Text Analytics for BI/DW Practitioners

> Seth Grimes Alta Plana Corporation 301-270-0795 -- *http://altaplana.com*

The Data Warehousing Institute San Diego August 20, 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.

Alta Plana

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.





Introduction (done).

- The "Unstructured Data" Challenge.
- Text Technologies.
- Examples & Applications.
- Best Practices.

The Market.



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 are an image and a quotation from a 1958 paper introducing BI as a method for processing documents and extracting knowledge...

Alta Plana

Document input and processing

Knowledge handling is key



Text-BI: Back to the Future

What is business intelligence (BI)? A 1958 definition, based on processing documents: 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 notion of intelligence is also defined here... 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 was BI redefined as work on DBs?

- "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)
- BI operates on data in relational tables that originated in transactional systems.
- Yet it's a truism that 80% of enterprise information is in "unstructured" form. Alta Plana OAlta Plana Corporation, 2008 TDWI World Conference, August 2008

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



Text Analytics for BI/DW Practitioners

The "Unstructured Data" Challenge

Traditional BI feeds off:



Alta Plana

©Alta Plana Corporation, 2008

Traditional BI produces:

	Lano" • barros indigrez"						Pentaho B	usiness Intelliger	nce Platforr Portal Dem
ome Getting	Started Reporting Business R	teles Printing	Borsting V	Vidgets DetaS	ource Secur	re Advanced			
3 Filters			74 DI	teadcount Dat	ta				
alart filters to	ands to other controls on this page								
creat inters to	appretto obres controles on one page		Pos	ition			Actual	Budget	Variano
REGION .	9 Central		SVF	Strategic Dev	elopment		\$383,242	\$403,405	\$20,1
DEPARTME	NT Executive Management		CEC	Parmerships			\$307,415	\$592,100	\$24,0
our restrict	Concorre interagrenterit (200		SVE	WW Operatio	205		\$476 000	\$725.887	\$249.8
Update)		Tota	ul internet		\$	1,776,282	\$2,043,642	\$267,3
Headcour	t Costs		VA 01	Actual Headcor	unt - % Var	iance from Bud	pet		
SVP 1 Devela 38	Rategia gazada 1240 P	SVP Infraction = 307,415		CEO	SVP	Partnerships	SVP Strategic	Development SVP W	W Operations
Actual Hea	adcount - % Variance from Br	sv# attachigs= attach	Measures	CEO	SVP	Partnerships	SVP Strategic	Development SVP W	W Operations
Actual Hea	ndcount = % Variance from Bi	sufficient states and	Measures • Actual	CEO	SVP	Partnerships	SVP Strategic	Development SVP W	W Operations
Actual Hea	Adcount - % Variance from Br	sterention - 307,415	Measures • Actual 143,639,982.00	• Budget	• Variance	Partnerships	SVP Strategic	Development SVP W	W Operations
