# Unstructured Text in Big Data
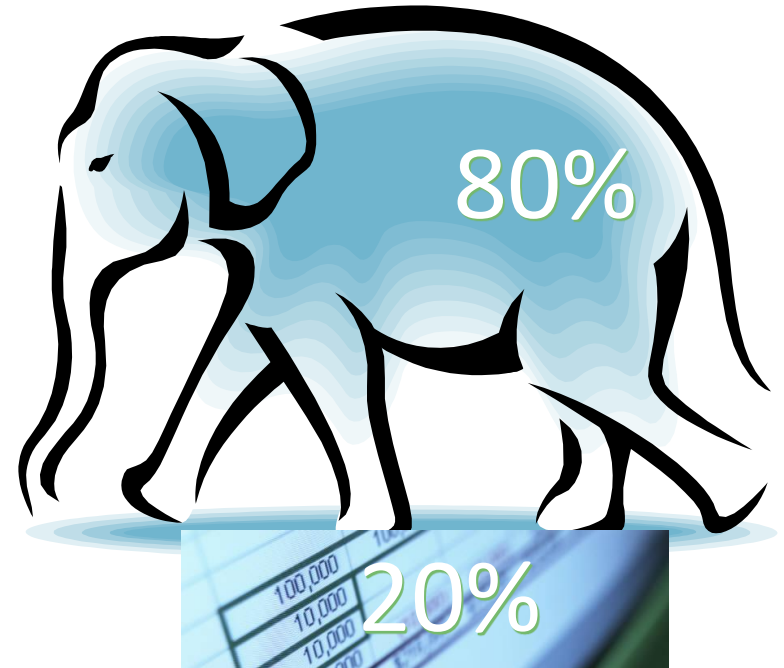## *The Elephant in  the Room*

David Milward

ICIC, October 2013

Linguamatics

# Unstructured Big Data

- Big Data
  - Volume, Variety, Velocity

- Estimated 80% of all data is unstructured
  - Need to be able to make decisions based on this data

- Ever increasing amount of data makes it harder and harder to filter by hand
  - Need more automation

# From Unstructured Data to Structured

**Identify concepts and relationships**

**Structured, relational data Normalized concept identifiers**

**Unstructured or Semi-Structured Content Sources**

Traditional text mining works in batch mode

Requires each extraction to be programmed in, or learnt from a large amount of annotated material
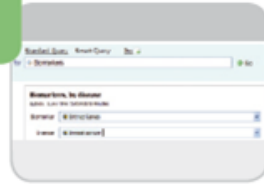
# I2E Agile Text Mining

# Increasing Range Of Applications



Sentiment Analysis in Social Media

Mining Electronic Medical Records

Vocabulary Development

Mining FDA Drug Labels

Conference Abstract Mining

Patent Analysis

Chemical Search

Extracting Numerical and Experimental Data

SAR Extraction

Competitive Intelligence

Protein-Protein Interactions

Biomarker Discovery

Systems Biology

Target Identification & Prioritization

Plant Metabolic Pathways

Key Opinion Leader Identification

In-licensing Opportunities

Safety/Tox

Patent FTO

Drug Repositioning

Linguamatics

# The Database Approach

- Can we predetermine what is required?

- One approach is to try to extract everything … but in practice have to make assumptions, whether by hand or automatically

    - Do we include negative statements?

    - Do we include speculation?

    - What context do we include?

- Need a good idea of how the data will be used to know what needs to be extracted

- There are cases where this works well e.g.

    - have useful structured data and want to add to it from free text

    - have new records structured, but want to bring in information from legacy free text

