

Preliminary Research into Internet Data Sources

Susan Williams and Martin Ralphs

Standards & Strategy Division

Office for National Statistics

26th June 2013

Topics

- Context: internet tools research at ONS
- Internet research topics in Beyond 2011
- Case Study: Google Insights and Migration

Internet tools research at ONS

- Internet search tools may be able to help ONS in several ways:
 - By providing quality assurance around other statistics – do patterns look the same?
 - By providing an indicator about possible change on the ground – e.g. Influx of new visitors to the UK, dispersal of migrant populations.
 - By informing statistical processes, e.g. direct inputs into models or to operational processes. This is a LONG WAY OFF!

Context: “Beyond 2011” Programme

- Identify the best way to provide small area population and socio-demographic statistics in future
- Provide a recommendation in September 2014
 - underpinned by full cost-benefit analysis
 - & high level design for implementation
- (subject to agreement) Implement the recommendation
- <http://www.ons.gov.uk/beyond2011>

Context: the caveats

- This is **VERY PRELIMINARY WORK**
- Very much an initial exploration of potential and pitfalls
- The results we have observed are interesting, but there are complex issues around meaningful interpretation
- A lot more work is needed before we can make any recommendations about the utility and robustness of Google Insights data
- **A long way from a “magic bullet”**

Our interest...

- What sort of information might we extract from internet transaction data like Google searches?
- Preliminary studies (e.g. Google Flu) suggest that these data can provide insight into emerging events
- Could we harness these data to validate other information, to flag potential events of interest and / or to improve the quality of socio-economic statistics?

Case Study: Google Insights

- The Internet is used by a large proportion of the UK population – smart phones/tablets;
- Search engines are used by most people
- The main search engine in the UK is Google
- Google have created an application to analyse the searches made through Google
- Google Insights for Search
www.google.com/insights/search/ - recently merged into www.google.com/trends/

Google Insights – basics

- Based on a **sample** of search queries (which changes each day) - sample size is not straightforward to obtain
- Weekly data (**usually**)
- Search Volume Index (SVI) – **not** actual transaction volumes!
 - Weekly proportions: Volume of searches of interest / All searches (within geographic area)
 - The week with the largest proportion is set to 100
 - All other proportions are scaled back appropriately

Compare by

- Search terms
- Locations
- Time Ranges

Search terms

Tip: Use the minus sign to exclude terms (wimbledon -tennis).

- All search terms
- + Add search term

Filter

- Web Search
- United Kingdom
- 2004 - present
- TV & Video Equipment
- All sub-regions

Search

Web Search Interest

United Kingdom, 2004 - present

All Categories > Shopping > Consumer Electronics > TV & Video Equipment

The categorisation taxonomy of Google Insights for Search has been updated during December 2011. [Learn more](#)

An improvement to our geographical assignment was applied retroactively from 1/1/2011. [Learn more](#)

Interest over time

forecast News headlines [?](#)

[How can I see numbers?](#)



The last value prior to the forecast is based on partial data and may change. [Learn more](#)
Future values are based exclusively on the extrapolation of past values. [Learn more](#)

Beyond 2011 application: in-migration

- EU expansion 2004 (EU8 + Malta/Cyprus) & 2007 (Bulgaria & Romania)
- Polish migrants the biggest EU8 population now in the UK: 500K+
- What can Google Insights tell us about national and sub-regional patterns of in-migration from EU8 countries?
- Will non-English speaking migrants be searching in Google in their native language??

Compare by

- Search terms
- Locations
- Time Ranges

Search terms

Tip: Use quotation marks to match an exact phrase ("table tennis").

- polski
- [+ Add search term](#)

Filter

- Web Search
- United Kingdom
- 2004 - present
- All sub-regions
- All Categories

Web Search Interest: polski

United Kingdom, 2004 - present

Categories: [Reference \(25-50%\)](#), [Arts & Entertainment \(0-10%\)](#), [Law & Government \(0-10%\)](#), [more...](#)

The categorisation taxonomy of Google Insights for Search has been updated during December 2011. [Learn more](#)

An improvement to our geographical assignment was applied retroactively from 1/1/2011. [Learn more](#)

Interest over time

forecast News headlines

[How can I see numbers?](#)