Actual Nea	Adcount - % Variance from Be Department -All Departments Executive Management	SVP "develop- 307,415 wdget Positions + All Positions + All Positions	Measures • Actual 143,639,982.00 6,299,022.00	* Budget 143,199,389.00 6,494,166.00	 Variance -440,593.00 195,144.00 	Partnerships	SVP Strategic	Development SVP W	W Operations
Actual Head	Department -All Departments Executive Management Finance	SvP patentian = 307.415 Positions +All Positions +All Positions	Measures • Actual 143,839,882.00 6,299,022.00 12,224,220.00	 Budget 143,199,389.00 6,494,166.00 12,087,406.00 	• Variance • 440,593.00 • 136,814.00	Partnerships	SVP Strategic	Development SVP W	W Operations
All Regions	Department -All Departments Executive Management Finance Human Resource	SVP addreships - 307/015 Positions +All Positions +All Positions +All Positions	Measures • Actual 143,639,982.00 6,299,022.00 12,224,220.00 13,075,463.00	CEO • Budget 143,199,389,00 6,494,166,00 12,087,406,00 12,989,341,00	Variance 440,593.00 195,144.00 -138,814.00 -66,122.00	Variance Perce	SVP Strategic	Development SVP W	W Operations
Actual Hei	Advent - % Variance from Br Department -All Departments Executive Management Finance Human Resource Marketing & Communication	DVP Jahonships= Jat2-et5 All Positions + All Positions + All Positions + All Positions + All Positions	Measures • Actual 143,639,982.00 6,299,022.00 12,224,220.00 13,075.463.00 13,910,753.00	 Budget 143, 199, 389, 00 6, 494, 166, 00 12, 967, 406, 00 12, 969, 341, 00 13, 770, 267, 00 	 Variance -440,593.00 195,144.00 -136,814.00 -86,122.00 -140,486.00 	Partnerships Partnerships Variance Perce -31 -30 -15 -66 -1.02 -	SVP Strategic	Development SVP W	W Operations
Actual Her	Addicount - % Variance from Bu Department -All Departments Executive Management Finance Human Resource Marketing & Communication Product Development	DVP Jaborahips= JUP ensitions + All Positions + All Positions + All Positions + All Positions	Measures • Actual 143,639,982.00 12,224,220.00 13,075,463.00 13,910,753.00 10,644,102.00	• Budget 143,199,389,00 6,494,166,00 12,087,406,00 12,989,341,00 13,770,267,00 10,786,611,00	 Variance -440,593.00 195,144.00 -138,814.00 -66,122.00 -140,496.00 142,509.00 	Partnerships Partnerships Variance Perce 31 30 132 -66 -100 1.32	SVP Strategic	Development SVP W	W Operations
Actual Hes	Advant - % Variance from Ba Department -All Departments Executive Management Finance Human Resource Marketing & Communication Product Development Professional Services	SVP Jaborabias Jaborabias Jaborabias All Positions - All Positions	Measures • Actual 143,839,982.00 6,299,022.00 12,224,220.00 13,075,463.00 10,844,102.00 10,844,102.00 10,844,102.00	CEO • Budget 143,199,389,00 6,494,166,00 12,087,406,00 12,989,341,00 13,770,267,00 10,786,611,00 76,088,206,00	 Variance -440,593.00 -195,144.00 -138,814.00 -140,486.00 -142,509.00 -219,443.00 	Partnerships Partnerships Variance Perce -33 -30 -113 -66 -100 1.33 -22	SVP Strategic	Development SVP W	W Operations
