**BANK OF ENGLAND**

# Advanced Analytics at the Bank of England

Presentation to the Riksbank conference on Big Data: Building Strategies for Central Banks in Light of the Data Revolution

9 September 2015

# Outline

- Why are we interested?

- What are we interested in?

  – Matched micro data sets

  – Text mining

  – Visualisation

  – Machine learning


- But… there are no free lunches, so what's the bill?

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**

# Why are we interested in Big Data?

- What do we mean by the term?
  - Fuzzy meaning, covering data, techniques and attitude
- Why are we interested?
  - Change of responsibilities
    - The arrival of the PRA
  - Change of opportunity
    - More data, increased computing power, technical advances
  - Change of circumstances
    - Lessons from the financial crisis
  - Change of philosophy
    - Inductive vs deductive reasoning

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**

# Why are we interested in Big Data?

- What do we mean by the term
  - Very loose meaning, covering data, techniques and attitude
- Why are we interested?
  - Change of responsibilities
    - The arrival of the PRA
  - Change of opportunity
    - More data, increased computing power, technical advances
  - Change of circumstances
    - Lessons from the financial crisis
  - Change of philosophy
    - Inductive vs deductive reasoning

BANK OF ENGLAND

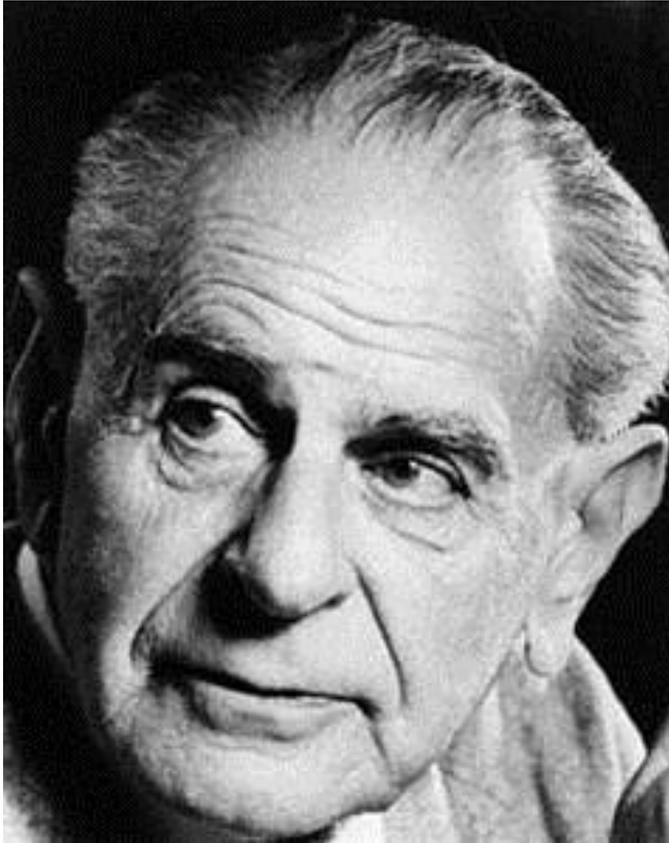**Advanced Analytics at the Bank of England**

# Why are we interested in Big Data?

- What do we mean by the term
  - Very loose meaning, covering data, techniques and attitude
- Why are we interested?
  - Change of responsibilities
    - The arrival of the PRA
  - Change of opportunity
    - More data, increased computing power, technical advances
  - Change of circumstances
    - Lessons from the financial crisis
  - Change of philosophy
    - Inductive vs deductive reasoning

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**

# Why are we interested in Big Data?

- What do we mean by the term
  - Very loose meaning, covering data, techniques and attitude
- Why are we interested?
  - Change of responsibilities
    - The arrival of the PRA
  - Change of opportunity
    - More data, increased computing power, technical advances
  - Change of circumstances
    - Lessons from the financial crisis
  - Change of philosophy
    - Inductive vs deductive reasoning

