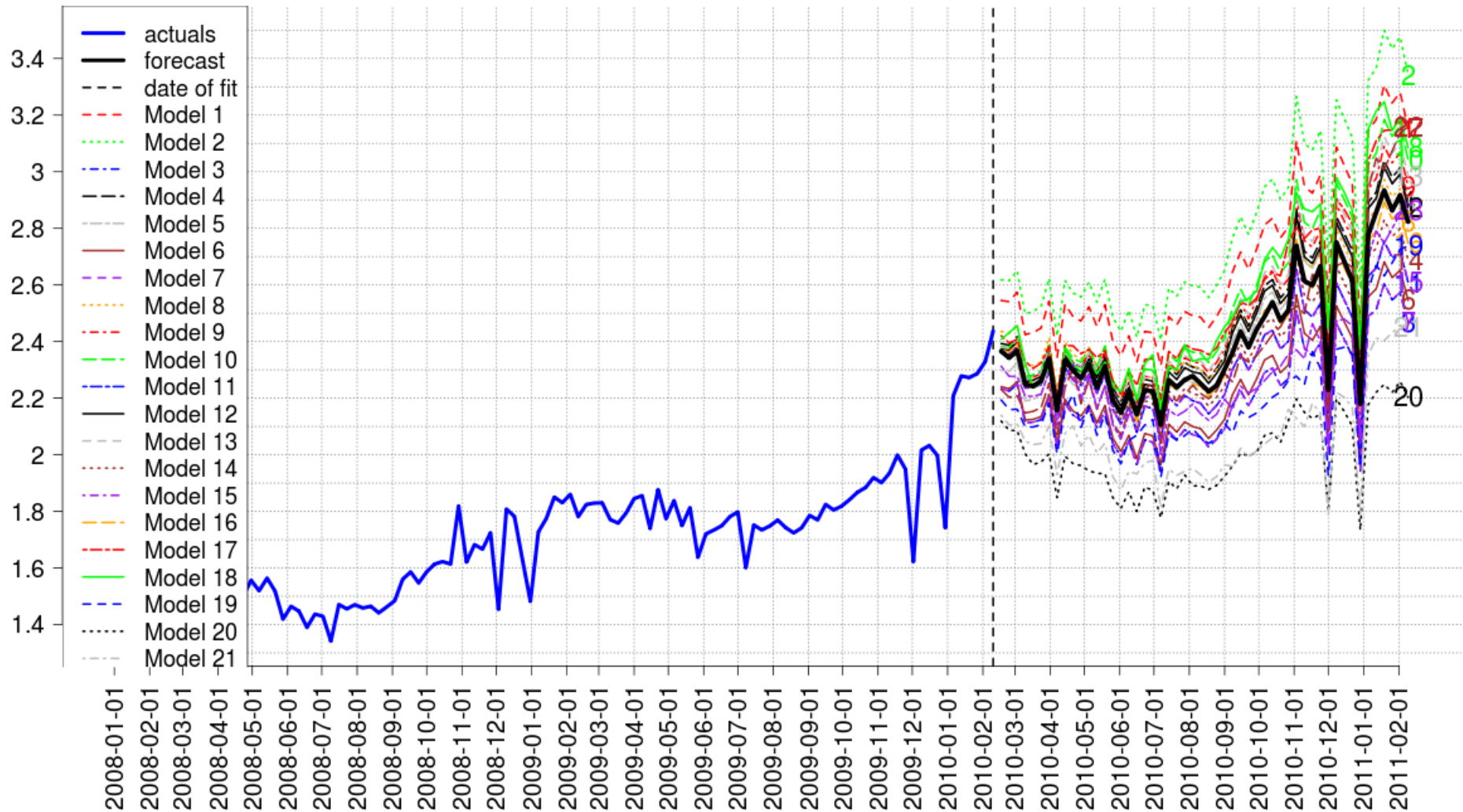
Many Models Approach

The trend for query 'pizza' in US from Google Trends



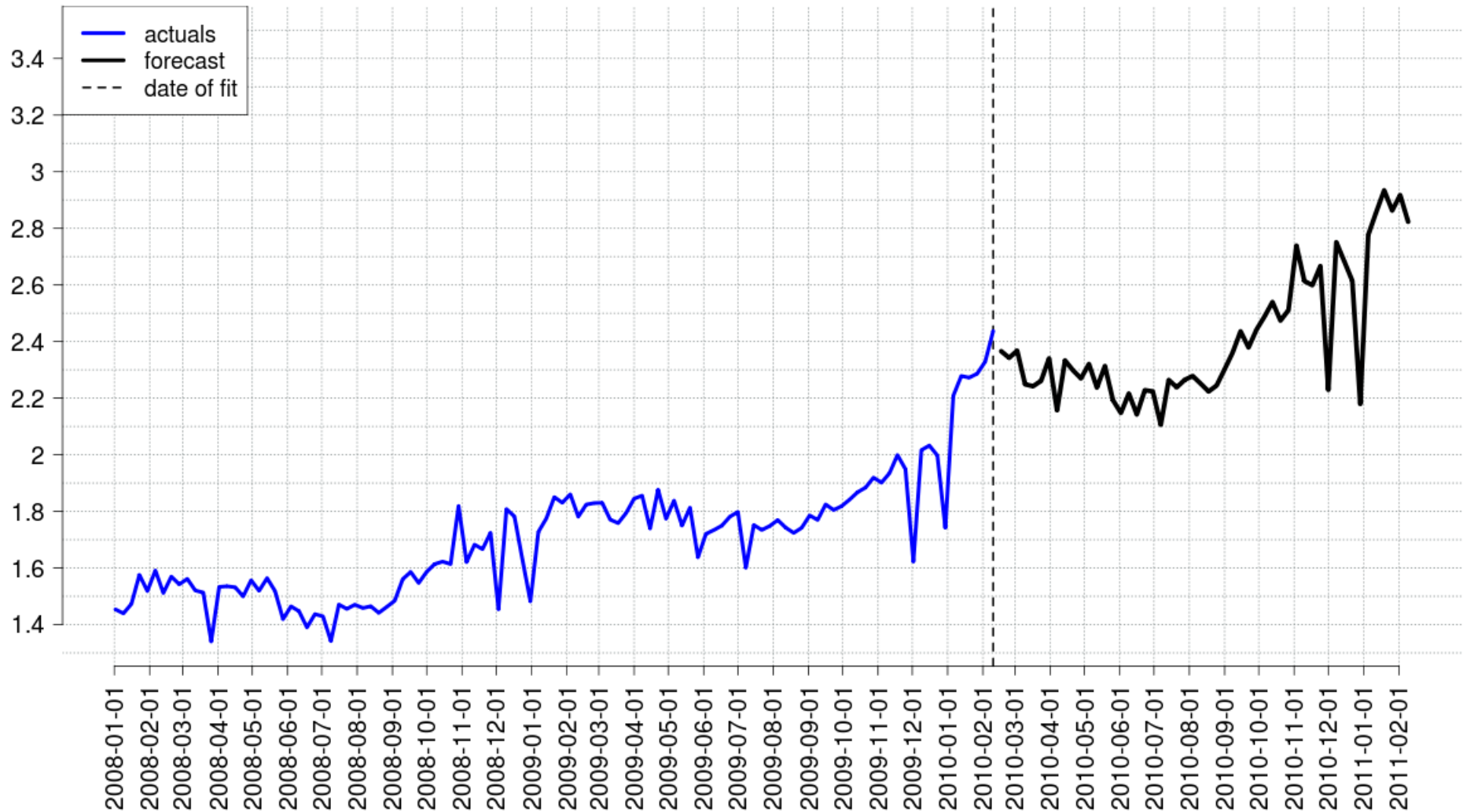
Many Models Approach

The trend for query 'pizza' in US from Google Trends



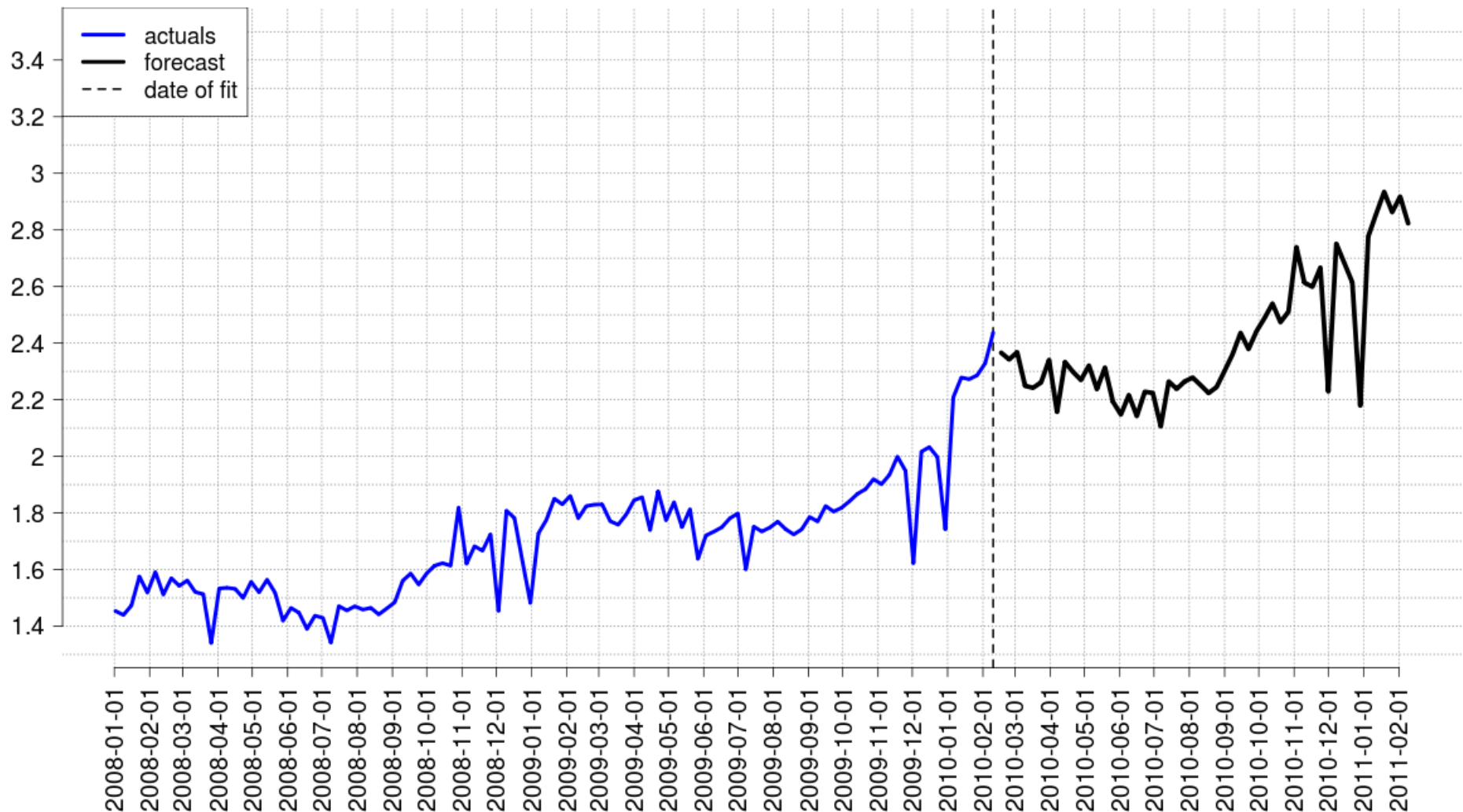
Many Models Approach

