**Attendance code 16th January 2020 is**

# Intelligent Data Analysis 2020

## Prof. Martin Russell

UNIVERSITY OF
BIRMINGHAM

# Course structure 2020

- 11 x 2 hour lectures - Thursday, 11am – 1pm
    - Text retrieval (1)   *Intro*
        - TF-IDF similarity, vectorization of documents
    - Maths revision 1: Vectors and linear algebra
    - Dimension reduction
        - Principal Components Analysis (PCA)
    - Visualization  of high-dimensional data
        - PCA,  Topographic maps, t-SNE
    - Clustering and vector quantization
    - Text retrieval (2)
        - Synonym relationships, Latent Semantic Analysis (LSA), Page Rank
    - Classification

Intelligent Data Analysis 2020 – Lecture 1

UNIVERSITY OF
BIRMINGHAM

# Assessment

- Assessment
  - 1.5 hour exam in May/June - answer 3 questions from 3
- Assignment for extended module

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# Exercise sheets

- Weekly exercise sheets on Canvas
- Solutions will be published on Canvas
- Not part of formal assessment

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# Course Canvas page & C code

- All materials will go on Canvas
- Canvas site for 2020 will contain:
  - Copies of all slides and C code
  - Details of assignment - for IDA (extended)
  - Weekly exercise sheets and solutions
  - Pointers to relevant websites
- C code
  - Simple ANSII C implementations of basic techniques from the course. Compile using the MS Visual Studio .NET command line C compiler.

Intelligent Data Analysis 2020

UNIVERSITY<sup>OF</sup>
BIRMINGHAM

# Office hours

- My office hours:
  - Tuesday 2.30-3.30pm
  - Thursday 2.30-3.30pm (until 13/2/2020 – week 5)
  - Friday 2.30-3.30pm (from 21/2/2020 – week 6)
- My office hours will be in my office in the **Gisbert Kapp building, GK-N429**

Intelligent Data Analysis 2020

**UNIVERSITY**OF
**BIRMINGHAM**

# Moore's Law and disk capacity

- Moore's Law
  - *Technology performance doubles and prices halve every 18 months*

- Implications for data storage
  - Applies to disk capacity
  - We have the potential to record and store online a big proportion of what we do.

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# How much speech fits on 1TB?

- How much speech data can be stored on a 1TB disk?

- Assume:
  - 16kHz sampling rate (16,000) samples per second (Nyquist) $\max N_{Hz} \to 2N\ Hz\ sampling\ rate$
  - 16 bits per sample

- Then:
  - 1 second of speech requires 32,000 bytes
  - 1TB = $3.125 \times 10^7$s = 520,833 mins = 8,681hrs = 362 days

Intelligent Data Analysis 2020

UNIVERSITY$^{OF}$ BIRMINGHAM

# Petabytes

- 1 petabyte of disk space costs
  - $2,000,000 in 2003
  - $25,000 in 2019 (based on Amazon, $100 for 4TB!)
- 1 petabyte = $10^{15}$ bytes
  - $10^6$ – (1MB), $10^9$ – (1GB), $10^{12}$ – (1TB)
  - $10^{15}$ – zillion
  - 1 zillion used to be synonymous with infinity – an unimaginably large number!

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# A Petabyte is a lot of data…

- 1PB =
  - 20 million 4-drawer filing cabinets filled with text
  - 13.3 years of HD-TV video
- 1.5PB =
  - Combined size of the 10B photos on Facebook
- 20PB =
  - The amount of data processed by Google per day

  (A Google search will find many similar examples)

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# Accessing data – "aboutness"

- Why store these huge corpora?
- Because **information** in them is potentially **useful**
- But, how can we find the **relevant** items?
  - AV recording of a meeting contains more information than conventional minutes, but only useful with good search functions
- Need to know:
  - What each item in a corpus is **about**
  - <u>Relationships</u> between different corpus items
  - <u>Relationships</u> between 'queries' and corpus items
- Manual **indexing** impossible - deal with 'raw' data.
- Need to determine **automatically** what a text is **about**

