# Text Analytics for Dummies

## Seth Grimes

Alta Plana Corporation

301-270-0795 -- *http://altaplana.com*

Text Analytics Summit 2008

Workshop

June 15, 2008

# Introduction

## Seth Grimes –

Principal Consultant with Alta Plana Corporation.

Contributing Editor, *IntelligentEnterprise.com*.

Channel Expert, *B-Eye-Network.com*.

**Founding Chair, Text Analytics Summit, *textanalyticsnews.com*.**

Instructor, The Data Warehousing Institute, *tdwi.org*.

*I am not paid to promote any vendor.*

# Perspectives

## Perspective #1: You're a business analyst or other "end user."

You have lots of text, and you want an automated way to deal with it.

## Perspective #2: You work in IT.

You support end users who have lots of text.

## Perspective #3:  Other?

You just want to learn about text analytics.

# Perspectives

# Perspective #1a, 2a: Extending analysis.

You want to extend an existing business intelligence (BI) / data-mining initiative to encompass information from textual sources.

# Perspective #1b, 2b: New to analysis.

You don't do traditional data analysis (yet).

# Perspectives

## What do people do with electronic documents?

1. Publish, Manage, and Archive.

2. Index and Search.

3. Categorize and Classify according to *metadata* & contents.

4. Information Extraction.

## For textual documents, text analytics enhances #2 and enables #3 & #4.

## Text analytics can be automated or interactive.

# Key Message -- #1

If you are not analyzing text – if you're analyzing only transactional information – you're missing opportunity or incurring risk...

"Industries such as travel and hospitality and retail live and die on customer experience." – *Clarabridge CEO Sid Banerjee*

This is the "Unstructured Data" challenge

# Key Message -- #2

Text analytics can boost business results...

Organizations embracing text analytics all report having an epiphany moment when they suddenly knew more than before." – *Philip Russom, the Data Warehousing Institute*

...via established BI / data-mining programs, or independently.

Text Analytics is an answer to the "Unstructured Data" challenge

# Key Message -- #3

Some folks may need to expand their views of what BI and business analytics are about.

Others can do text analytics without worrying about BI.

Let's deal with text-BI first.  Here's an image and a quotation from a 1958 paper introducing BI as a method for processing documents and extracting knowledge...

# Business Intelligence

## What is business intelligence (BI)?

In this paper, business is a collection of activities carried on for whatever purpose, be it science, technology, commerce, industry, law, government, defense, et cetera. The communication facility serving the conduct of a business (in the broad sense) may be referred to as an intelligence system. The notion of intelligence is also defined here, in a more general sense, as "the ability to apprehend the **interrelationships of presented facts** in such a way as **to guide action towards a desired goal**."
– *Hans Peter Luhn,* A Business Intelligence System*, IBM Journal, October 1958*

## Why does BI not focus on textual documents?

**Text Analytics Summit 2008 – Workshop**

# The "Unstructured Data" Challenge

"The bulk of information value is perceived as coming from data in relational tables. The reason is that data that is structured is easy to mine and analyze."

> – *Prabhakar Raghavan, Yahoo Research, former CTO of enterprise-search vendor Verity (now part of Autonomy)*

That's where BI operates, on data in a relational table that originated in transactional systems.

Yet it's a truism that 80% of enterprise information is in "unstructured" form.