The trend for query 'pizza' in US from Google Trends



Many Models Approach

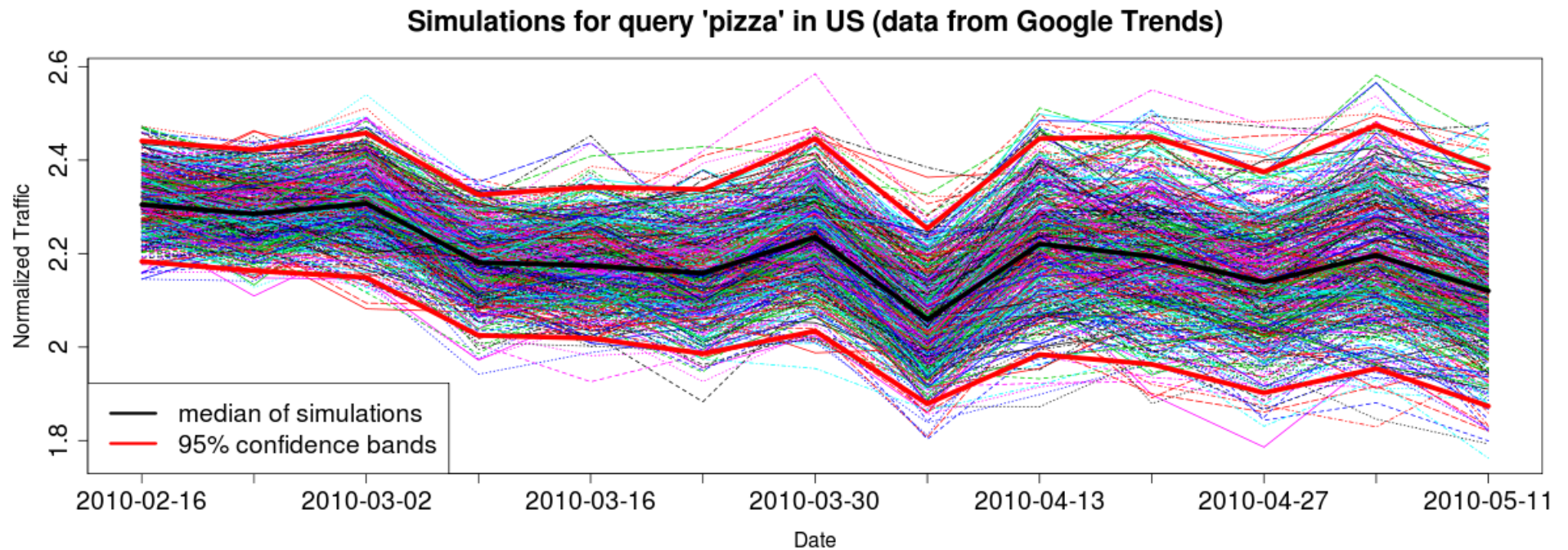
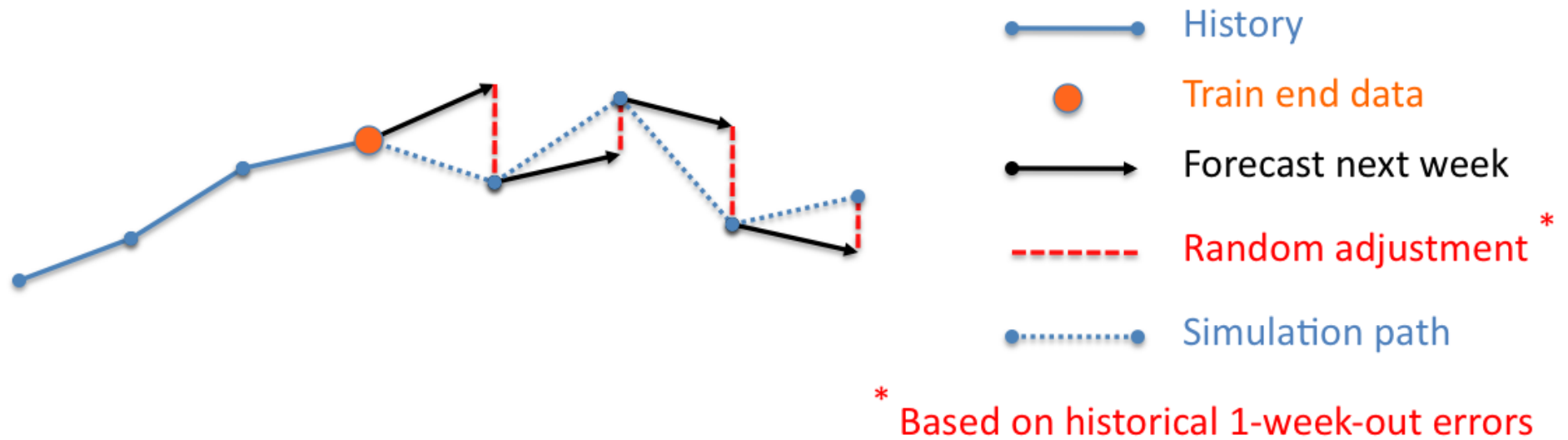
The trend for query 'pizza' in US from Google Trends



Armstrong, J. S. **Combining forecasts.** *Principles of forecasting: A handbook for researchers and practitioners*, 417–439.