Linguamatics

# Extracting for a Relational Database

- Extracting all relevant information from pathology reports

| | | | | |
|---|---|---|---|---|
| ▼Aspirate | Eosinophils | 1% | 1 | Bone marrow aspirate is hypocellular with suboptimal stain... Differential shows 8% myeloblasts, 12% PMN/bands, 1% eos, 6% mono, 65% ... |
| | Lymphocytes | 65% | 1 | Bone marrow aspirate is hypocellular with suboptimal stain... Differential shows 8% myeloblasts, 12% ... eos, 6% mono, 65% lymph, 8% normoblasts. |
| | Monocytes | 6% | 1 | Bone marrow aspirate is hypocellular with suboptimal stain... Differential shows 8% myeloblasts, 12% ... bands, 1% eos, 6% mono, 65% lymph, 8% ... |
| | Myeloblasts | 8% | 1 | Bone marrow aspirate is hypocellular with suboptimal stain... Differential shows 8% myeloblasts, 12% PMN/bands... |
| | Normoblasts | 8% | 1 | Bone marrow aspirate is hypocellular with suboptimal stain... Differential shows 8% myeloblasts, 12% ... mono, 65% lymph, 8% normoblasts. |
| Core Biopsy | ▶ hypocellular | 5% | 1 | Bone marrow biopsy is markedly hypocellular (5% cellularity) with no evidence of fibrosis or granuloma. |
| ▼Peripheral Blood | Blasts | 19% | 1 | Our manual differential on the peripheral blood smear shows 19% blasts, 33% PMNs,11% monos, ... |
| | Lymphocytes | 37% | 1 | Our manual differential on the peripheral blood smear shows 19% blasts, 33% PMNs,11% monos, and 37% lymphs. |

Linguamatics

# Extracting for a Semantic Store

- Output as a set of triples

- Because all the parts have URIs, the parts are *web-addressable*

© 2013 Linguamatics Ltd.

# The Ad Hoc Approach

- No need to build a database

- Regard text mining queries as ways to create a database on the fly

- Keep all the unstructured free text, and query for what you need when you need it

- Even index structured data with text mining to link information together

- Used very successfully by knowledge professionals:

  – Just like a search engine, there are thousands of questions people want to ask of the data

  – You can't predict them all beforehand

# Ad Hoc Querying

- Find relationships that you want for any particular information request

- Treat the text mining as an "agile" database to answer thousands of different questions e.g. metabolites from fruit

| Strawberry | contain | Polyphenols | 2 22304566 |
|---|---|---|---|
| | | Catechin | 1 16661581 |
| | | Anthocyanins | 1 18258692 |
| | | Flavonols | 1 12381146 |

Strawberry (Fragaria × ananassa) contains several polyphenols with strong antioxidant and anti-inflammatory activities.

Douglas fir preparations contained the most complex set of procyanidins and consisted of oligomers of catechin and epicatechin, whereas strawberry and avocado contained mainly (+)-catechin and (-)-epicatechin derivatives, respectively.

Strawberry (Fragaria x ananassa) fruit contains several anthocyanins that give the ripe fruits their attractive red color.

Strawberry fruit contains flavonols as well as other phenolic compounds such as anthocyanins and phenolic acids.

Linguamatics

# Do We Need More?

- I2E Agile Text Mining has opened up text mining to the end user, but it is still mainly used by knowledge professionals

- Given the importance of unstructured big data, can we bring the benefits of text mining to an even wider audience?

Linguamatics

# Workflow Automation

- Successful queries converted into regular workflows using Pipeline Pilot or KNIME

- Real-time workflows for up-to-date dashboards, alerts

- Further analysis, visualization, integration with other data



**Pipeline Pilot**

# Clinical Trials Analysis from I2E Text Mining Results



Colours indicate Phase

Using Pipeline Pilot visualization

Sponsors

Dates

# Embedding I2E within Web Apps

# Form-Based Querying

- Create web interfaces connecting to I2E server

- Information professionals develop sophisticated queries

- Queries are parameterized to allow end users to customize

- Javascript allows portability to mobile devices such as iPads



**Select Index**

standard

**Select Query**

Gene Known as a Biomarker for Disease

**Edit Query**

Gene known as a biomarker for a disease

Find gene biomarkers for a disease

Class1  Entrez Genes

Class1 Class2 Reln  Disease Biomarker

Class2  breast ca

Breast calcifications
Breast cancer
Breast cancer stage II
Breast cancer in situ
Breast cancer recurrent
Breast cancer metastatic

Version: default 3ea2e48ca98b (2012-04-12 19:00 +0100)

Enable default limits ☑   Run Query Task: ➡

Linguamatics

# Enhancing Enterprise Search

- Use text mining to annotate concepts to feed into a search engine to:

  – provide concept search

  – provide concepts for facets

  – provide intelligent thumbnails for documents, pulling out key information

Linguamatics

# Semantic annotation of documents



Value

Keywords

Concepts

Relationships

NLP-based Text Mining

Search Engine

Linguamatics

# Provide Concepts as Facets for Enterprise Search



© 2013 Linguamatics Ltd.