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**

# Why are we interested in Big Data?

- What do we mean by the term
  - Very loose meaning, covering data, techniques and attitude
- Why are we interested?
  - Change of responsibilities
    - The arrival of the PRA
  - Change of opportunity
    - More data, increased computing power, technical advances
  - Change of circumstances
    - Lessons from the financial crisis
  - Change of philosophy
    - Inductive vs deductive reasoning

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**

# Is Correlation the New Causality?



**Karl Popper**
(Source: http://en.wikipedia.org/wiki/Karl_Popper)



**Hal Varian**
(Source: http://en.wikipedia.org/wiki/Hal_Varian)

# What are we interested in?

- Gaining a richer understanding of the phenomenon of interest
  - Can help disentangle cause and effect…
  - …and identify the underlying issue that needs to be addressed
- Getting a speedier reading of developments in the economy and financial system
  - 'Nowcasting' and 'nearcasting'
  - This might be particularly important when the system is undergoing rapid changes
- Quantifying previously purely qualitative data
  - Eg text

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**

# Matched micro data sets

# Loan-to-income multiple ≥ 4.5



2009 Q1

0%  15%

2014 Q2

0%  15%

Source: Data are based on the Bank of England's internal Product Sales Database collected by the FCA.

**BANK OF ENGLAND**

# Discount from last asking price by buyer type



Sources: WhenFresh (Zoopla listings), Land Registry Price Paid, Land Registry Cash/Mortgage data, FCA Product Sales Data on mortgages, ONS Postcode Directory.
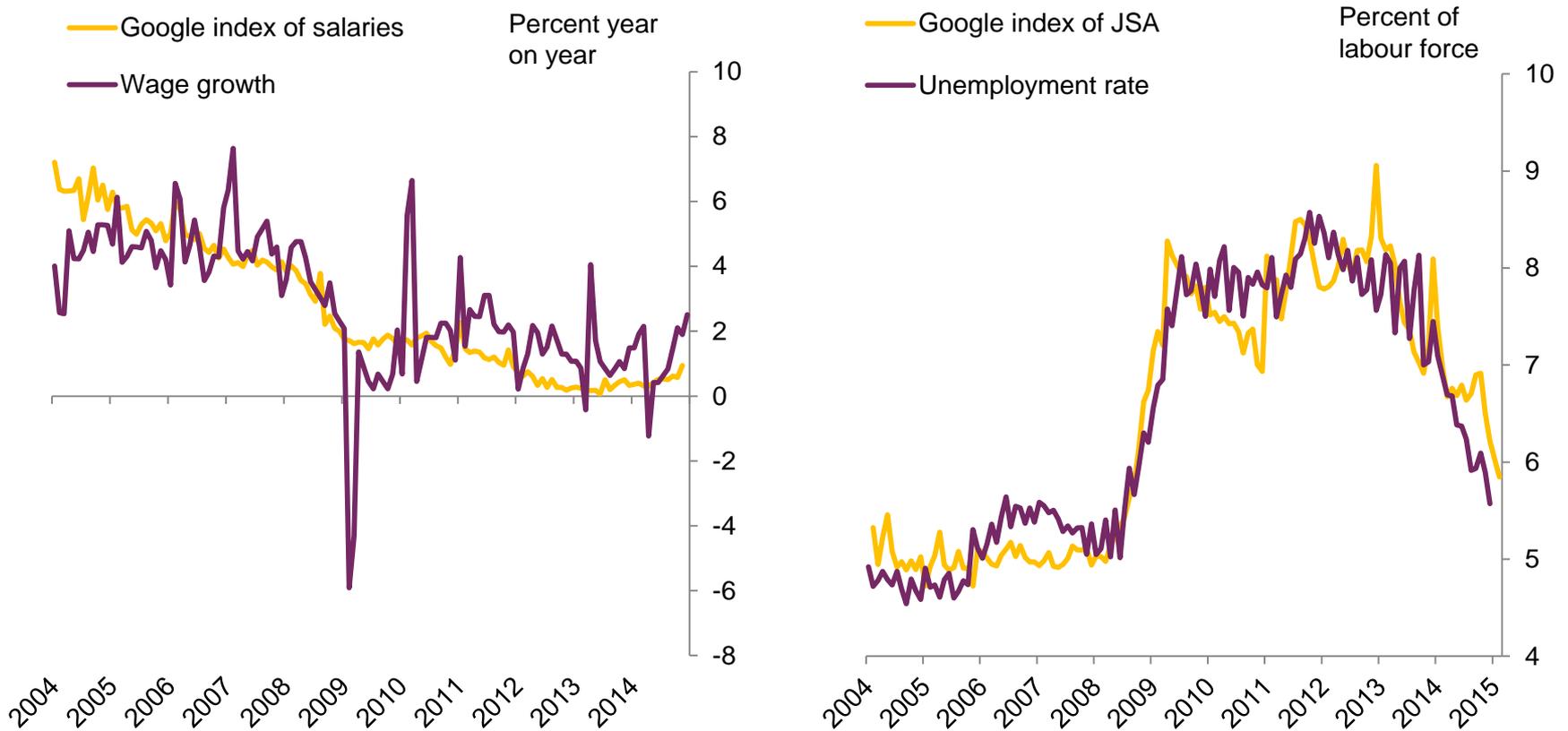
BANK OF ENGLAND