Clemen, R. **Combining forecasts: A review and annotated bibliography.** *International Journal of forecasting* 5, 559-583

Confidence Intervals



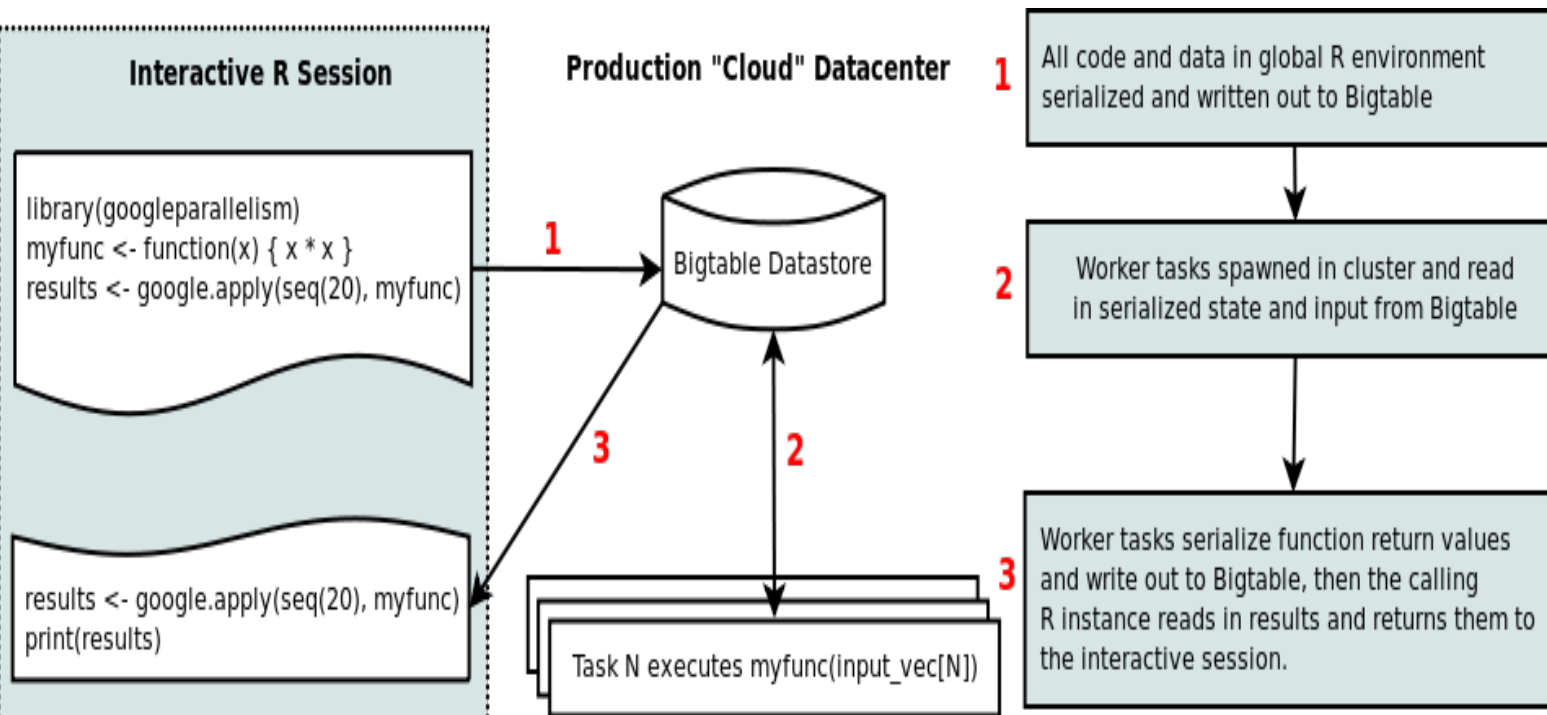
Parallelization in R: Why?