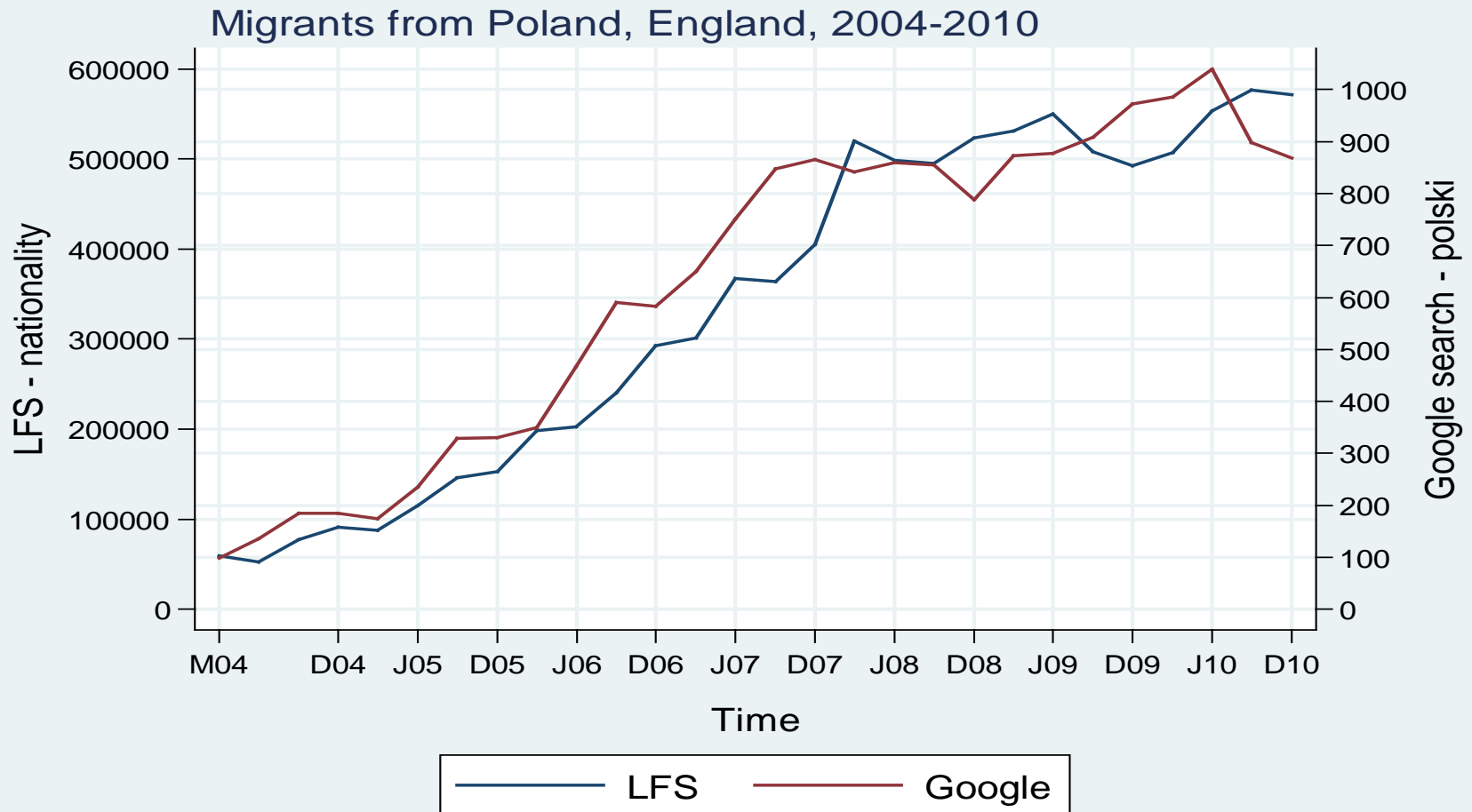
The last value on the graph is based on partial data and may change. [Learn more](#)

[Embed this chart](#)

Regional interest

[Sub-region Town/City](#)