# The "Unstructured Data" Challenge

## Traditional BI feeds off:

```
"SUMLEV","STATE","COUNTY","STNAME","CTYNAME","YEAR","POPESTIMATE",
50,19,1,"Iowa","Adair County",1,8243,4036,4207,446,225,221,994,509
50,19,1,"Iowa","Adair County",2,8243,4036,4207,446,225,221,994,509
50,19,1,"Iowa","Adair County",3,8212,4020,4192,442,222,220,987,505
50,19,1,"Iowa","Adair County",4,8095,3967,4128,432,208,224,935,488
50,19,1,"Iowa","Adair County",5,8003,3924,4079,405,186,219,928,495
50,19,1,"Iowa","Adair County",6,7961,3892,4069,384,183,201,907,472
50,19,1,"Iowa","Adair County",7,7875,3855,4020,366,179,187,871,454
50,19,1,"Iowa","Adair County",8,7795,3817,3978,343,162,181,841,439
50,19,1,"Iowa","Adair County",9,7714,3777,3937,338,159,179,805,417
```

# The "Unstructured Data" Challenge

## Traditional BI feeds off:

```
"SUMLEV","STATE","COUNTY","STNAME",
50,19,1,"Iowa","Adair County",1,824
50,19,1,"Iowa","Adair County",2,824
50,19,1,"Iowa","Adair County",3,821
50,19,1,"Iowa","Adair County",4,809
50,19,1,"Iowa","Adair County",5,800
50,19,1,"Iowa","Adair County",6,796
50,19,1,"Iowa","Adair County",7,787
50,19,1,"Iowa","Adair County",8,779
50,19,1,"Iowa","Adair County",9,771
```

**CUSTOMER_DIM**

| PK | SHIP_TO_ID |
|----|------------|
| | SHIP_TO_DSC |
| | ACCOUNT_ID |
| | ACCOUNT_DSC |
| | MARKET_SEGMENT_ID |
| | MARKET_SEGMENT_DSC |
| | TOTAL_MARKET_ID |
| | TOTAL_MARKET_DSC |
| | WAREHOUSE_ID |
| | WAREHOUSE_DSC |
| | REGION_ID |
| | REGION_DSC |
| | ALL_CUSTOMERS_ID |
| | ALL_CUSTOMERS_DSC |

**CHANNEL_DIM**

| PK | CHANNEL_ID |
|----|------------|
| | CHANNEL_DSC |
| | ALL_CHANNELS_ID |
| | ALL_CHANNELS_DSC |

**UNITS_HISTORY_FACT**

| PK,FK4 | CHANNEL_ID |
|--------|------------|
| PK,FK2 | ITEM_ID |
| PK,FK3 | SHIP_TO_ID |
| PK,FK1 | MONTH_ID |
| | UNITS |

## It runs off:

**PRICE_AND_COST_HISTORY_FACT**

| PK,FK1 | ITEM_ID |
|--------|---------|
| PK,FK2 | MONTH_ID |
| | UNIT_PRICE |
| | UNIT_COST |

**PRODUCT_DIM**

| PK | ITEM_ID |
|----|---------|
| | ITEM_DSC |
| | ITEM_PACKAGE_ID |
| | FAMILY_ID |
| | FAMILY_DSC |
| | CLASS_ID |
| | CLASS_DSC |
| | TOTAL_PRODUCT_ID |
| | TOTAL_PRODUCT_DSC |

**TIME_DIM**

| PK | MONTH_ID |
|----|----------|
| | MONTH_DSC |
| | QUARTER_ID |
| | QUARTER_DSC |
| | YEAR_ID |
| | YEAR_DSC |
| | MONTH_TIMESPAN |
| | QUARTER_TIMESPAN |
| | YEAR_TIMESPAN |
| | MONTH_END_DATE |
| | QUARTER_END_DATE |
| | YEAR_END_DATE |

*Alta Plana*

# The "Unstructured Data" Challenge

Traditional BI produces:

*Alta Plana*

©Alta Plana Corporation, 2008

**Text Analytics Summit 2008 – Workshop**

# The "Unstructured Data" Challenge

## Some information doesn't come from a data file.



*www.stanford.edu/%7ernusse/wntwindow.html*

**Axin and Frat1 interact with dvl and GSK, bridging Dvl to GSK in Wnt-mediated regulation of LEF-1.**
Wnt proteins transduce their signals through dishevelled (Dvl) proteins to inhibit glycogen synthase kinase 3beta (GSK), leading to the accumulation of cytosolic beta-catenin and activation of TCF/LEF-1 transcription factors. To understand the mechanism by which Dvl acts through GSK to regulate LEF-1, we investigated the roles of Axin and Frat1 in Wnt-mediated activation of LEF-1 in mammalian cells. We found that Dvl interacts with Axin and with Frat1, both of which interact with GSK. Similarly, the Frat1 homolog GBP binds Xenopus Dishevelled in an interaction that requires GSK. We also found that Dvl, Axin and GSK can form a ternary complex bridged by Axin, and that Frat1 can be recruited into this complex probably by Dvl. The observation that the Dvl-binding domain of either Frat1 or Axin was able to inhibit Wnt-1-induced LEF-1 activation suggests that the interactions between Dvl and Axin and between Dvl and Frat may be important for this signaling pathway. Furthermore, Wnt-1 appeared to promote the disintegration of the Frat1-Dvl-GSK-Axin complex, resulting in the dissociation of GSK from Axin. Thus, formation of the quaternary complex may be an important step in Wnt signaling, by which Dvl recruits Frat1, leading to Frat1-mediated dissociation of GSK from Axin.

*www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=PubMed&cmd=Retrieve&list_uids=10428961&dopt=Abstract*

# The "Unstructured Data" Challenge

## Consider:

E-mail, news & blog articles, forum postings, and other social media.

Contact-center notes and transcripts.

Surveys, feedback forms, warranty claims.

And every kind of corporate documents imaginable.

## These sources may contain "traditional" data.

The Dow fell 46.58, or 0.42 percent, to 11,002.14. The Standard & Poor's 500 index fell 1.44, or 0.11 percent, to 1,263.85, and the Nasdaq composite gained 6.84, or 0.32 percent, to 2,162.78.

*Alta Plana*

# Search

So there's data and other interesting information in text.  How do we get at it?

Search is not the answer.  It returns documents.

Analysts want facts, answers to questions.  And what if you're unsure what question to ask?

All the same, let's think about searches and answers...

# Search

## Search involves –

  Words & phrases: search terms & natural language.

  Qualifiers: include/exclude, and/or, not, etc.

## Answers involve –

  Entities: names, e-mail addresses, phone numbers

  Concepts: abstractions of entities.

  Facts and relationships.

  Abstract attributes, e.g., "expensive," "comfortable"

  Opinions, sentiments: attitudinal data.

  ... and sometimes BI objects.

# Search

## Q&A may involve hidden knowledge:

What was the population of Paris in 1848?



## Concepts and complexity:

What's the best price for new laptop that I'll use for business trips and around the office?

## Opinion:

What do people think of the *Iron Man* movie?

## Calculation and structuring:

Who were the top 4 sales people for each product line, region, and quarter for the last two years?

**Text Analytics Summit 2008 – Workshop**

# Search

## Search is not enough.

*Search helps you find things you already know about.  It doesn't help you* **discover** *things you're unaware of.*

*Search results often lack* **relevance**.

*Search finds documents, not* **knowledge**.

*Search doesn't enable* **unified analytics** *that links data from textual and transactional sources.*

## Text analytics can make it better…

*Alta Plana*                    **Text Analytics Summit 2008 – Workshop**

# Beyond Search: Analysis

Text analytics enables results that suit the information and the user, e.g., answers –



Now on to knowledge discovery, to discerning *interrelationships of presented facts*...

# Beyond Search: Analysis



*www.washingtonpost.com/wp-srv/politics/daily/graphics/527Diagram_101704.html*

# Text Mining

|  | Search/Query (goal-oriented) | Discovery (opportunistic) |
|---|---|---|
| **Fielded Data** | Data Retrieval | Data Mining |
| **Documents** | Information Retrieval | Text Mining |

Based on Je Wei Liang, *www.database.cis.nctu.edu.tw/seminars/2003F/TWM/slides/p.ppt*

# Text Mining

## Text Mining = Data Mining of textual sources.

Clustering and classification.