1,000 Trends t.s. X

1,000 realization/t.s. X

~10 sec/realization = 10^7 sec ~ 4 months

But 9 hours in parallel, cut by factor of ~300

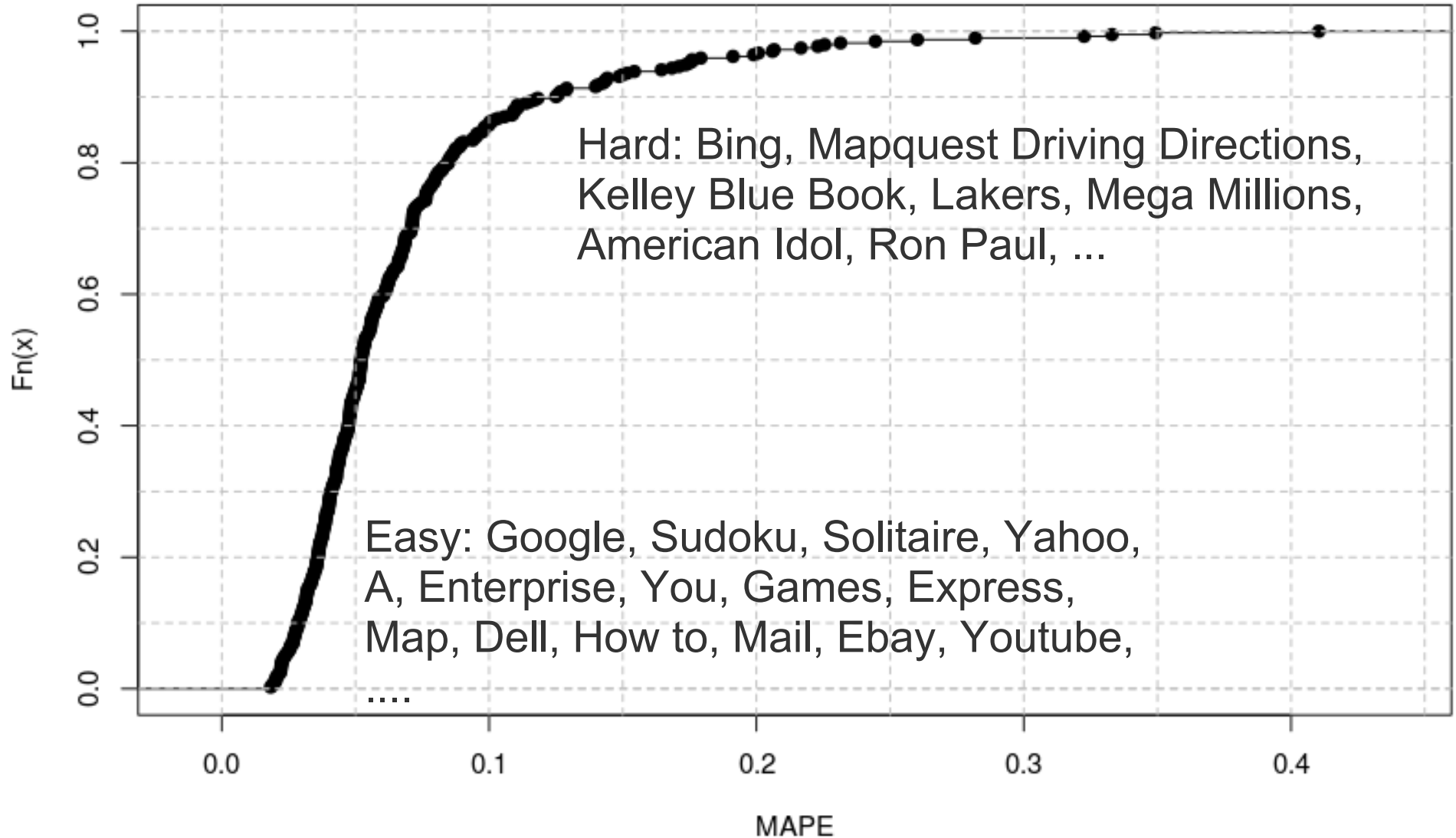


Overall Performance



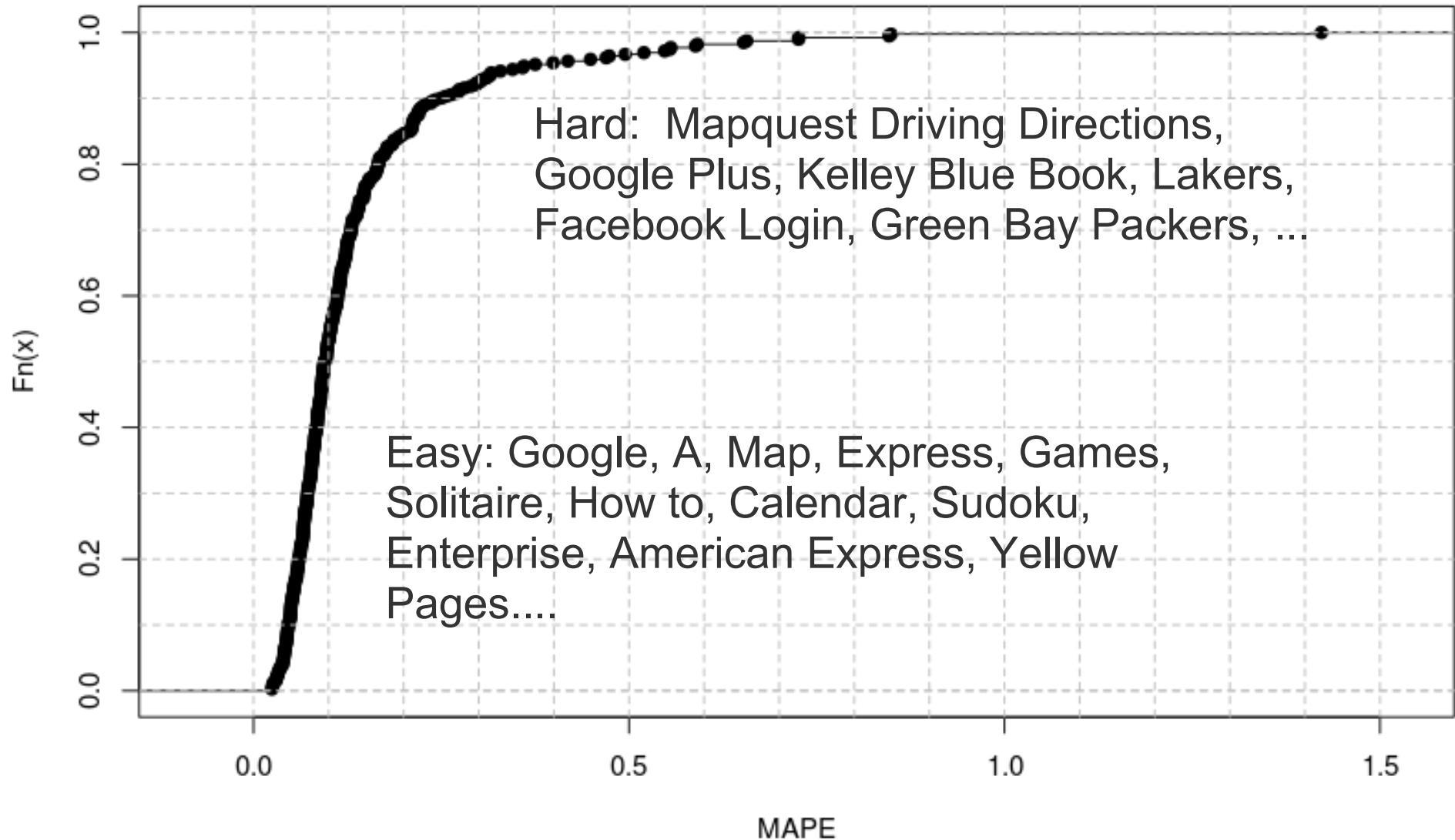
One-Week-Out MAPE

One-Week-Out MAPE Distribution

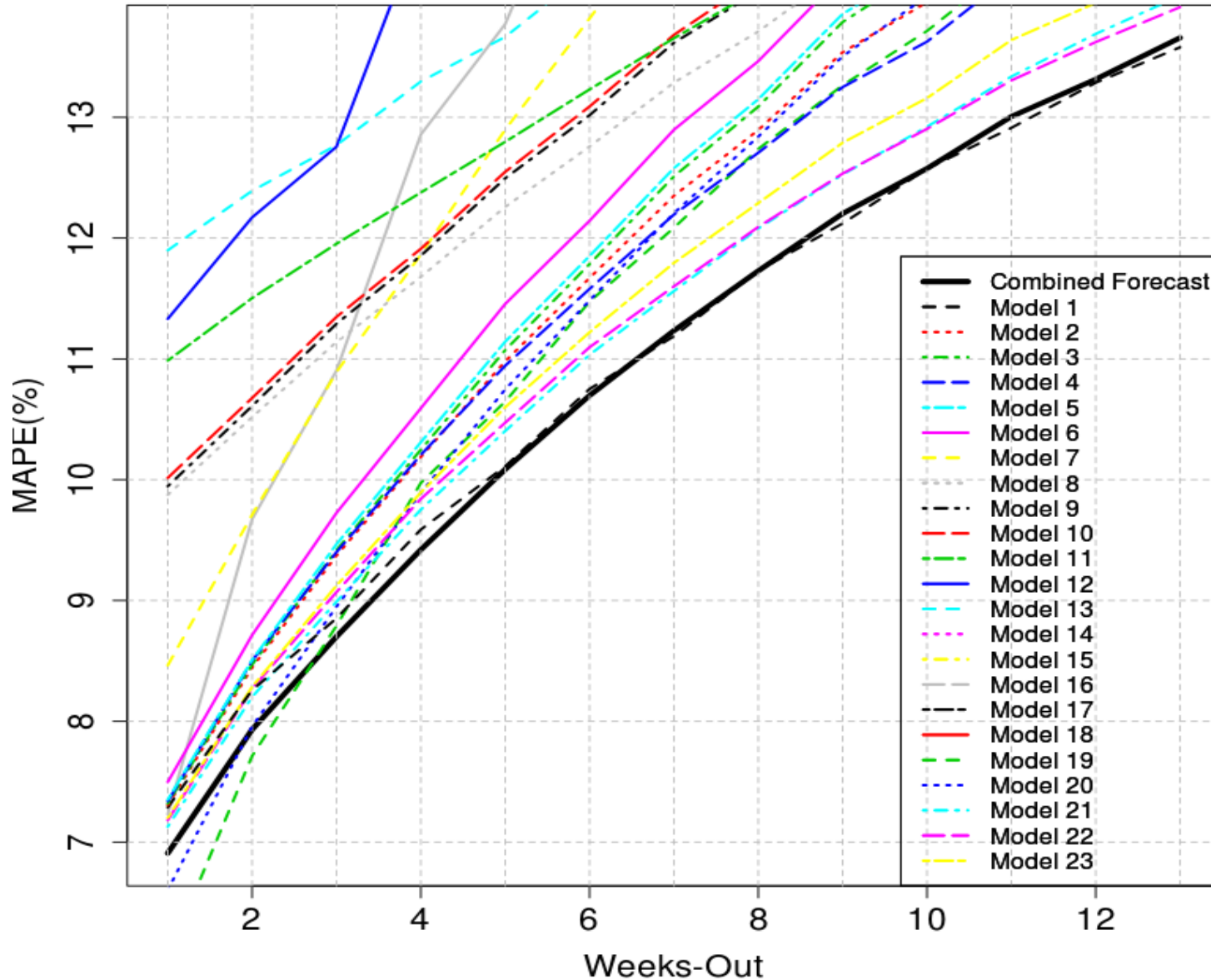


One-Quarter-Out MAPE

One-Quarter-Out MAPE Distribution



Performance of Individual Models (Overall, based on all queries)