*query — corpus*

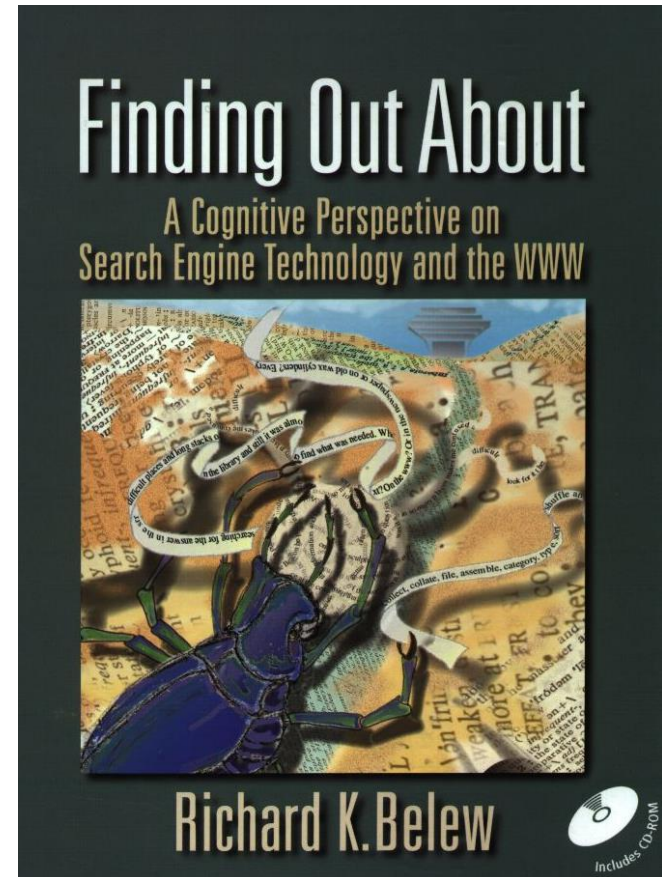Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# The problem of "Aboutness"

- What is a text, audio signal, or image **about**?
- This is a problem in **semantics**
- This is exactly the type of problem which:
  - Humans are good at, but
  - Computer programmes are particularly bad at!
- For example – "is this image about dogs?"

UNIVERSITY OF
BIRMINGHAM

# "Aboutness"

*information retrieval*

- Richard K Belew
- *Finding Out About: A Cognitive Perspective on Search Engine Technology and the WWW*
- Cambridge University Press, 2001
- Includes CD-ROM & website



Finding Out About
A Cognitive Perspective on
Search Engine Technology and the WWW

Richard K. Belew

Includes CD-ROM

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# When is an image "about" dogs

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# The problem of "aboutness"

- Intuitively, if we focus on text things should be more straightforward

- But even human interpretation of texts may be ambiguous...

- Simple example:
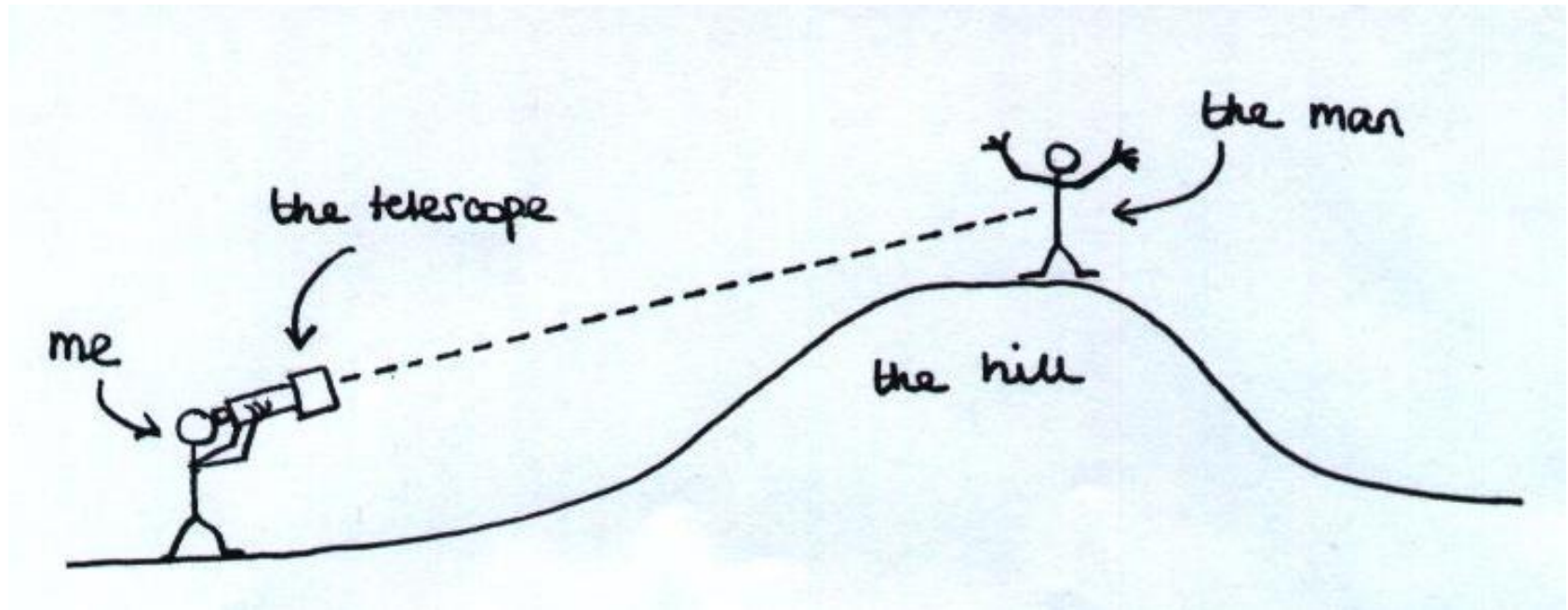
*I saw the man/on the hill/with the telescope*