Link Analysis.

Prediction.

Association rules.

~~Regression.~~

~~Forecasting.~~



## Text Mining = Knowledge Discovery in Text.

# Search can be pretty smart.

This slide and the next show dynamic, clustered search results from Grokker…



*live.grokker.com/grokker.html?query=text%20analytics&Yahoo=true&Wikipedia=true&numResults=250*

…with a zoomable display.

Clustering here utilizes statistical (text) data mining techniques to identifying cohesive groupings of retrieved documents.

A dynamic
network viz.:
the Touch-
Graph
Google-
Browser
applet

*touchgraph.com/*
*TGGoogleBrowser.php*
*?start=text%20analytics*

# Text Analytics

So text analytics enhances search, a.k.a. Information Retrieval.

It recognizes patterns in search queries to enable basic question answering.

It recognizes patterns in search results to enable clustering of results.

We want to get beyond IR to Information Extraction (IE).

First, *time out* to summarize and provide some definitions...

# Glossary

Text analytics automates what researchers, writers, scholars, and all the rest of us have been doing for years.  Text analytics –

**Applies linguistic and/or statistical techniques to extract concepts and patterns** *that can be applied to categorize and classify documents, audio, video, images.*

**Transforms "unstructured" information into data** *for application of traditional analysis techniques.*

**Unlocks meaning and relationships** *in large volumes of information that were previously unprocessable by computer.*

# Glossary

***Text Analytics*** is perhaps a superset of ***Text Mining***.

***Information Extraction (IE)*** involves pulling features – entities & their attributes, facts, relationships, etc. – out of textual sources.

***Entity***: Typically a name (person, place, organization, etc.) or a patterned composite (phone number, e-mail address).

***Concept:*** An abstract entity or collection of entities.

***Fact***: A relationship between two entities.

***Sentiment***: A valuation at the entity or higher level.

***Opinion***: A fact that involves a sentiment.

# Glossary

***Semantics***: A fancy word for meaning, as distinct from ***Syntax,*** which is structuring.

***Natural Language Processing (NLP)***: Computers hear humans.

***Parsing***: Evaluating the contents of a document.

***Tokenization***: Identification of distinct elements within a text.

***Stemming/ Lemmatization*** : Reducing variants of word bases created by conjugation, declension, case, pluralization, etc.

***Tagging***: Wrapping XML tags around distinct text elements, a.k.a. ***text augmentation***.

***POS Tagging***: Specifically identifying parts of speech.

# Glossary

**Categorization**: Specification of ways like items can be grouped.

**Clustering**: Creating categories according to statistical criteria.

**Taxonomy**: An exhaustive, hierarchical categorization of entities and concepts, either specified or generated by clustering.

**Classification**: Assigning an item to a category, perhaps using a taxonomy.

**Taxonomy**: A hierarchical categorization of entities and concepts.

**Accuracy**: How well an IE or IR task has been performed, computed as an **F-score** weighting **Precision** & **Recall**.

# Text Analytics

Typical steps in text analytics include –

Retrieve documents for analysis.

Apply statistical &/ linguistic &/ structural techniques to **identify, tag, and extract** entities, concepts, relationships, and events (features) within document sets.

Apply statistical pattern-matching & similarity techniques to **classify** documents and organize extracted features according to a specified or generated categorization / taxonomy.

– via a *pipeline* of statistical & linguistic steps.

# Text Analytics

So text analytics looks for structure that is inherent in the textual source materials.  Let's look at some of the steps.

First, we'll do a lexical analysis of a text file, essentially a basic statistical analysis of the words and multi-word terms...