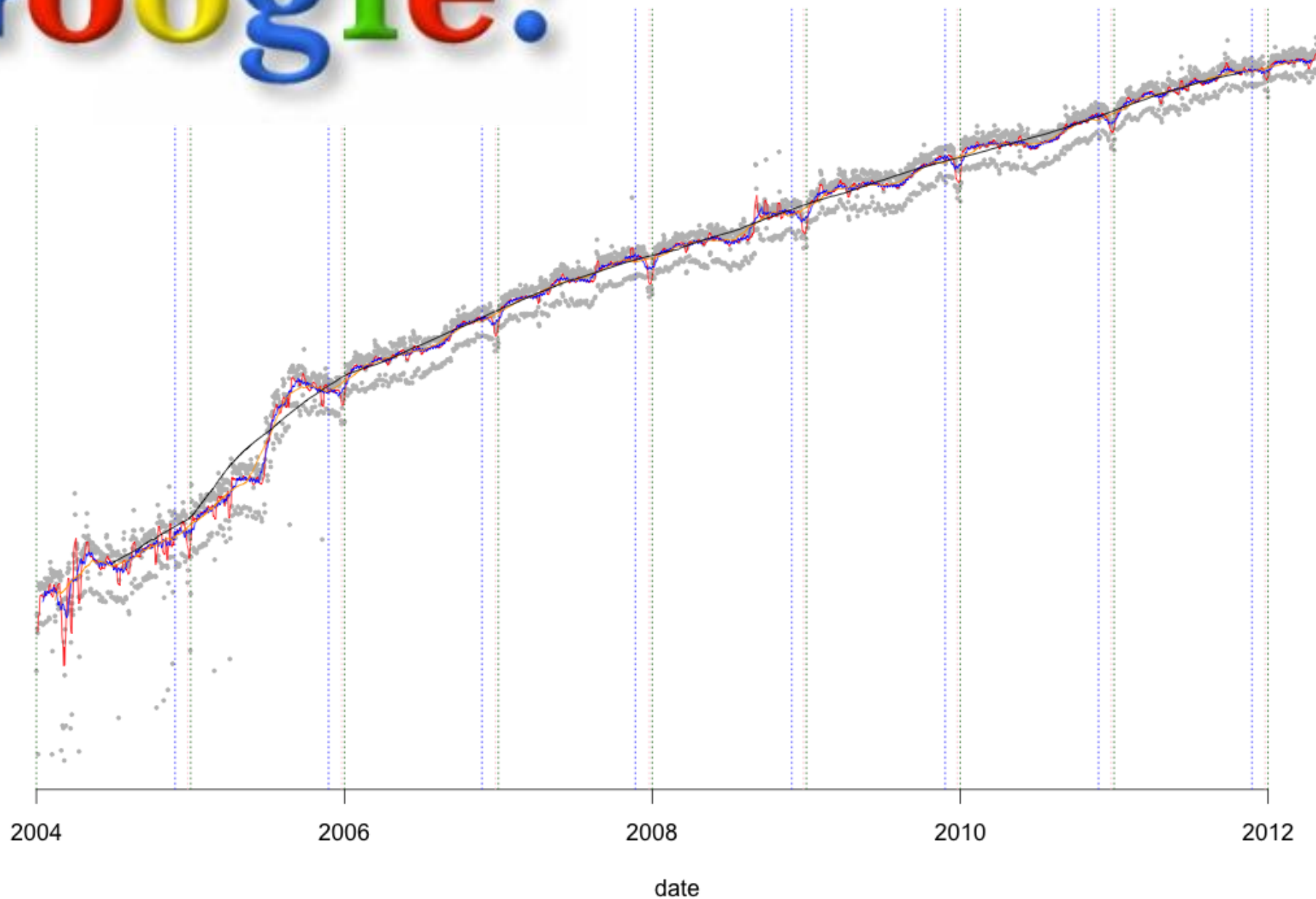
Performance on Individual Time Series



Google!

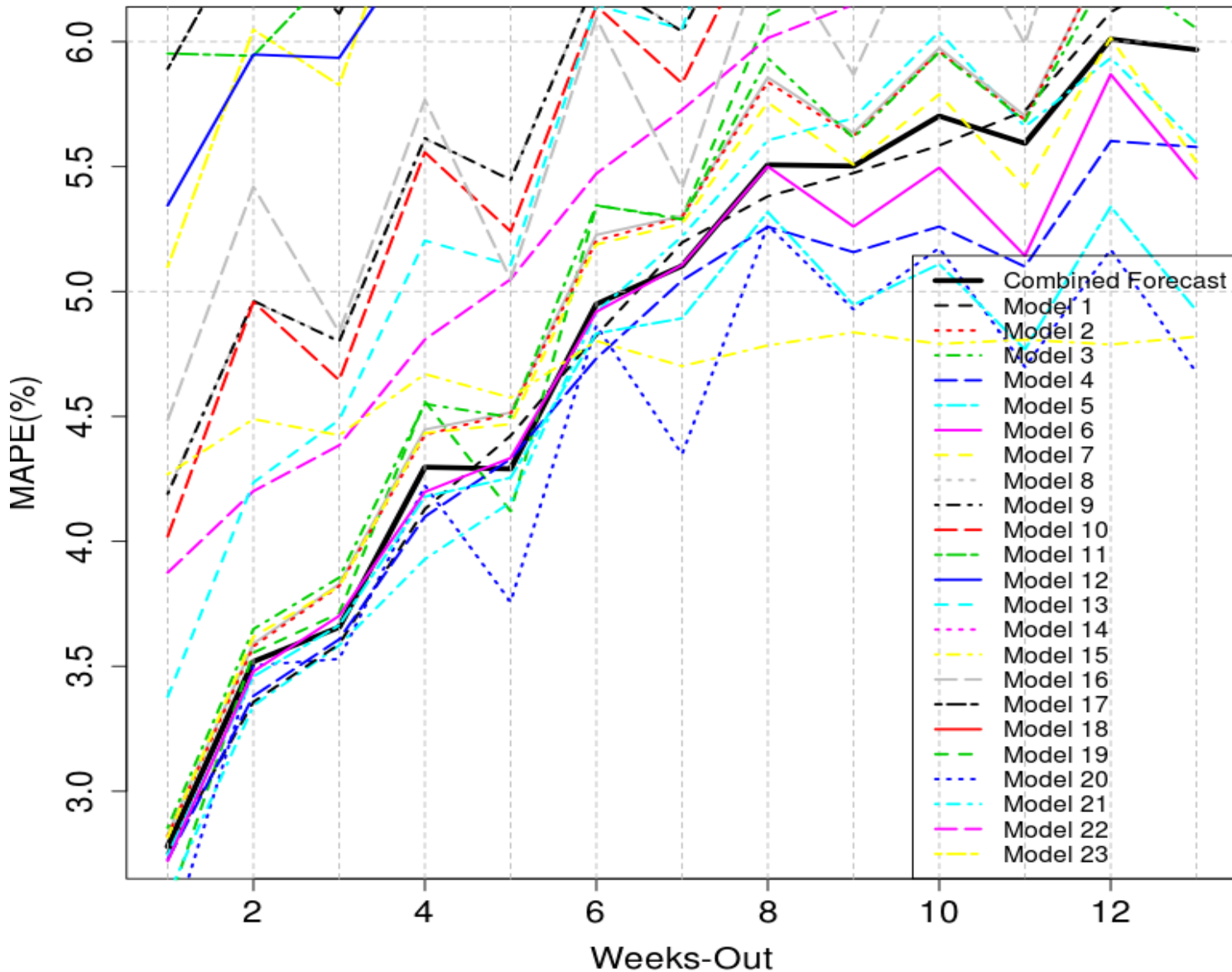
$\log(\text{google})$

value

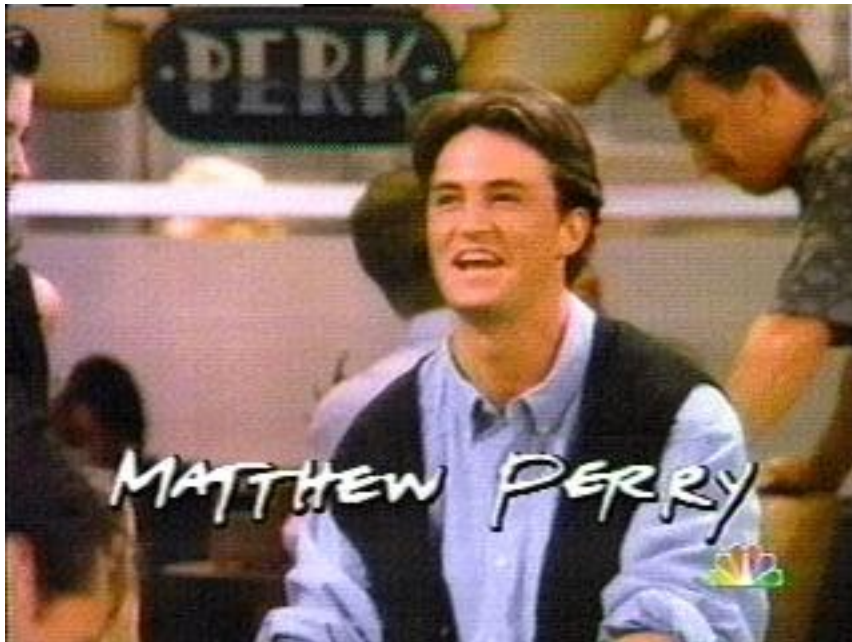


Google™

Performance on query "google"