# Home owners

# Investors

Discount  -3.5%  ███████████████████████  3.5%

Sources: WhenFresh (Zoopla listings), Land Registry Price Paid, Land Registry Cash/Mortgage data, FCA Product Sales Data on mortgages, ONS Postcode Directory.

BANK OF ENGLAND

**Advanced Analytics at the Bank of England**

# Highly granular data sets

# EMIR Data

Positions in
outstanding CHF-
denominated FX
derivatives
positions on
15/1/15



BANK

# Text Analytics

# Googling the Labour Market



Source: ONS; Google. Notes: The Google indices are mean and variance adjusted to put on the same scale as the unemployment rate and wage growth. The Google indices are drawn from searches containing the terms "salaries" and "job seekers allowance". See Mclaren and Shanbhogue (2011) for further details.

Correspondence Analysis of Themes & Passive Variables
MPC Minutes & Speeches, Carney Governorship (Jul 2013-Jun 2014)

| Factor 1 | Association | Cumulative |
|---|---|---|
| Factor 1 | 33.5% | 33.5% |
| Factor 2 | 26.3% | 59.8% |

Legend:
- □ Class
- ✕ Attribute
- ○ Internal Member
- △ External Member

Class 1 (27%) — Forward Guidance
Class 2 (27%) — Productivity
Class 3 (10%) — Credit & Lending in the Economy
Class 4 (24%) — Real Economy, Inflation & Labor Markets
Class 5 (12%) — Financial Markets

BANK OF E

# Visualisation

Asset class correlation heatmap

**Advanced Analytics at the Bank of England**

# BoE communication

# Machine learning

# Issues with analysing 'Big Data'

- Example: CPI micro-data
- The ONS has produced a data set comprising:
  - 215 months (Feb 1996-Dec 2013)
  - ~110,000 prices collected per month (not the same number each month)
  - 1,113 items (not the same items each year)
  - 71 COICOP classes
  - various other meta-data (eg type of shop, region etc)
  - in total: 24,442,988 records with 25 fields
  - 611,074,700 pieces of data