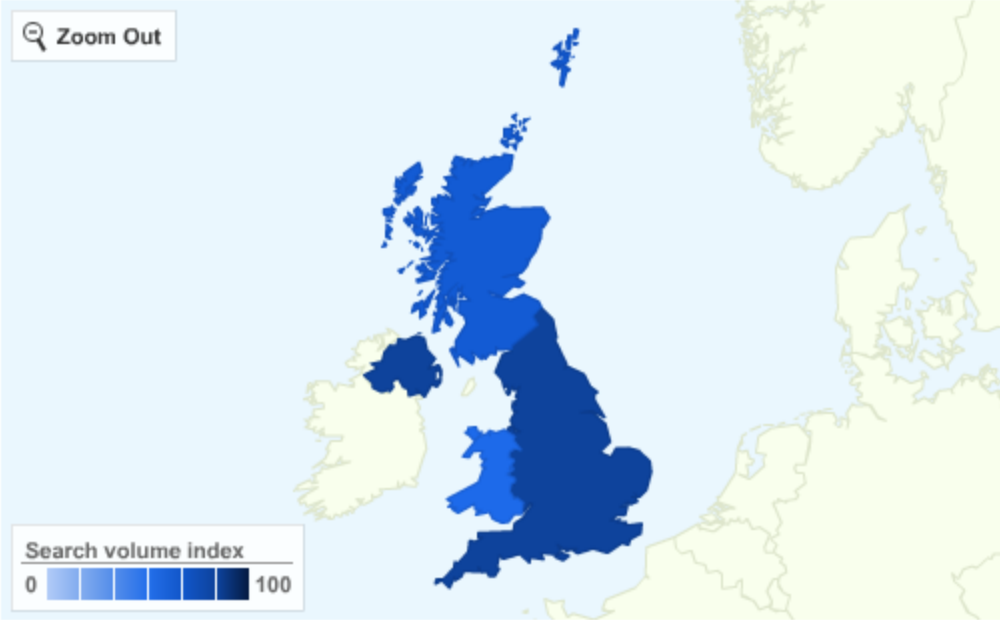
Insights SVI compared to Labour Force Survey: Polish Nationals in England



Source: Labor Force Survey, Google Insight

Regional interest

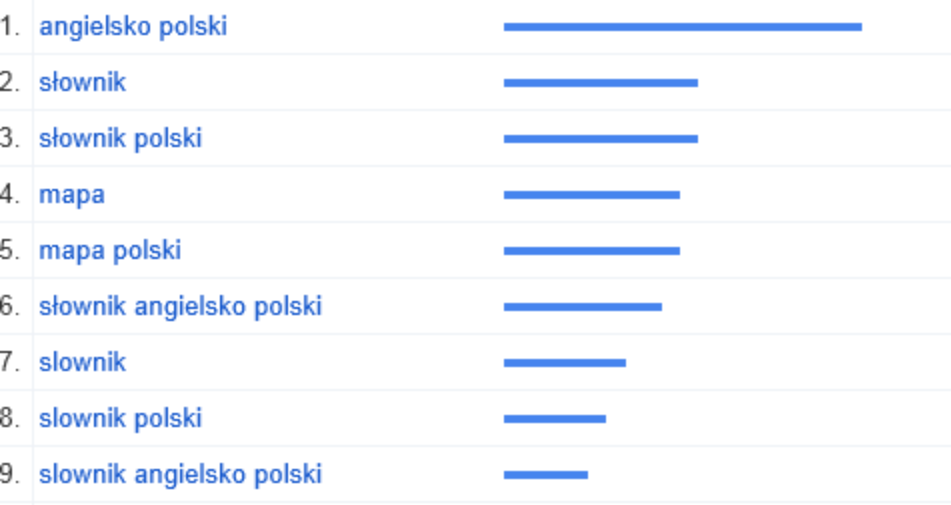
Sub-region [Town/City](#)



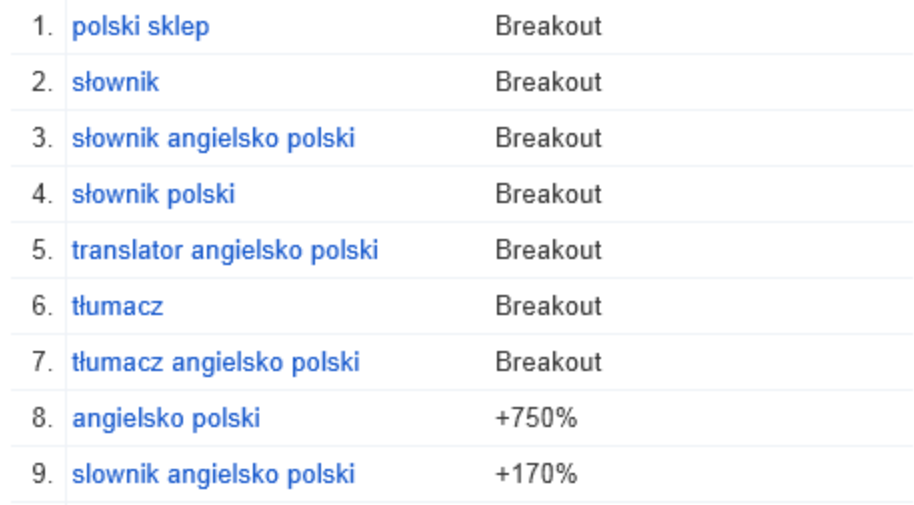
[View change over time](#)