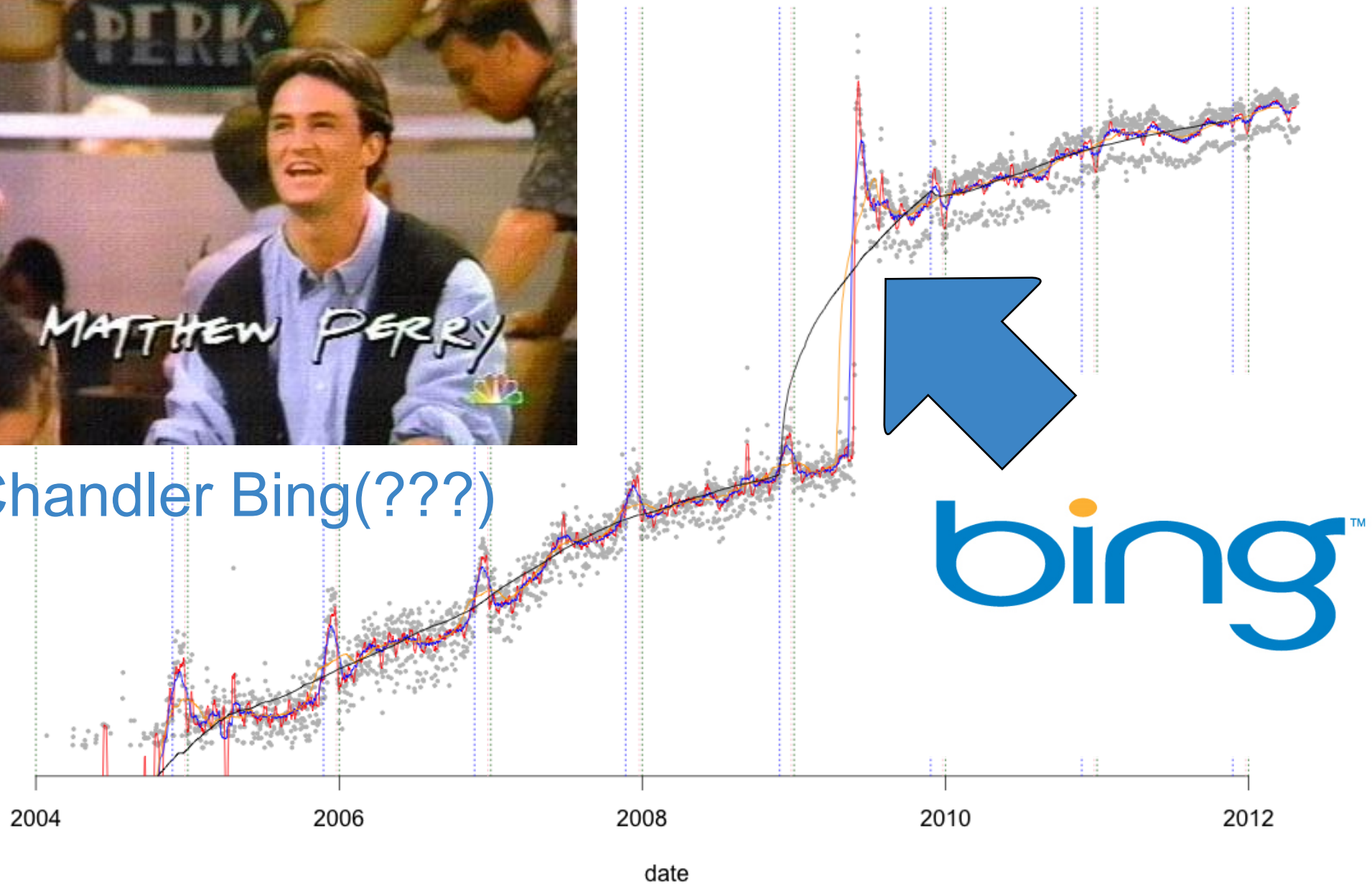


log(bing)



Chandler Bing(???)