Actual Hei	Department -All Departments Executive Management Finance Human Resource Marketing & Communication Product Development Professional Services Sales	Positions All Positions - All Positions	Measures • Actual 143,639,982.00 6,299,022.00 12,224,220.00 13,075,643.00 13,910,753.00 10,644,102.00 76,317,649.00 11,168,773.00	CEO * Budget 143,199,389,00 12,989,341,00 12,989,341,00 13,770,267,00 10,786,611,00 76,098,208,00 10,973,392,00	 Variance -440,593.00 195,144.00 -36,122.00 -140,496.00 142,509.00 -219,443.00 -195,381.00 	Partnerships Partnerships Variance Perce -30 -115 -66 -100 1.33 -22 -21 -176 -12 -22 -21 -22 -22 -22 -22 -22 -22 -22	SVP Strategic	Development SVP W	W Operations
Actual Hei	Advector of the second	Positions All Positions • All Positions	Measures • Actual 143,639,982.00 6,299,022.00 13,075.463.00 13,910,753.00 10,844,102.00 76,317,649.00 11,168,773.00 37,893,162.00	CEO • Dudget 143,199,389,00 6,494,166,00 12,989,341,00 13,770,267,00 10,786,611,00 76,098,206,00 10,973,392,00 38,397,600,00	 Variance 440,593.00 195,144.00 196,122.00 140,485.00 142,509.00 219,443.00 195,381.00 504,438.00 	Partnerships Partnerships Variance Perce Variance 73 300 -13 -60 13 -71 13 -71 13 -31 -32 -32 -33 -32 -32 -33 -32 -33 -32 -33 -33	SVP Strategic	Development SVP W	W Operations
Actual Hes Actual Hes Region All Regions	Advent - % Variance from But Department -All Departments Executive Management Finance Human Resource Marketing & Communication Product Development Professional Services Sales +All Departments +All Departments	All Positions All Positions	Measures • Actual 143,639,982.00 6,299,022.00 12,224,220.00 13,910,753.00 10,644,102.00 76,317,649.00 11,168,773.00 37,893,162.00 35,248,940.00	CEO • Budget 143, 199, 389, 00 6, 494, 166, 00 12, 987, 406, 00 12, 989, 341, 00 10, 786, 611, 00 76, 098, 206, 00 10, 973, 392, 00 38, 397, 600, 00 38, 397, 600, 00 35, 487, 861, 00	 Variance 440,593.00 195,144.00 138,814.00 -86,122.00 140,496.00 140,496.00 142,509.00 -219,443.00 -94,438.00 504,438.00 238,921.00 	 Partnerships Partnerships Variance Perce 300 -31 -30 -40 -100 1.30 -25 -177 1.31 -61 61 61 	SVP Strategic	Development SVP W	W Operations
Actual Her Region -All Regions Central Eastern Southern	Adicount - % Variance from Bit Department -All Departments Executive Management Finance Human Resource Marketing & Communication Professional Services Sales +All Departments +All Departments +All Departments +All Departments +All Departments	DVP Jaborations Jaborations All Positions All Positions	Measures • Actual 143,639,962.00 12,224,220.00 13,910,753.00 13,910,753.00 10,644,102.00 76,317,649.00 11,168,773.00 37,893,162.00 35,248,940.00 35,248,940.00	CEO • Budget 143, 199, 389, 00 6, 494, 166, 00 12, 087, 406, 00 12, 089, 341, 00 13, 770, 267, 00 10, 786, 611, 00 76, 098, 206, 00 10, 973, 392, 00 38, 397, 600, 00 34, 803, 861, 00 34, 803, 805, 00 34, 805, 805, 00 34, 805, 805, 00 34, 805, 805, 00 34, 805, 805, 0	 Variance -440,593.00 195,144.00 -138,814.00 -66,122.00 -140,496.00 -219,443.00 -195,381.00 504,438.00 -538,921.00 -445,079.00 	Partnerships Partnerships Variance Perce -31 -30 -41 -10 -10 -13 -28 -17 -13 -28 -17 -13 -31 -31 -31 -31 -31 -31 -31 -31 -31	SVP Strategic	Development SVP W	W Operations