# Dictionary Matching - Companies

# Pattern Matching - Institutions

| Organizations by Typ.. | | Doc | | Hit |
|---|---|---|---|---|
| Department | ▶ 576 | 21146741 | 1 | Department of Ophthalmology, Poznań City Hospital, ... |
| Center | ▶ 84 | 21146733 | 1 | Nail Disease Centre, 06400 Cannes, France... |
| Institute | ▶ 68 | 21056094 | 1 | Bioscience Institute, AmorePacific Corporation R&D Center... |
| Division | ▶ 66 | 21135170 | 1 | Division of Immunology, Allergy and Infectious Diseases... |
| University | ▶ 53 | 20971496 | 1 | Université Paris Descartes, faculté de Médecine Paris ... |
| School | ▶ 32 | 20685182 | 1 | School of Psychology, Murdoch University, Perth... |
| Corporation | ▶ 31 | 21142838 | 1 | Amygdala Ltd, Letchworth Garden City, ... |
| Laboratory | ▶ 30 | 20934246 | 1 | State Key Laboratory of Oral Diseases, West China College of ... |
| Hospital | ▶ 26 | 21109520 | 1 | Krankenhaus Hagen-Haspe, Brusebrinkstrasse 20, 58135 ... |
| Clinic | ▶ 10 | 21029565 | 1 | Paediatric Clinic, Copenhagen University Hospital, ... |
| College | ▶ 9 | 20447006 | 1 | College of Nursing, East Tennessee State University... |
| Faculty | ▶ 8 | 21158933 | 1 | Faculty of Life Sciences, University of Manchester, ... |
| Academy | ▶ 2 | 20718779 | 1 | Manchester Academic Health Science Centre, NIHR Translational Research Facility ... |
| Foundation | ▶ 2 | 20211144 | 1 | Fundação de Medicina Tropical do Amazonas, Manaus, ... |
| Network | ▶ 2 | 20865273 | 1 | Information Network of Departments of Dermatology, University of Göttingen, ... |
| Ministry | 1 | 20015772 | 1 | Ministry of Health Ankara Training and Research Hospital, 2nd E.N.T. |

Tree panel (left):

- Organizations by Type
  - Academy
  - Center
  - City
  - Clinic
  - College
  - Corporation
  - Department
  - Division
  - Government
  - Faculty
  - Foundation
  - Hospital
  - Institute
  - Laboratory
  - Ministry
  - Network
  - School
  - University

- Not a fixed list: can find previously unknown institutions

Linguamatics

# Pattern Matching - Mutations

| Mutation | Doc | | | | Hit |
|---|---|---|---|---|---|
| L858R | ▶ 72 | 20552223 | ▶ 7 | | ... , exon 19 deletions and L858R mutations, and clinical outcomes ... |
| T790M | ▶ 71 | 21233402 | ▶ 9 | | Pretreatment EGFR T790M mutation and BRCA1 mRNA expression ... |
| C8092A | ▶ 7 | 21827803 | ▶ 4 | | Polymorphisms of ERCC1 C118T/C8092A and MDR1 C3435T predict outcome ... |
| C118T | ▶ 6 | 21827803 | ▶ 4 | | Polymorphisms of ERCC1 C118T/C8092A and MDR1 C3435T ... |
| Lys751Gln | ▶ 6 | 21129812 | ▶ 4 | | Assessment of XPD Lys751Gln and XRCC1 T-77C polymorphisms in ... |
| Asn118Asn | ▶ 6 | 20354815 | ▶ 3 | | ... A60G), ERCC1 (Asn118Asn), APE1 (Asn148Glu... |
| Arg399Gln | ▶ 6 | 21805378 | ▶ 2 | | ... ERCC2 Lys751Gln, and XRCC1 Arg399Gln. |
| V600E | ▶ 5 | 21825258 | ▶ 4 | | ... mutations (56.8%) were V600E, and 16 (43.2%... |

EXPERIMENTAL DESIGN (METHODS): We assessed the T790M mutation in pretreatment diagnostic specimens from 129 erlotinib-treated advanced NSCLC patients with EGFR mutations. The expression of eight genes and two proteins involved in DNA repair and four receptor tyrosine kinases was also examined.

RESULTS (RESULTS): The EGFR T790M mutation was observed in 45 of 129 patients (35%). Progression-free survival was 12 months in patients with and 18 months in patients without the T790M mutation (P = 0.05). Progression-free survival was 27 months in patients with low BRCA1 mRNA levels, 18 months in those with intermediate levels, and 10 months in those with high levels (P = 0.02). In the multivariate analysis, the presence of the T790M mutation (HR, 4.35; P = 0.001), intermediate BRCA1 levels (HR, 8.19; P < 0.0001), and high BRCA1 levels (HR, 8.46; P < 0.0001) emerged as markers of shorter progression-free survival.