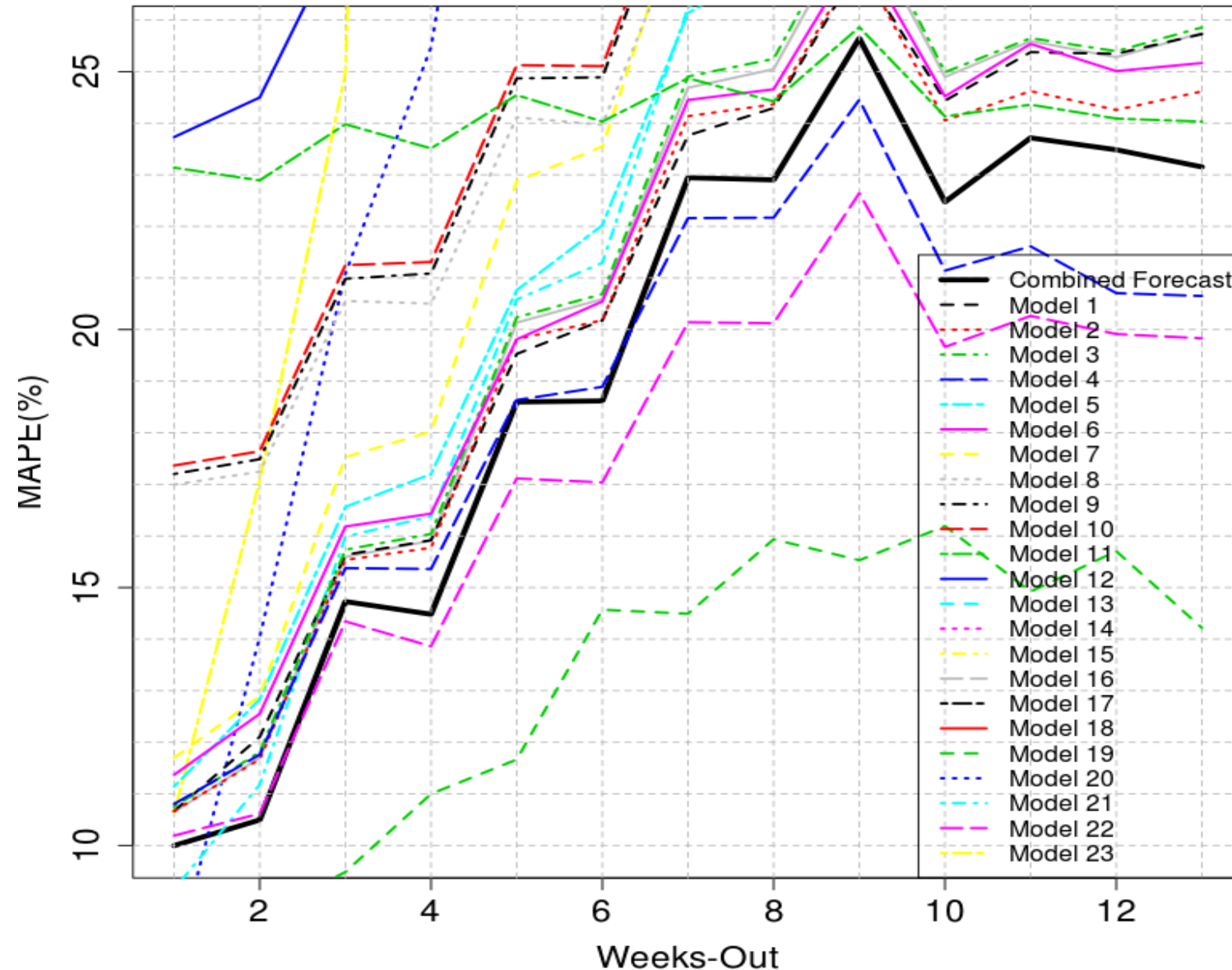
value



bing™

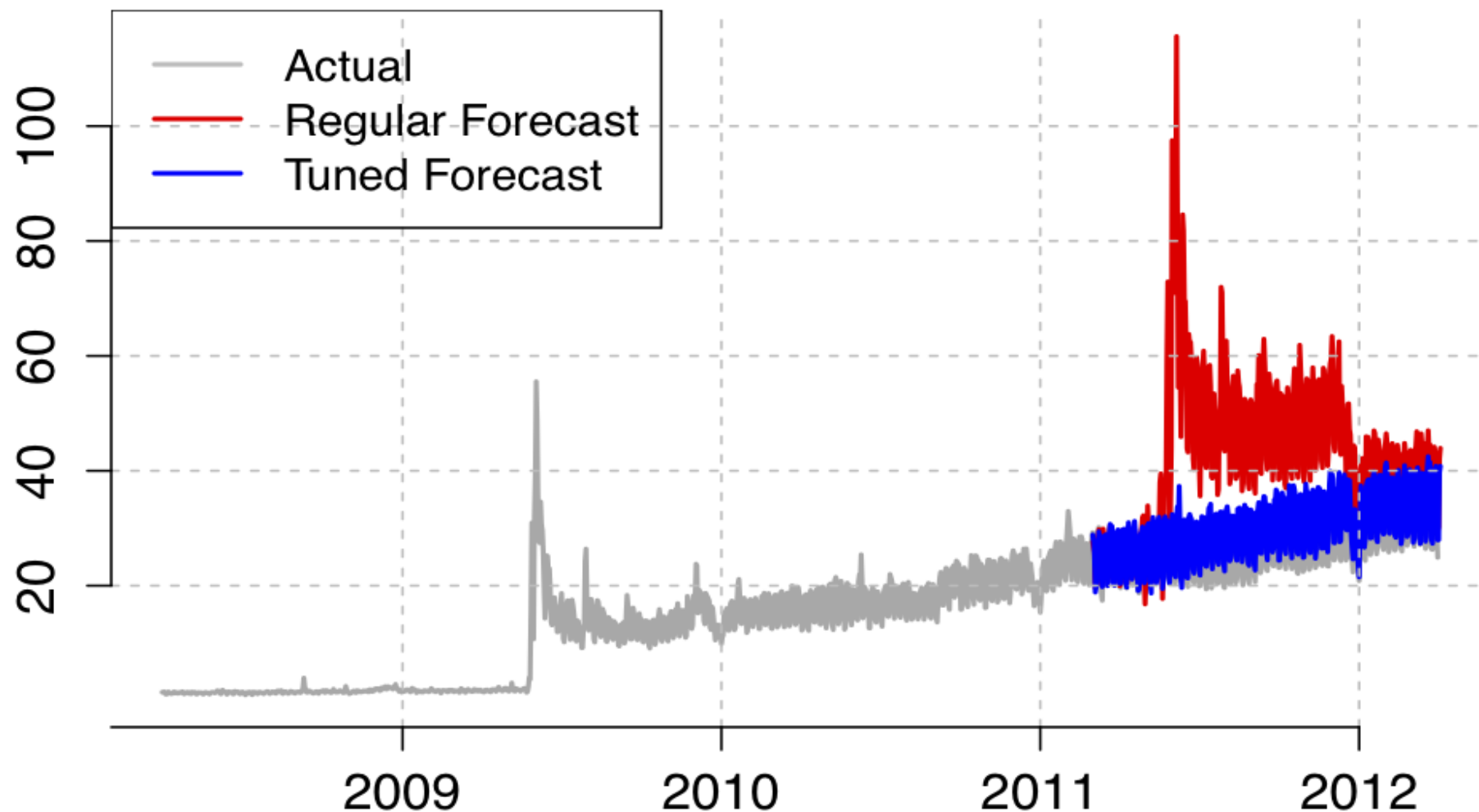
Google™

Performance on query "bing"



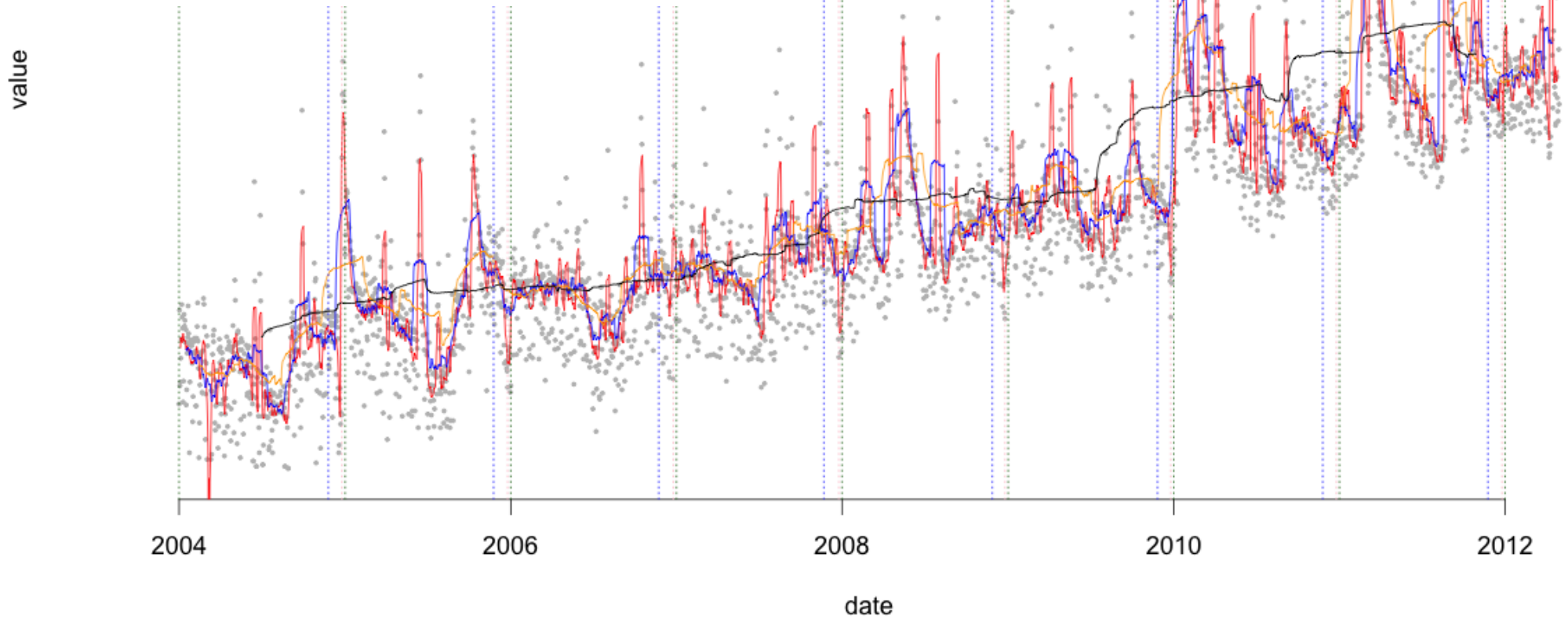
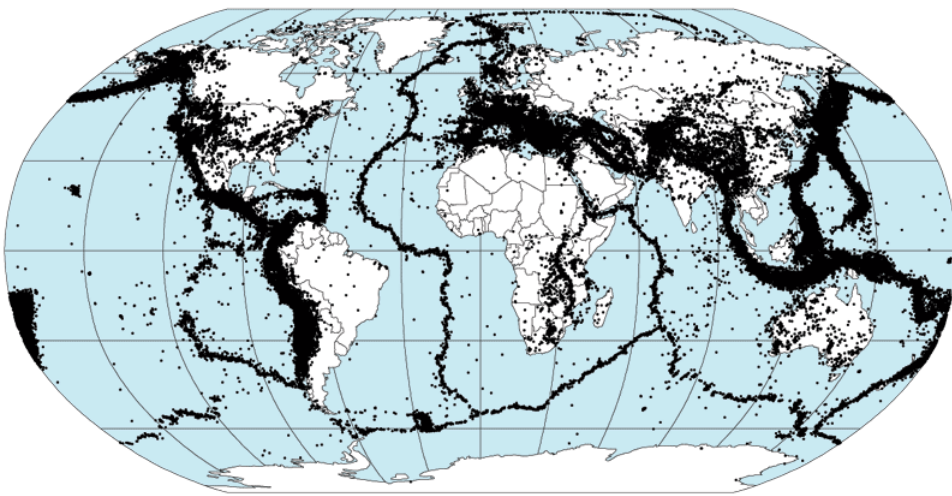
Tuned forecast (normalized values)