Alta Plana

http://www.pentaho.com/products/dashboards/

©Alta Plana Corporation, 2008

Text Analytics for BI/DW Practitioners

The "Unstructured Data" Challenge

Some information doesn't come from a data file.

April 2006, Roel Nusse

These diagrams display interactions between proteins in Wnt signaling and the approximate sites of binding. The partners are <u>hyper-linked</u> to one literature reference in PubMed. From there, one can retrieve more literature.



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

Alta Plana

©Alta Plana Corporation, 2008

Axin and Frat1 interact with dvl and GSK, bridging Dvl to GSK in Wnt-mediated regulation of LEF-1.

What proteins transduce their signals through dishevelled (DvI) 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 DvI 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 Dyl. The observation that the Dyl-binding domain of either Frat1 or Axin was able to inhibit Wnt-1-induced LEF-1 activation suggests that the interactions between DvI 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-DvI-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&Plist_uids=10428961&dopt=Abstract

What do you do when your source information looks like this?

When you walk in the foyer of the hotel it seems quite inviting but the room was very basis and smelt very badly of stale cigarette smoke, it would have been nice to be asked if we wanted a non smoking room, I know the room was very cheap but I found this very off putting to have to sleep with the smell, and it was to cold to leave the window open. Excellent location for restaurants and bars

Overall I would never sell/buy a Motorola V3 unless it is demanded. My life would be way better without this phone being around (I am being 100% serious) Motorola should pay me directly for all the problems I have had with these phones. :-(

Alta Plana

Consider again –

The purpose of intelligence is to "guide action towards a desired goal." (Luhn)

"Industries such as travel and hospitality and retail live and die on customer experience." (Banerjee)

Exercise...



Rereading the first text, what 1) goals and 2) useful information do you see? 3) How might you "structure" that information?

When you walk in the foyer of the hotel it seems quite inviting but the room was very basis and smelt very badly of stale cigarette smoke, it would have been nice to be asked if we wanted a non smoking room, I know the room was very cheap but I found this very off putting to have to sleep with the smell, and it was to cold to leave the window open. Excellent location for restaurants and bars



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.



Exercise: Organize the information in this paragraph into a table –

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.



What is the primary key? Is there a derived field?

Is there another way to represent this data?

Alta Plana

©Alta Plana Corporation, 2008

How about a (fabricated) XML representation -

```
...
</rdf:RDF>
```

Who can tell us about the mark-up here?



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...

24



Alta Plana

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.

Alta Plana

Q&A may involve hidden knowledge:

What was the population of Paris in 1848?

Concepts and complexity:

ge:

Search



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?

Alta Plana

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 search better...

Text analytics enables results that suit the information and the user, e.g., answers map massachusetts - Google Search - Mozilla Firefox

<u>File Edit View History Bookmarks Tools H</u> elp		
👍 🕶 🛶 👻 🐼 🏠 🗴 http://www.google.com/search?hl=en&q=map+massachusetts&i 🔻 🕨	G• Google	
Web Images Maps News Shopping Gmail more •	🕲 population peru - Live Search - Mozilla Firefox	
Google map massachusetts Search Preferences	File Edit View Higtory Bookmarks Tools Help Image: The state of the sta	
Web Maps Images Results 1 - 10 of about 6,500,000 for map mass brac Massachusetts maps.google.com Coogle.it Massachusetts imaps.google.com Coogle.it Coogle.it	Web Images Video Local Shopping more ~ 20*30 Search Ontions ~ Customize 1 - 10 of 11,000,000,000 for	(20°-30 (<u>About</u>) - 0.40 s SearchScan ²⁶⁷⁴
File Edit View Higtory Bookmarks Iools Help Connecticut Bindopport Bindopport Image: Connecticut Image: Connecticut	20*30 = 600 Yahool Shortcut - About Renault 20/30 - Wikipedia, the free encyclopedia The Renault 20 and Renault 30 are two executive cars produced by the French Renault 20 had two single rectangular headlights whereas the Renault 30 had en.wikipedia.org/wiki/Renault_20/30 - Cached Active 20-30 USA & Canada for young adults between the ages of 20 and 39. Provides young adults with an	SPONSOR RESU 30 20 Save on 30 20 and More! Buy, Bid, or Make an Offer now. www.eBayMotors.com 20 30 Find Bargain Prices On 20 30. www.BizRate.com
Peru Population, total: 29,041,593 2008 estimate · United States Census International Programs Center Is this useful? Yest № Peru: Population Estimated at 22 million in 1990, Peru's population has more than tripled last 50 years (it was slightly more than 7 million at 1940 census), more th over the www.ddg.com/LIS/aurelia/perpop.htm · Cached page Done	Peru State Peru Flag over the Peru Map an doubled Peru Newspapers Refugee Population Peru	
Alta Plana Orporatio	TDW/I World Conference	August 2008

Contributions to 527s active in federal

elections have not kept pace with soft

money donations to national party

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

Soft Money Game

Democrats initially ran into difficulty getting corporate chieftains and their companies to donate soft money to their upstart 527 groups, America Coming Together, The Media Pund and their fundraising arm, the Joint Victory Campaign 2004. Fundraisers turned to maverick donors, many of whom had given soft money to the Democratic Party in the past. This chart shows most donations and transfers of more than \$1 million to Democratic 527s through Sept. 30.



annucs aroanic by Sann conex, sames a annucli or the insemicion rost, not the eventse for polici integents annuc by lower by lower washing to annuc by lower the same by lower by annucli or the insemicion rost. Www.washingtonpost.com/wp-srv/politics/daily/graphics/527Diagram_101704.html



©Alta Plana Corporation, 2008

Exercise: Association rules.