# Issue 1: the stability of annual inflation



UK CPI inflation 12m rate

Percentage change over 12 months

# Issue 2: over-fitting: what to do?

- $T$ is normally >> $N$; here $N$ >>$T$ (indeed this is occasionally used as a definition of 'big data' by statistically-minded analysts)
- Aggregate the data? – eg by type of good
  - But that tends to obviate the point of using the micro data
- Shrink the dimensionality of the matrix of explanatory variables:
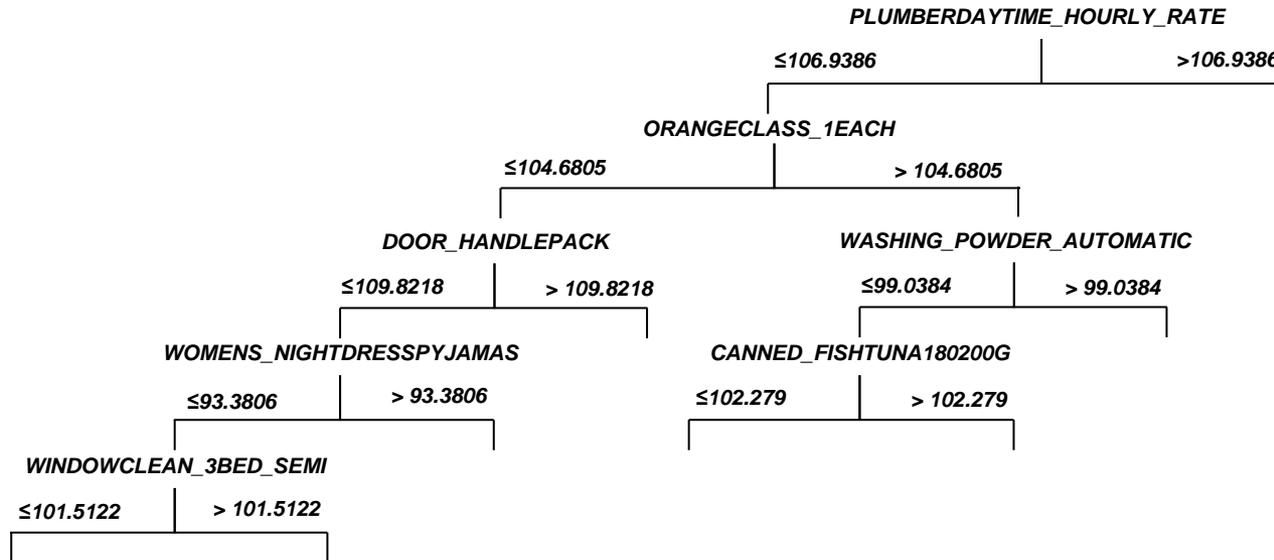  - PCA/factor models extract combinations of the variables that explain most of the variance of the dependent variable
  - So you can keep much of the information that is contained in the data set without getting into large statistical problems
  - A key issue here though is how to interpret the resulting model, and the components/factors may be unstable over time
- Penalised regressions

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**

# Issue 3: explaining non-linear functions



- Try explaining the intuition behind this relationship to busy policy makers…

# Issue 4: Stability

- An issue that is closely linked to over-fitting is the stability of the models

- This is a particularly important issues when there is no strong *a priori* reason to think that the world works in this way

- (Though *a priori* thinking can also be misleading at times)



**% false positives over 30 random samples** — % of the total number of test cases — Run number

**Positives correctly identified over 30 random samples** — % of true positives — Run number

# Issue 5: Confidentiality / 'Big Brother' state

- This was not relevant to the CPI work

- In general, the more detailed and granular the data set is, the more likely it is to contain confidential information

- We must ensure that:

  - we only use data for appropriate reasons

  - the minimum number of people are able to see any confidential data given the needs of the situation

  - data are stored securely and professionally

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**

# Issue 6: Practical issues

- Hiring!
- IT!

# Conclusion

- Do these issues mean that 'Big Data' is likely to be a passing fad?

    – No!

- The data exist and the BoE has the responsibility and opportunity to use them to help us understand economic developments and the structure of the economy and financial system

- But the issues do mean that this is no panacea

- Just as with any other empirical work the data need to be cleaned and understood (both of which are more difficult with larger data sets) prior to analysis

- And then analysed carefully using appropriate methods

**BANK OF ENGLAND**

**Advanced Analytics at the Bank of England**