Search terms

Top searches

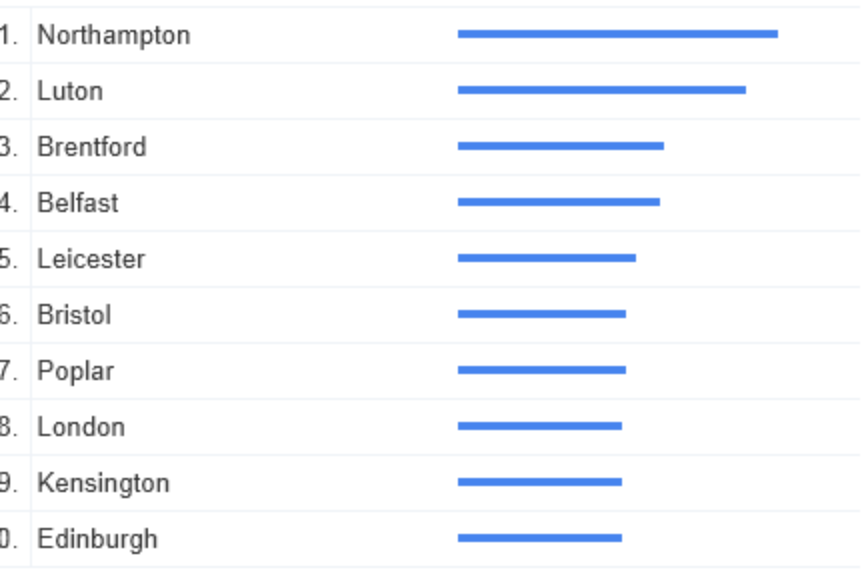


Rising searches



Regional interest

Sub-region Town/City

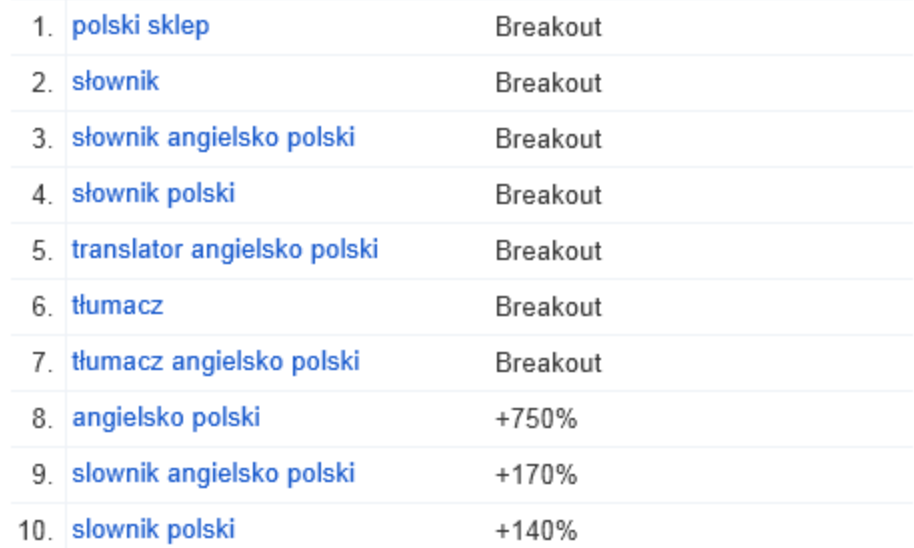


Search terms

Top searches



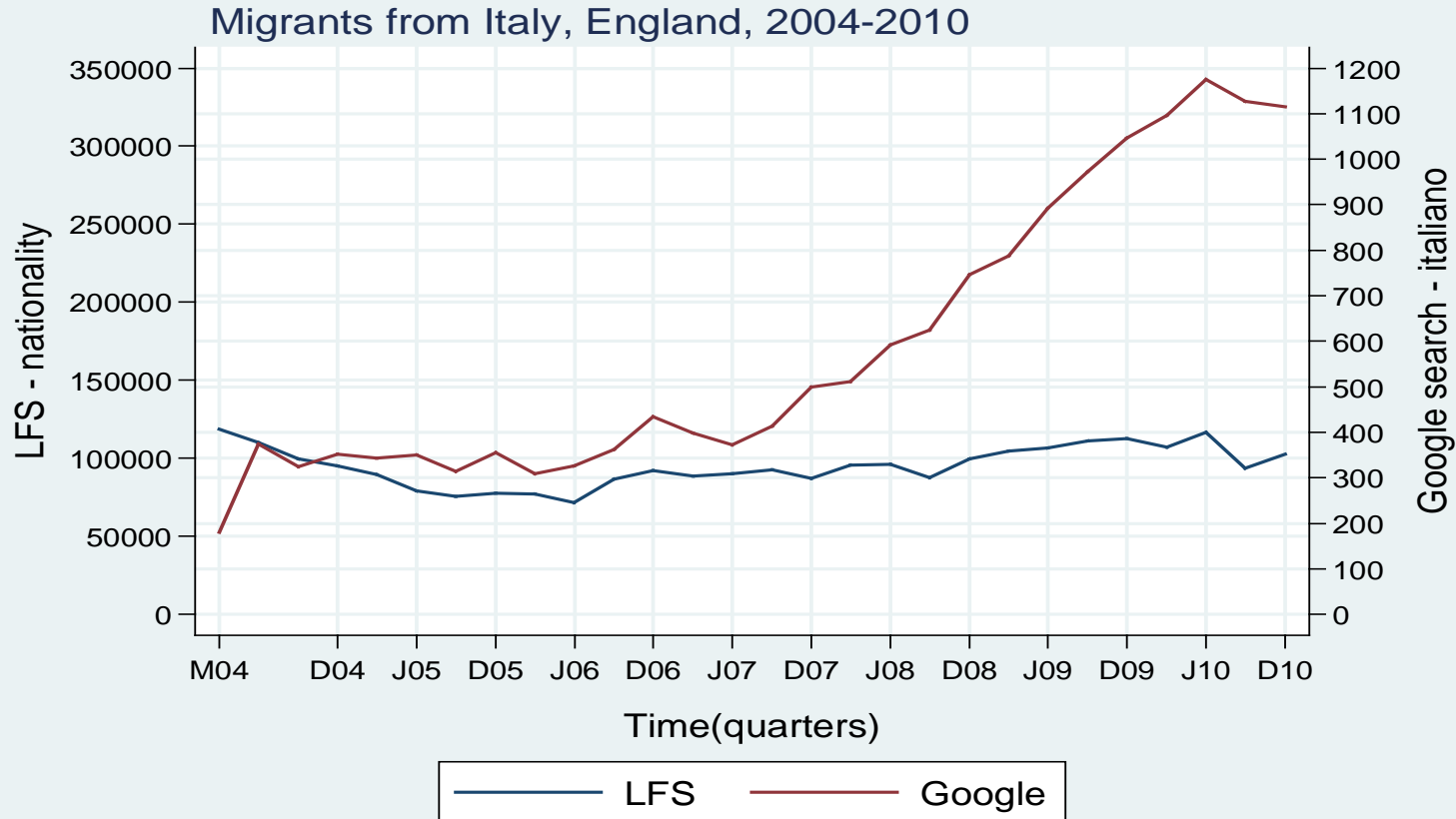
Rising searches



Investigation on other nationalities within the UK

- Promising results for Lithuanians and Romanians
- Less success with smaller EU8 populations (less than 50K)
- More established populations such as Indians, Bangladeshis and Pakistanis more difficult (English commonly spoken, multiple languages; unable to differentiate from individuals with British nationality).
- Searches on “Italiano” showed diverging timeseries from LFS up to 1 Jan 2011 where there was a steep fall in the SVI:
 - Increased short-term migration? More interest in things “Italiano”?
 - Discontinuities (at 1 Jan 2011) - corresponding to methodological change on location?

Searches on “Italiano” in England



Source: Labor Force Survey, Google Insight

Research within rest of Europe

- Searches originating in other European countries : using “polski” term
- Large drops in SVI at 1 Jan 2011 (possibly due to methodological change on location) – esp. Italy and Spain
- Germany and Austria opened their borders to the EU8 on May 1 2011.