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# Text Understanding

- How can a machine **understand** what this sentence is about?

- Traditionally this involves:

  - Finding the <mark>grammatical</mark> role and meaning of each word

  *Analyse* — <mark>Parsing</mark> the word sequence – applying a set of rules to identify the <mark>structure</mark> of the word sequence relative to a grammar

  - A **grammar** is a <mark>model</mark> that encodes all of the valid word sequences (sentences) in a language
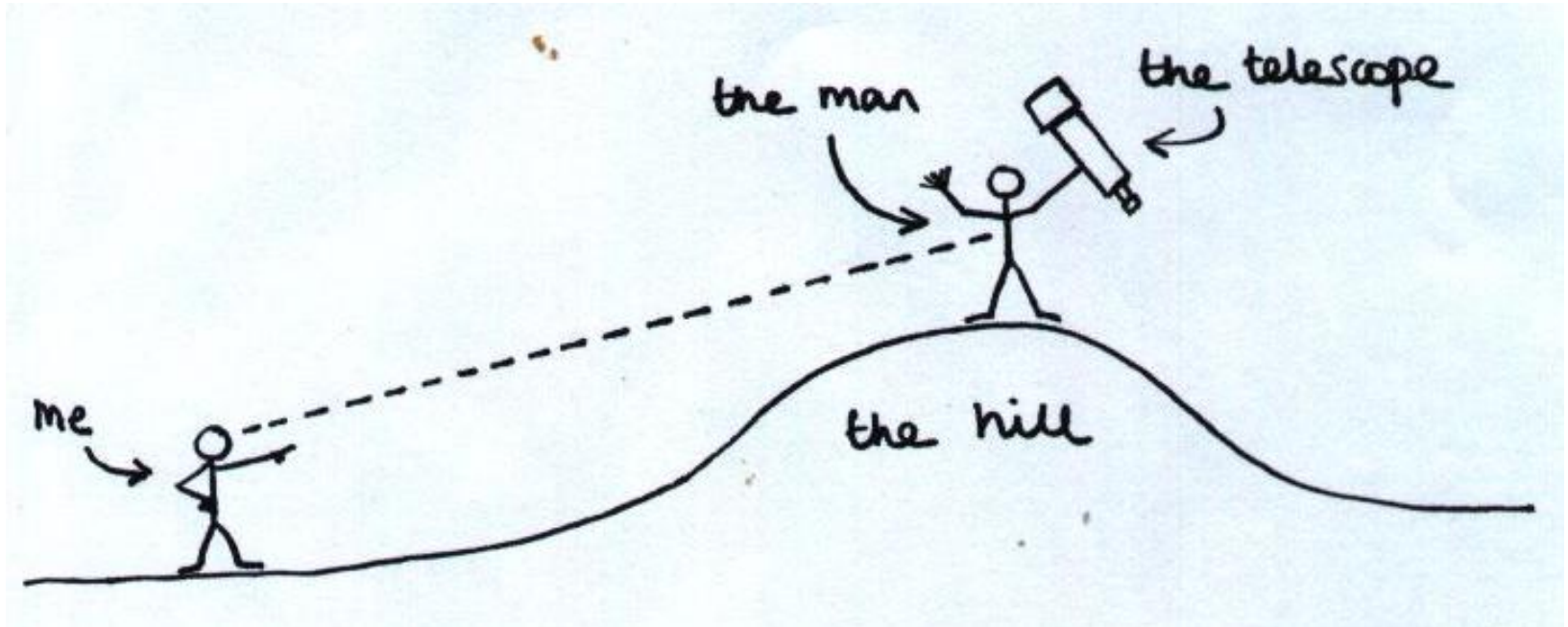
Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# Text Understanding