Do you have a child or children living at home? Do you live in an apartment or a house?

Are there statistically significant correlations?

YesNoApt.IHouseI

Exercise: Link analysis, discovery & search.

Find something you and a person next to you have in common, e.g., school attended, the industry you work in, favorite sport, other.

(Next step is mine.)

What rule can we derive?

Text Mining = Data Mining of textual sources.

Clustering and Classification.

Link Analysis.

Association Rules.

Predictive Modelling.

Regression.

Forecasting.

Soft Money Game

Democrats initially ran into difficulty gatting corporate chieftains and their companies to donate soft mones to their upstart 252 groups, America Coming Together, The Media Pand and their fundrationizar and, the Joint Victory Campaign 2004, Fundraisers truncal to maverick donors, many of whom had given soft mones to the Democratic Party in the past. This chart shows most donations and transfers of more than \$1 million to Democratic Sci? through \$29, 50.



CS REPORTING BY SARAH COHEN, JAMES V. GRIMALDI OF THE WASHINGTON POST, AND THE CENTER FOR PUBLIC INTEGRITY. GRAPHIC BY LOUIS SPIRITO-THE WASHINGT

Text Mining = Knowledge Discovery in Text.

Alta Plana

©Alta Plana Corporation, 2008

TDWI World Conference, August 2008

Contributions to 527s active in federal

elections have not kept pace with soft

money donations to national party



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



©Alta Plana Corporation, 2008

Presentation of search results can be enhanced by discovery.

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

Alta Plana

©Alta Plana Corporation, 2008

...with a zoomable display.

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

Alta Plana



http://live.grokker.com/grokker.html?query=text analytics&Yahoo=true&Wikipedia=true&numResults=250

©Alta Plana Corporation, 2008

More results clustering...

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

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





Fi	ter-					1	e e	EXE	luoita	6
s⊦		lidden	Name				n The Bottom	$\langle \rangle$		Google
		Name	URL	Sim#	⊽		/	Text		Set
9		Data Mining, Text Mining	megaputer.c	1	-			<u> </u>		
9	2	SAS Data Mining and Te	sas.com/tec	1			QRAA			
0		National Centre for Text M	nactem.ac.uk	1						
0		text mining and web-bas	filebox.vt.edu	1						
0		Q&A: A Summary of Text	users.ox.ac	1						Enterprise
0		The BioMinT project hom	biomint.org	1					Entity	\mathcal{Y}
0	P	Data Mining and Analytic	thearling.com	1						
0		Text Analysis Info	textanalysis.i	1						
0		Text Mining Research Gr	cs.waikato.a	1				\frown		
0	W	Text analytics - Wikipedia,	en.wikipedia	10				lext-mining or	U Ven	KDnuggelss
0	ĩ	The New York Times - Br	nytimes.com	1	=					
0	1.	Slashdot: News for nerds	slashdot.org	1						
0	HDb	The Internet Movie Datab	imdb.com	1				Data		
0	 	BBC NEWS News Front	news.bbc.co	1					Breditive	Read Textmining
0	Θ	Blogger: Create your Blog	blogger.com	1						
0		MediaWiki - MediaWiki	mediawiki.org	1			(A)			/ _ \
0	ON	CNN.com - Breaking Ne	cnn.com	1			Aditonomy	T		
0	••	Welcome to Flickr - Photo	flickr.com	1				ETTOLIER H		
0	G	Google News	news.googl	1						upen
0		what does this mean	help.blogger	1						
l o		LICAI 2007 Workshop on	research iho	10						

Alta Plana


So text analytics enhances results of search, a.k.a. Information Retrieval (IR).