Bing

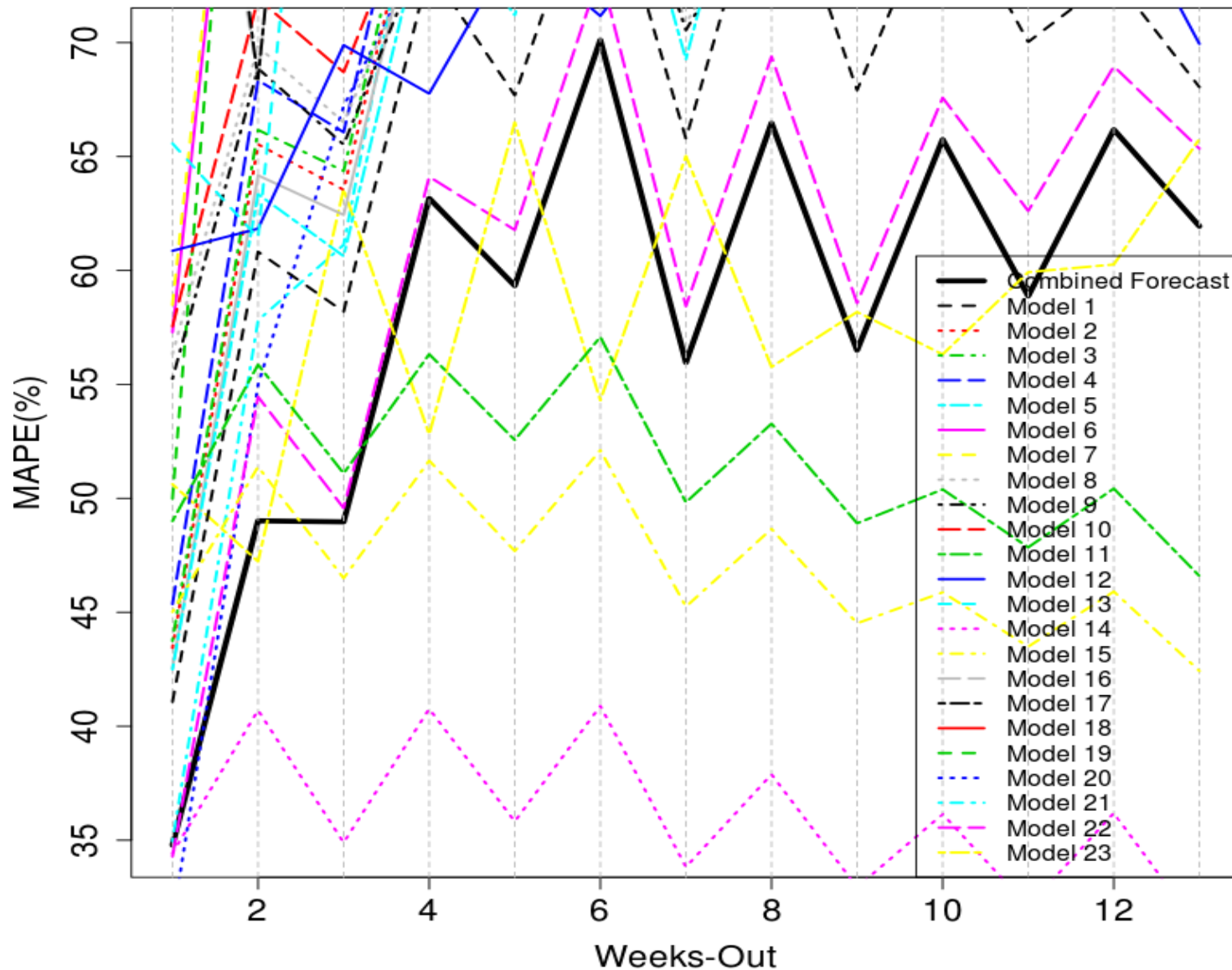


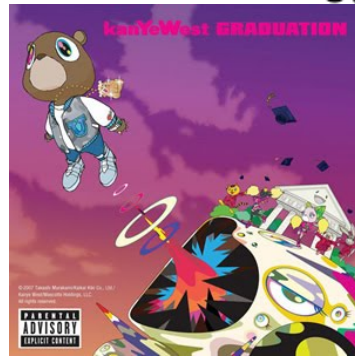
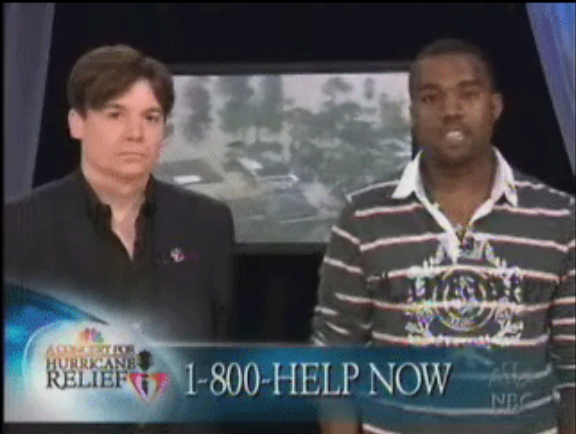
Preliminary Determination of Epicenters
358,214 Events, 1963 - 1998

$\log(\text{earthquake})$



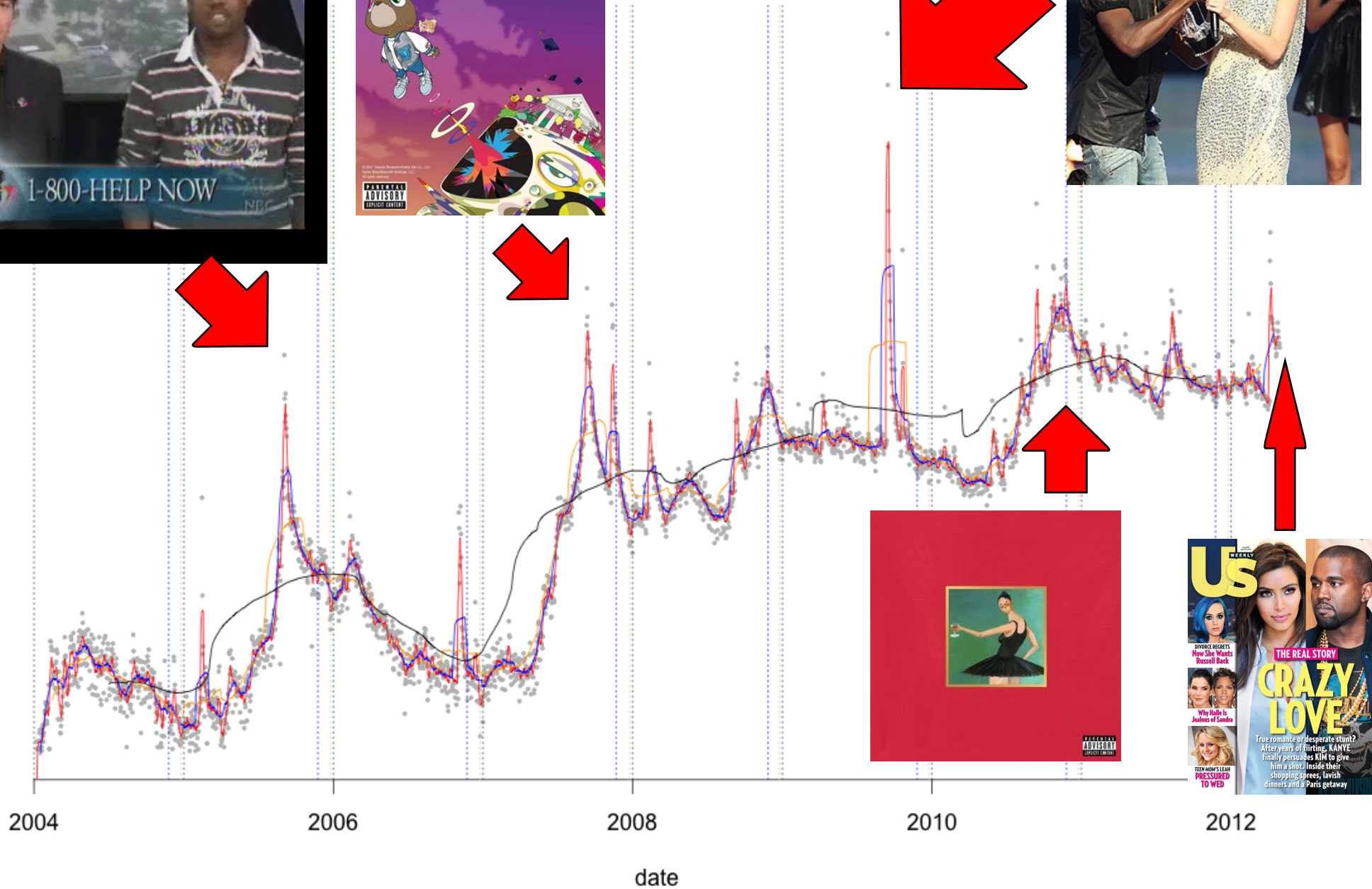
Performance on query "earthquake"





log(kanye west)

value



Celebrity earthquakes



Wrapping Up...



Conclusions, Statistical Methods

- Automatic, robust method needed for large number of diverse time series
- Many Models (MM) approach does well for large percentage of 1000 Trends queries
- Overall, MM beats individual models
- Some models do better on particular queries, but how to know *a priori*?

Conclusions, Computing

- Large number of time series & statistical method => large number of forecasts
- Number of forecasts => parallelization
- Parallelization: `google.apply()` internally, external options available
- Our version cut run time by factor of 300 (4 months -> 9 hrs)

See also:

JSM2011 Paper, <http://research.google.com/pubs/pub37483.html>) & JSM2012 Presentation (K. Millar, "Scaling R to Internet Scale Data", **upcoming** 8/1/12 @ 11:05am)



Future Work

- Classification?
- Correlations / multivariate methods?
- Aggregate / disaggregate forecasting?
- Dimension reduction? Then forecast?
- More models?

That's All, Folks!