- Words have different meanings and grammatical roles (e.g. "lead" (verb or noun))

- A word sequence may have multiple interpretations relative to the grammar

- A grammatical word sequence may not occur in the given grammar (under generation)

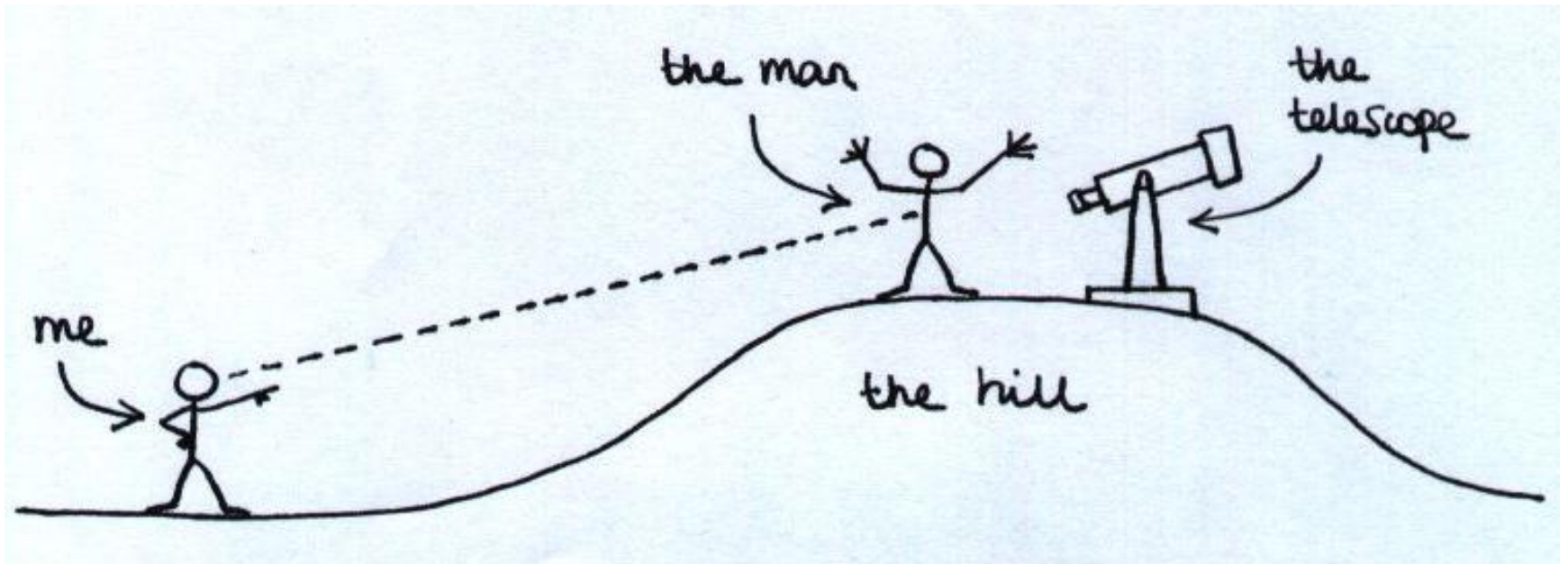- Conversely, an ungrammatical sentence may be in the grammar (over-generation)

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# *I saw the man on the hill with the telescope*