CONCLUSIONS (CONCLUSIONS): Low BRCA1 levels neutralized the negative effect of the T790M mutation and were associated with longer progression-free survival to erlotinib. We advocate baseline assessment of the T790M mutation and BRCA1 expression to predict outcome and provide alternative individualized treatment to patients based on T790M mutations and BRCA1 expression.

Linguamatics

# Pattern Matching - Chemicals

- Assign structure to novel chemicals

- Distinguish exemplified compounds

- Distinguish compounds given properties

| Chemical | Feature | Value | Hit |
|----------|---------|-------|-----|
| | mp | 209-210° | Example 31 cis-4-Hydroxy-5-phenyl-2,3,4,5-tetrahydro spiro(benzothiepine-3,1'-cyclohexane)-1,1-dioxide (60) ... mg of white crystal, mp 209-210° C. Proton and carbon NMR ... |
| | mp | 154-155° | 8b-Phenyl-1a,2,3,8b-tetrahydrospiro(benzothiepino[4,5-b]oxirene-2,1'-cyclohexane)-4,4-dioxide (58) ... ) of yellow solid, mp 154-155° C. Proton and carbon NMR ... |
| | mp | 99-100° | trans-4-Hydroxy-5-phenyl-2,3,4,5-tetrahydro spiro(benzothiepine-3,1'-cyclohexane)-1,1-dioxide (59) ... as a white solid, mp 99-100° C. Proton NMR showed this ... |