**Text Analytics Summit 2008 – Workshop**

**Keyword Density & Prominence Tool v1.5b - Mozilla Firefox**

File  Edit  View  History  Bookmarks  Tools  Help  del.icio.us

http://www.ranks.nl/cgi-bin/ranksnl/spider/spider.cgi?lang=

G▾ Google

## ⊟ Phrase repeats

### Total 2 word phrases : 102 - Total Repeats : 246

| phrase | repeats | density | Prominence |
|---|---|---|---|
| text analytics | 9 | 1.26 % | 58.87 |
| of the | 6 | 0.84 % | 46.49 |
| and the | 4 | 0.56 % | 48.45 |
| e mail | 4 | 0.56 % | 62.86 |
| from sources | 4 | 0.56 % | 88.12 |
| influence networks | 4 H | 0.56 % | 76.00 |
| notes and | 4 | 0.56 % | 52.11 |
| of text | 4 | 0.56 % | 52.37 |
| to the | 4 | 0.56 % | 60.17 |
| to understand | 4 | 0.56 % | 63.55 |
| by the | 3 | 0.42 % | 34.65 |
| call center | 3 | 0.42 % | 68.96 |
| can be | 3 | 0.42 % | 81.68 |
| customer experience | 3 H | 0.42 % | 52.99 |
| enterprise feedback | 3 H | 0.42 % | 52.73 |
| experience management | 3 H | 0.42 % | 52.92 |
| feedback management | 3 H | 0.42 % | 52.66 |
| in the | 3 | 0.42 % | 41.79 |
| of opinion | 3 | 0.42 % | 69.97 |
| real time | 3 | 0.42 % | 17.01 |
| seek to | 3 | 0.42 % | 28.58 |
| sentiment analysis | 3 L,I | 0.42 % | 69.52 |
| sentiment extraction | 3 | 0.42 % | 37.29 |
| the results | 3 | 0.42 % | 33.45 |
| triggered by | 3 | 0.42 % | 26.00 |
| a decision | 2 | 0.28 % | 20.41 |
| a new | 2 | 0.28 % | 65.21 |
| analytics can | 2 | 0.28 % | 97.15 |
| analytics vendor | 2 | 0.28 % | 55.02 |
| analyze attitudinal | 2 | 0.28 % | 96.66 |
| and analyze | 2 | 0.28 % | 96.73 |
| and other | 2 | 0.28 % | 37.70 |

### Total 3 word phrases : 45 - Total Repeats : 93

| phrase | repeats | density | Prominence |
|---|---|---|---|
| customer experience management | 3 H | 0.63 % | 52.99 |
| enterprise feedback management | 3 H | 0.63 % | 52.73 |
| of text analytics | 3 | 0.63 % | 46.78 |
| analytics can be | 2 | 0.42 % | 97.15 |
| analyze attitudinal information | 2 | 0.42 % | 96.66 |
| and analyze attitudinal | 2 | 0.42 % | 96.73 |
| and survey responses | 2 | 0.42 % | 95.54 |
| applied to extract | 2 | 0.42 % | 96.94 |
| articles blog postings | 2 | 0.42 % | 96.10 |
| as articles blog | 2 | 0.42 % | 96.17 |
| as varied as | 2 | 0.42 % | 96.31 |
| attitudinal information from | 2 | 0.42 % | 96.59 |
| be applied to | 2 | 0.42 % | 97.01 |
| blog postings e | 2 | 0.42 % | 96.03 |
| call center notes | 2 | 0.42 % | 95.75 |
| can be applied | 2 | 0.42 % | 97.08 |
| center notes and | 2 | 0.42 % | 95.68 |
| ceo of text | 2 | 0.42 % | 55.24 |
| cries for help | 2 | 0.42 % | 7.70 |
| e mail call | 2 | 0.42 % | 95.89 |
| experience management enterprise | 2 H | 0.42 % | 62.65 |
| extract and analyze | 2 | 0.42 % | 96.80 |
| focus on applications | 2 | 0.42 % | 97.96 |
| from linguamatics to | 2 | 0.42 % | 81.52 |
| from sources as | 2 | 0.42 % | 96.45 |
| information from sources | 2 | 0.42 % | 96.52 |
| mail call center | 2 | 0.42 % | 95.82 |
| management enterprise feedback | 2 H | 0.42 % | 62.58 |
| notes and survey | 2 | 0.42 % | 95.61 |
| of opinion leadership | 2 | 0.42 % | 80.43 |
| online consumer forums | 2 | 0.42 % | 55.90 |
| postings e mail | 2 | 0.42 % | 95.96 |
| real time two | 2 | 0.42 % | 18.58 |