*I saw* <span style="color:red">*the man on the hill*</span> *with the telescope*

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# *I saw the man on the hill with the telescope*



*I saw the man on the hill with the telescope*

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# *I saw the man on the hill with the telescope*



*I saw the man on the hill with the telescope*

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# *I saw the man on the hill with the telescope*



*I saw the man on the hill with the telescope*

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# Analysis

- Example illustrates two different problems
  - Different <mark>grammatical parses</mark> of same word sequence
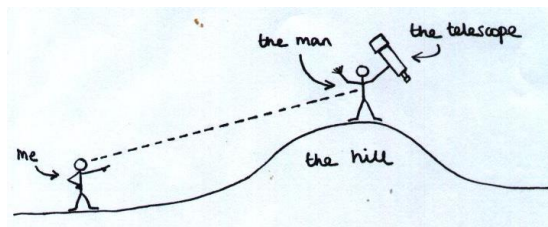
    *I saw the man on the hill with the telescope*

    *vs*

    *I saw the man on the hill with the telescope*

  - Identical parses but different <mark>interpretations of words</mark>

    *I saw the man on the hill with the telescope*



- Move towards **Machine Learning**

Intelligent Data Analysis 2020

# What is Data Mining?

- Mining
  - *Digging deep into the earth, to find hidden, valuable materials*

- Data Mining
  - Analysis of large data corpora: biomedical, acoustic, video, text,... to discover **structure**, **patterns** and **relationships**
  - Corpora too large for human inspection
  - Patterns and structure may be hidden

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# Related "hot" topics

- "Big Data"
- Pattern recognition/processing
  - As a prerequisite for Data Mining (e.g. ASR for spoken data retrieval)
  - As a consequence of Data Mining
- Machine learning
  - (Deep) Neural Networks, "Deep Learning"
- Data Visualization
  - Dimension reduction

Intelligent Data Analysis 2020
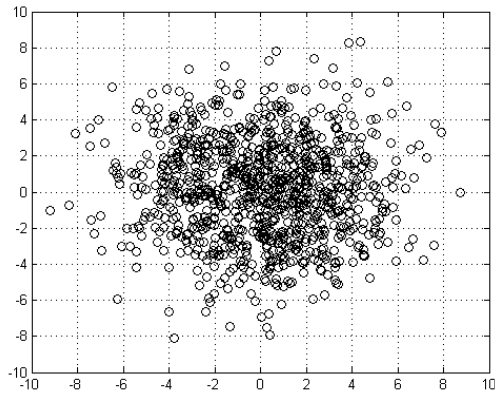
UNIVERSITY OF
BIRMINGHAM

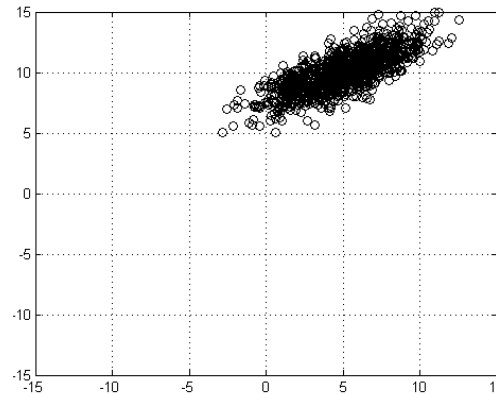# Some example data



Fig 1: Single, spherical cluster centred at origin.

Correlation X

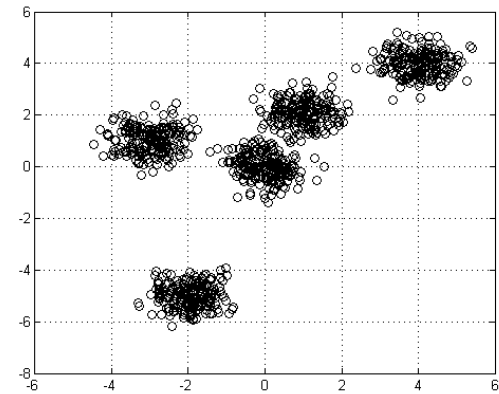Fig 2: Single, arbitrary elliptical cluster

Correlation ✓
diff. variance

Fig 3: Multiple, arbitrary elliptical clusters

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# What is Information Retrieval (IR)?