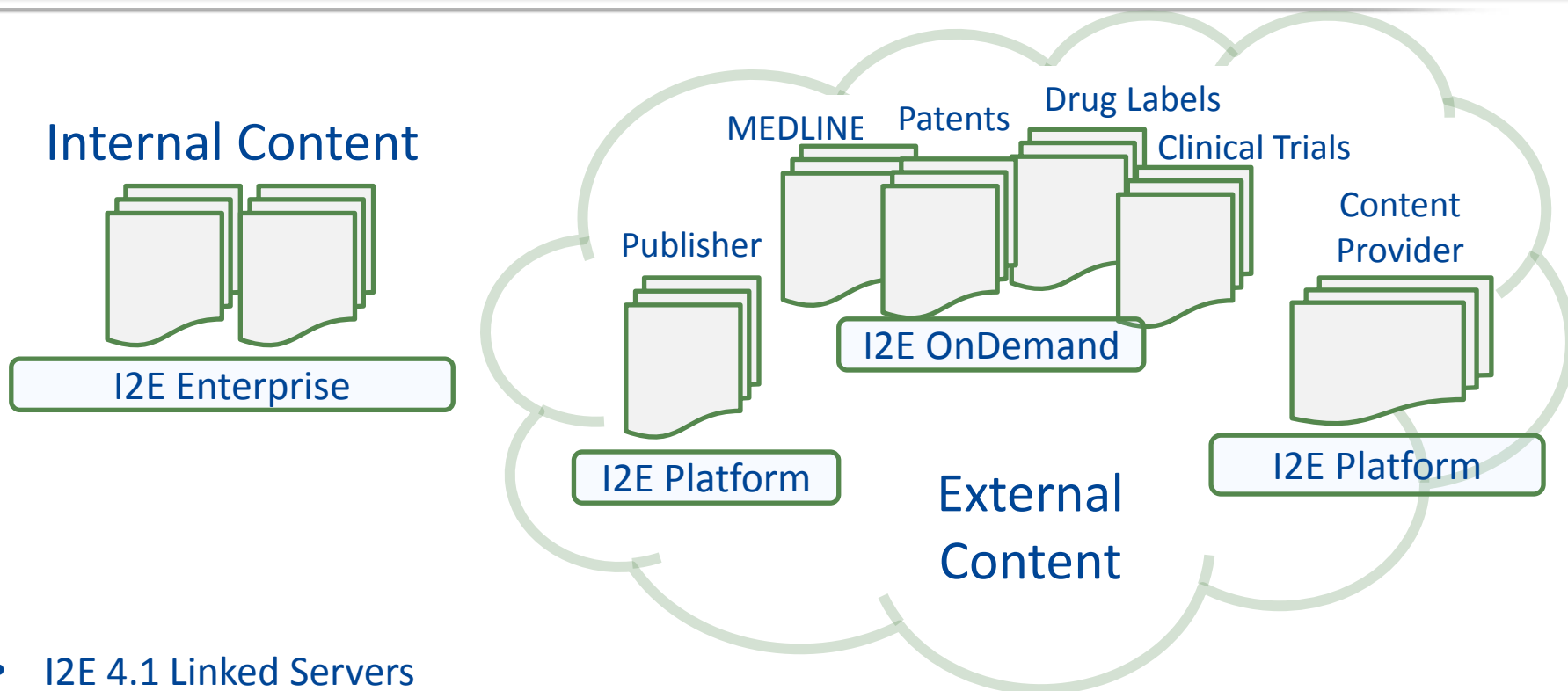# Whatever the Content…

- … I2E can mine and extract with precision

| Gene/Protein | Biomarker | Disease | | Doc | | Hit |
|---|---|---|---|---|---|---|
| ▼TRIM25 | as a new potential biomarker for | Breast cancer | 1 | 16144914 | 1 | Estrogen-responsive finger protein as a new potential biomarker for breast cancer. |
| | is a significant prognostic factor in | Breast cancer | 1 | 16144914 | 1 | CONCLUSIONS: Our data suggest that Efp immunoreactivity is a significant prognostic factor in breast cancer patients. |
| | was significantly correlated with | Breast cancer | 1 | 16144914 | 1 | Moreover, Efp immunoreactivity was significantly correlated with poor prognosis of breast cancer patients, and multivariate analyses of disease-free survival and overall survival for 151 breast cancer patients showed that Efp immunoreactivity was the independent marker. |
| ▼BCL2 | is a prognostic marker in | Breast cancer | 1 | 16638854 | 1 | Bcl-2 is a prognostic marker in breast cancer independently of the Nottingham Prognostic Index. |
| | is an independent predictor of | Breast cancer | 1 | 16638854 | 1 | CONCLUSION: Bcl-2 is an independent predictor of breast cancer outcome and seems to be useful as a prognostic adjunct to the NPI, particularly in the first 5 years after diagnosis. |

| Intention | Treatment | User | | Doc | | Hit |
|---|---|---|---|---|---|---|
| ▶ Just got | ▶ my flu shot | @4momama | | 1 217 | 1 | @4momama… Just got my flu shot and in meeting for the … |
| ▶ I got | ▶ my flu shot | @liilii89 | ▶ 2 | 657 | 1 | @liilii89… I got my flu shot today! |
| ▶ got | ▶ my flu shot | @aaaaaaannie | | 1 658 | 1 | @aaaaaaannie… got my flu shot but now dizzy from the … |
| ▼Getting | ▶ my flu shot | @amytnguyen | | 1 737 | 1 | @amytnguyen… Getting my flu shot on |
| | ▶ a flu shot | @oregonian | ▶ 3 | 650 | 1 | @oregonian… Getting a flu shot (or two)? |
| | ▶ flu shot | @bobbartz | | 1 767 | 1 | @bobbartz… Getting flu shot |
| | ▶ the flu shot | @akronbabygirl | | 1 262 | 1 | @akronbabygirl… Getting the flu shot? |
| | ▶ flu shots | @lauraann9 | | 1 180 | 1 | @lauraann9… Getting flu shots and having blood drawn is … |
| | ▶ my free flu shot | @brittag10 | | 1 752 | 1 | @brittag10… Getting my free flu shot! |

Linguamatics

# …. and Wherever



**Internal Content**

**I2E Enterprise**

**External Content**

Publisher

**I2E Platform**

MEDLINE  Patents  Drug Labels  Clinical Trials

**I2E OnDemand**

Content Provider

**I2E Platform**

- I2E 4.1 Linked Servers

- Users can query local information or information on the cloud

- Proprietary content can reside within Enterprise

- Standard sources e.g. patents can be hosted on the cloud, and shared by all users

- Data from content providers can reside on their sites

Linguamatics

# Bringing the Elephant Down to Size

- Data warehouses for data that you know you need

- Embedding text mining in other applications e.g. Enterprise Search to provide benefits of semantic search

- Building specialized new interfaces to provide self-service applications for end-users

- Continuing role for ad hoc text mining

    - text mining can ask almost any question of the data

    - you can never predict every question

Linguamatics