Done

![Alta Plana]

©Alta Plana Corporation, 2008

# Text Analytics

Those "tri-grams" are pretty good at describing the *Whatness* of the source text.

Lesson: "Structure" may not matter.

Shallow parsing and statistical analysis can be enough, for instance, to support classification. (But that's not BI.)

It can help you get at meaning, for instance, by studying co-occurrence of terms.

But statistical pattern matching – the bag/vector of words approach – may fall short.

# The Need for Linguistics

## Consider –

The Dow *fell* 46.58, or 0.42 percent, to 11,002.14. The Standard & Poor's 500 index fell 1.44, or 0.11 percent, to 1,263.85, and the Nasdaq composite *gained* 6.84, or 0.32 percent, to 2,162.78.

The Dow *gained* 46.58, or 0.42 percent, to 11,002.14. The Standard & Poor's 500 index fell 1.44, or 0.11 percent, to 1,263.85, and the Nasdaq composite *fell* 6.84, or 0.32 percent, to 2,162.78.

Example from Luca Scagliarini, Expert System.

## Let's try syntactic analysis of a bit of text...

*Alta Plana*

# Information Extraction

Let's see tagging in action.  We'll use GATE, an open-source tool...

**Text Analytics Summit 2008 – Workshop**
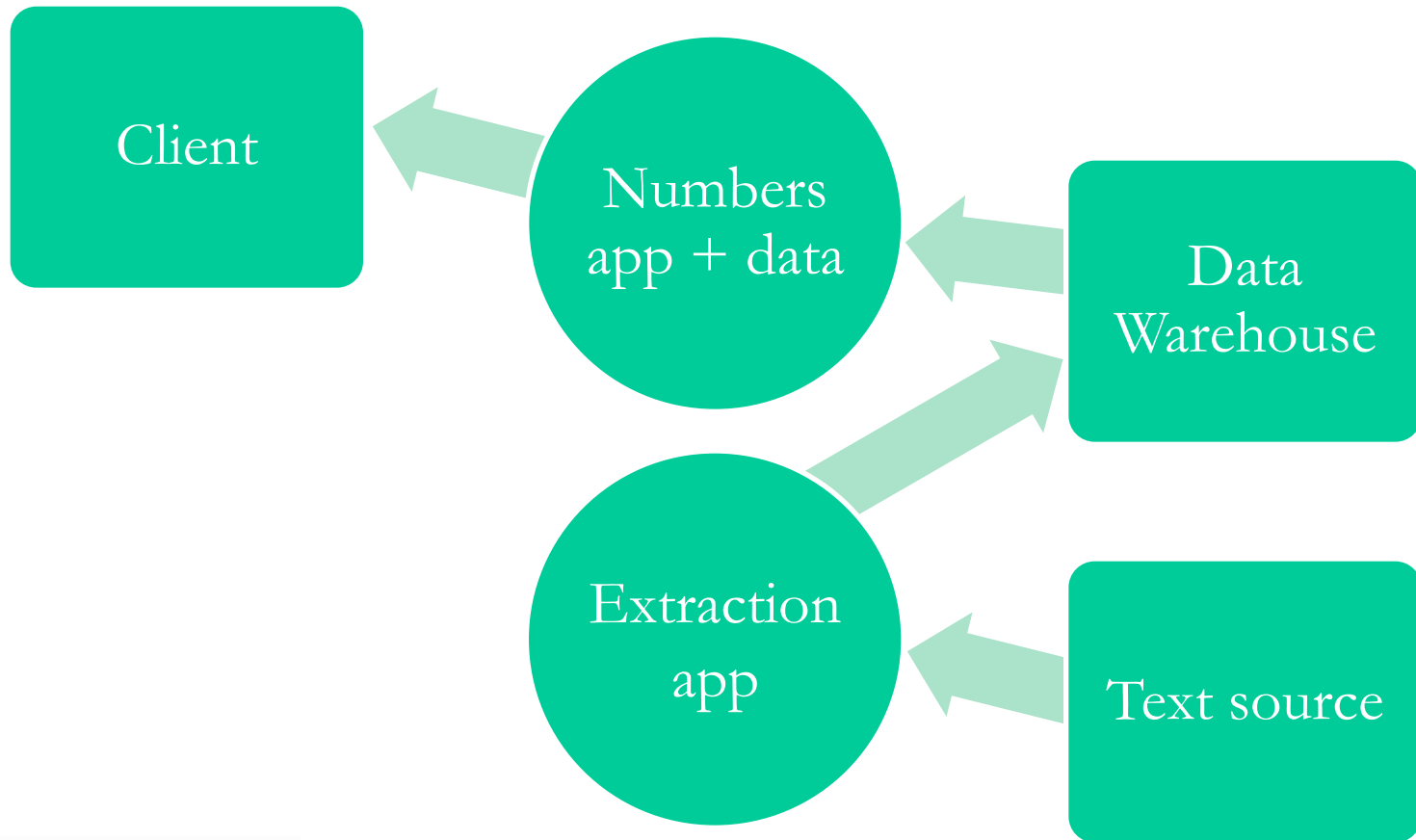
# Information Extraction

For content analysis, key in on extracting information to databases.

Entities and concepts (features) are like dimensions in a standard BI model. Both classes of object are hierarchically organized and have attributes.