Compare by

- Search terms
- Locations
- Time Ranges

Search terms

Tip: Use the plus sign to indicate OR (tennis + squash).

- polski
- [+ Add search term](#)

Filter

- Web Search
- Austria
- 2004 - present
- All Categories

All sub-regions

Web Search Interest: polski

Austria, 2004 - present

Categories: [Reference \(10-25%\)](#), [Arts & Entertainment \(0-10%\)](#), [People & Society \(0-10%\)](#), [more...](#)

The categorisation taxonomy of Google Insights for Search has been updated during December 2011. [Learn more](#)

An improvement to our geographical assignment was applied retroactively from 1/1/2011. [Learn more](#)

Interest over time

forecast News headlines

[How can I see numbers?](#)



The last value on the graph is based on partial data and may change. [Learn more](#)

[Embed this chart](#)

Regional interest

[Sub-region Town/City](#)

Compare by

- Search terms
- Locations
- Time Ranges

Search terms

Tip: Use a comma as shorthand to add comparison items. (tennis, squash)

- polski
- [+ Add search term](#)

Filter

- Web Search
- Germany
- 2004 - present
- All sub-regions
- All Categories

Web Search Interest: polski

Germany, 2004 - present

Categories: [Reference \(10-25%\)](#), [Arts & Entertainment \(0-10%\)](#), [Business & Industrial \(0-10%\)](#), [more...](#)

The categorisation taxonomy of Google Insights for Search has been updated during December 2011. [Learn more](#)

An improvement to our geographical assignment was applied retroactively from 1/1/2011. [Learn more](#)

Interest over time

forecast News headlines

[How can I see numbers?](#)



The last value prior to the forecast is based on partial data and may change. [Learn more](#)
Future values are based exclusively on the extrapolation of past values. [Learn more](#)

Pros of the approach

- Simple to use
- Free
- Time lag only 2-3 days
- Local areas highlighted
- SOME time series do appear to follow the Official Statistics on nationals (Polish, Lithuanian, Romanian) quite closely

Cons of the approach

- Minimal methodology and metadata available – seriously limits our ability to understand and validate what we are seeing
- Use of sample restricts its use at small geographies – and quality of local data is limited by geographical precision issues (e.g. IP address location)
- Frequent changes to methodology and tools
- Discontinuities?
- Search indices instead of empirical volumes
- Lack of consistency with official definitions
- Observed results may be caused by external factors that are irrelevant to topic of interest!

Other thoughts on applications

- There may be some application in identifying potential for sudden changes in migration – before they actually happen!
- How this might be translated into actual projections of numbers is far from clear at this stage

Compare by

- Search terms
- Locations
- Time Ranges

Search terms

Tip: Use quotation marks to match an exact phrase ("table tennis").

- london
- + Add search term

Filter

- Web Search
- Nepal
- 2004 - present
- All Categories
- All sub-regions

Search

Web Search Interest: london

Nepal, 2004 - present

The categorisation taxonomy of Google Insights for Search has been updated during December 2011. [Learn more](#)

An improvement to our geographical assignment was applied retroactively from 1/1/2011. [Learn more](#)

Interest over time

forecast News headlines

[How can I see numbers?](#)



The last value on the graph is based on partial data and may change. [Learn more](#)

Google Embed this chart

Regional interest

Town/City

Compare by

- Search terms
- Locations
- Time Ranges

Search terms

Tip: Use quotation marks to match an exact phrase ("table tennis").

- londres
- [+ Add search term](#)

Filter

- Web Search
- Spain
- 2004 - present
- All Categories
- All sub-regions

Web Search Interest: londres

Spain, 2004 - present

Categories: [Travel \(25-50%\)](#), [News \(0-10%\)](#), [Business & Industrial \(0-10%\)](#), [more...](#)

The categorisation taxonomy of Google Insights for Search has been updated during December 2011. [Learn more](#)

An improvement to our geographical assignment was applied retroactively from 1/1/2011. [Learn more](#)

Interest over time

forecast News headlines

[How can I see numbers?](#)



The last value prior to the forecast is based on partial data and may change. [Learn more](#)
Future values are based exclusively on the extrapolation of past values. [Learn more](#)

Initial thoughts on findings

- Promising early research - but limited by the lack of information about Insights methods and data quality
- May have some potential use as part of QA for dispersal of populations across the UK / sudden changes to migration patterns
- Could be enhanced by more information about language in which search terms were submitted, rather than using simple words or phrases
- Lots more work needed!