- Underlying principles of Search Engine technology
- Finding out About… [Belew]
- Retrieving Information from text sources
- Retrieving Information from other sources
  - Spoken Data Retrieval
  - Bio-informatics
- In IDA we will focus on **text retrieval**

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# IR vs Database Retrieval

- IR is <u>not</u> 'database retrieval'

- Databases are characterised by:
  - Strong prior assumptions about
    - Salient properties of data
    - Format
    - Logical relations between data items
    - Likely user queries
  - Formal, restrictive query syntax
  - Need for dedicated maintenance to keep it up-to-date
  - Gives specific replies to specific queries

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# Expedia.co.uk

**Expedia Sale! Family trips from £405 per family**

home | **deals** | flights & charters | hotels | cars | holidays | cruises | guides | maps | insurance | customer support

*SALE*

Site Map | My Trips | My Profile

**Welcome to Expedia.**

Sign in - Sign up - It's Free!

London Hotels from £52

*SALE*

Fun Breaks from £121

## BUILD YOUR PERFECT TRIP

○ Flight only
○ Hotel only
○ Car only

○ Flight + Hotel
○ Flight + Hotel + Car
○ Flight + Car

**Book together and save!**
💬 Tell me more

Departing from:
Birmingham

Going to:
Edinburgh

Depart:
5/2/2004   Evening

Return:
8/2/2004   Midday

Adults: (12-64)   2
Seniors: (65+)   0
Children: (2-11)   0
Infants: (under 2)   0

**Search**

More flight search options: Additional airports, multiple destinations...

Bag your perfect trip... *SALE*

### City Breaks
**Flight + Hotel from £99**
Including: New York, Prague, Rome...

### Luxury Breaks
**Flight + Hotel from £231**
Including: Los Angeles, Miami, St. Lucia...

### Fun Breaks
WELCOME TO LAS VEGAS NEVADA
**Flight + Hotel from £121**
Including: Amsterdam, Barcelona, Las Vegas...

### Family Trips

### Relaxing Getaways

### Ski & Snow

**TRAVELLER TOOLS**

- Airport Guides
- Arrivals/Departures
- Flight Timetables
- Currency Converter
- World Guide
- Weather

| HOTEL DEALS | from |
|---|---|
| Las Vegas | £13 |
| London | £52 |
| New York | £49 |
| Orlando | £38 |
| Paris | £37 |
| Rome | £51 |
| Amsterdam | £44 |
| UK & Ireland | £34 |
| More hotel deals... | |

| FLIGHT+HOTEL | from |
|---|---|
| Shopping Breaks | £111 |
| Caribbean Deals | £456 |
| Last Minute Deals | £97 |
| Regional Departures | £96 |
| Fly-Drive Deals | £124 |
| Florida Holidays | £283 |
| More flight+hotel deals | |

| FLIGHT DEALS | from |
|---|---|
| Barcelona | £81 |
| New York | £184 |
| Orlando | £224 |
| More flight deals... | |

| CAR DEALS | |
|---|---|
| Price per day | from |
| Portugal | £10 |
| Florida | £17 |

Start | C:\My Documents\Publica... | Expedia.co.uk Travel ... | Microsoft PowerPoint - [O...    12:52

# IR vs Database Retrieval

- IR (Finding Out About)
  - No prior assumptions about:
    - Salient properties of data
    - Format of data
    - Logical relations between data items
  - Less specific 'natural language' queries
  - Source information remains up-to-date
  - Much less focussed replies

Intelligent Data Analysis 2020

UNIVERSITY OF
BIRMINGHAM

# Relevant topics in mathematics

- **Vectors and matrices**
  - Data that we analyse is generally vector data
  - A single data point may comprise multiple measurements
  - Words or documents typically represented as vectors
  - A basic understanding of the mathematics of vectors (linear algebra) is crucial for intelligent data analysis and text retrieval

- **Probability**

- **Next lecture – linear algebra revision**

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM

# Summary

- Introduction to course components
  - Background mathematics
  - Data visualization and data mining
  - Information retrieval
- Motivation
  - Availability of huge corpora of raw data
- Problems
  - Aboutness

Intelligent Data Analysis 2020

UNIVERSITY OF BIRMINGHAM