We can have both discovered and predetermined classifications (taxonomies) of text features.

# Information Extraction

Data integration via information extraction.

# Information Extraction

## XML-annotated text is an intermediate format.

```
<?xml version='1.0' encoding='windows-1252'?>
<GateDocument>
<!-- The document's features-->

<GateDocumentFeatures>
<Feature>
      <Name className="java.lang.String">MimeType</Name>
      <Value className="java.lang.String">text/html</Value>
</Feature>
<Feature>
      <Name className="java.lang.String">gate.SourceURL</Name>
      <Value className="java.lang.String">http://altaplana.com/SentimentAnalysis.html</Value>
</Feature>
</GateDocumentFeatures>
<!-- The document content area with serialized nodes -->

<TextWithNodes><Node id="0" />Sentiment<Node id="9" /> <Node id="10" />Analysis<Node id="18" />:<Node
id="19" /> <Node id="20" />A<Node id="21" /> <Node id="22" />Focus<Node id="27" /> <Node id="28"
/>on<Node id="30" /> <Node id="31" />Applications<Node id="43" />
<Node id="44" />
<Node id="45" />by<Node id="47" /> <Node id="48" />Seth<Node id="52" /> <Node id="53" />Grimes<Node
id="59" />
<Node id="60" />Published<Node id="69" />:<Node id="70" /> <Node id="71" />February<Node id="79" />
<Node id="80" />19<Node id="82" />,<Node id="83" /> <Node id="84" />2008<Node id="88" />
<Node id="89" />Text<Node id="93" /> <Node id="94" />analytics<Node id="103" />
                                                                    <material cut>
</TextWithNodes>
```

# Information Extraction

## XML-annotated text...

```
<!-- The default annotation set -->
<AnnotationSet>

                                                                    <material cut>
<Annotation Id="67" Type="Token" StartNode="48" EndNode="52">
      <Feature>
            <Name className="java.lang.String">length</Name>
            <Value className="java.lang.String">4</Value>
      </Feature>
      <Feature>
            <Name className="java.lang.String">category</Name>
            <Value className="java.lang.String">NNP</Value>
      </Feature>
      <Feature>
            <Name className="java.lang.String">orth</Name>
            <Value className="java.lang.String">upperInitial</Value>
      </Feature>
      <Feature>
            <Name className="java.lang.String">kind</Name>
            <Value className="java.lang.String">word</Value>
      </Feature>
      <Feature>
            <Name className="java.lang.String">string</Name>
            <Value className="java.lang.String">Seth</Value>
      </Feature>
</Annotation>
                                                                    <material cut>
</AnnotationSet>
</GateDocument>
```

# Example: E-mail

## What else can we extract?  Let's look at an e-mail message –

Date: Sun, 13 Mar 2005 19:58:39 -0500

From: Adam L. Buchsbaum <alb@research.att.com>

To: Seth Grimes <grimes@altaplana.com>

Subject: Re: Papers on analysis on streaming data

seth, you should contact divesh srivastava, divesh@research.att.com

regarding at&t labs data streaming technology.

adam

# Example: E-mail

An e-mail message is "semi-structured."

   Semi=half.  What's "structured" and what's not?

   Is augmentation/tagging and entity extraction enough?

   What categorization might you create from that example message?

From semi-structured text, it's especially easy to extract metadata.

There are many forms of s-s information...

# Example: Survey

# Example: Survey

In analyzing surveys, we typically look at
   frequencies and distributions:



There may be fields that indicate what
   product/service/person the coded rating applies to.

Comments may be linked to coded ratings.

# Example: Survey

The respondent is invited to explain his/her attitude:

| | | | | | |
|---|---|---|---|---|---|
| My overall experience was positive. | ○ | ○ | ○ | ○ | ○ |
| **Please complete the section below if your contact with us involved permitting/licensing/registration assistance.** | | | | | |
| The regulations were understandable. | ○ | ○ | ○ | ○ | ○ |
| The application instructions were understandable. | ○ | ○ | ○ | ○ | ○ |
| The terms and conditions of the permit, license, or registration were understandable. | ○ | ○ | ○ | ○ | ○ |

**Please indicate the name(s) of any staff person you would like to commend:**

**Comments:**

**If you feel we fell short in meeting your service expectations, please describe the situation, including name of the staff person involved and the date the incident occurred:**

# Example: Survey

A survey of this type, like an e-mail message, is "semi-structured."

   Exploit what is structured in interpreting and using the free text.

   Use the *metadata* that describes the information and its provenance.

   Sentiment extraction comes into play for Voice of the Customer / Customer Experience Management applications.

# Sentiment Extraction

Sentiment (opinion) extraction –

Applications include:

Reputation management.

Competitive intelligence.

Quality improvement.

Trend spotting.

Sources include:

Wikis, blogs, forums, and newsgroups.

Media stories and product reviews.

Contact-center notes and transcripts.

Customer feedback via Web-site forms and e-mail.

Survey verbatims.

# Sentiment Extraction

We need to –

  Identify and access candidate sources.

  Extract sentiment to databases.

  Correlate expressed sentiment to measures such as:

    Sales by product, location, time, etc.

    Defects by part, circumstances, etc.

  And information such as –

    Customer information and customer's transactions.

  Correlation depends on semantic agreement: are we talking about the same things?

# Unified Analytics

Approaches build on familiar BI tools and approaches...

Adding data and text mining...

Extracting entities, facts, sentiment, etc....

Relying on semantic integration...

...for true, 360° enterprise views.

You'll learn about lots of applications over the next two days. Good luck.

*Alta Plana*

**Text Analytics Summit 2008 – Workshop**

Questions?

Discussion?

Thanks!

Seth Grimes

Alta Plana Corporation

301-270-0795 – *http://altaplana.com*