- It recognizes patterns and "named entities" in search queries to enable basic question answering.
- It recognizes patterns in search results to enable clustering and classification 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.



©Alta Plana Corporation, 2008

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.

Co-reference: Multiple expressions that describe the same thing.

Fact: A relationship between two entities.

Sentiment: A valuation at the entity or higher level.

Opinion: A fact that involves a sentiment.

- *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.

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.
- **Ontology**: In practice, a classification of a set of items in a way that represents knowledge, e.g., Assigning an item to a category, perhaps using a taxonomy.

A rose is a flower. A deer is an animal. A sparrow is a bird. Russia is our fatherland. Death is inevitable.

-- P. Smirnovskii, A Textbook of Russian Grammar

Alta Plana

©Alta Plana Corporation, 2008

Precision: The proportion of decisions (e.g., classifications) that are correct.

Recall: The proportion of actual correct decisions (e.g., classifications) relative to the total number of correct decisions. Find the even numbers:

9 17 **12 4** 1 **6 2** 20 **7** 3 **8** 10

Exercise: What is my Precision? What is my Recall?

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

f = 2*(precision * recall) / (precision + recall)

Alta Plana

©Alta Plana Corporation, 2008

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 documents, 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...



Keyword Density & Pr	ominence Tool v1.	5b - Mozilla Fire	=fox				
∃ile <u>E</u> dit ⊻iew Hi <u>s</u> torγ	y <u>B</u> ookmarks <u>T</u> oo	ols <u>H</u> elp del <u>.</u> i	cio.us				
存 • 🧼 • 🥑 🕻	8 🏠 🛃	TAG 🔺 http:	//www.ranks.nl/cgi-bin	/ranksnl/spider/spider.cgi?la	ng=	• 🕨 🖸	Google
							Rooks Exignds Log in
DANKS			KEVWOR		POMINEN	CE v1 5h	New Penort
NAININE			RETHOR			02 1100	non hoport
		10			— More	Domain / LIE	3L info —
Url tested : htt	p://altapla	na.com/S	entimentAna	alysis.html		Domain, or	
🗉 Details							
Comparison form							
🗉 Header data							
HTML							
E Totals, count	s, special w	ords					
1423 total words	in the file.	orda					
644 unique word	s in the file, sh	ort words inc	luded				
5 possible StopWo	rd(s) : an and i	the with www	v				
∃ Page elements							
□ Single word i	repeats						
word	repeats	density	Prominence	word	repeats	density	Prominence
sentiment	18 L,I	1.26%	46.93	for	17 L	1.19%	34.44
that	15	1.05%	55.22	text	15 L	1.05%	58.77
analytics	12 L	0.84%	52.83	from	10	0.70%	71.16
management	9 H	0.63%	50.37	analysis	9 L,I	0.63%	50.61
our	8	0.56%	20.36	are	8	0.56%	56.38
influence	7 H	0.49%	78.46	customer	7 H	0.49%	33.75
which	6	0.42%	63.18	understanding	6	0.42%	47.34
she	6	0.42%	68.22	notes	6	0.42%	51.18
have	6	0.42%	35.14	can	6	0.42%	55.43
been	6	0.42%	28.93	understand	5	0.35%	57.77
they	5	0.35%	54.28	sources	5	0.35%	87.31
not	5	0.35%	37.68	more	5	0.35%	42.90
		0.250%	55.84	mail	5	0.35%	63.50
mining	5	0.33%0					
mining extraction	5	0.35%	40.15	enterprise	5 H	0.35%	40.59
mining extraction way	5 5 4	0.35%	40.15 23.61	enterprise time	5 H 4	0.35% 0.28%	40.59 20.59
mining extraction way take	5 5 4 4	0.35% 0.28% 0.28%	40.15 23.61 14.78	enterprise time surveys	5 H 4 4 L	0.35% 0.28% 0.28%	40.59 20.59 50.39
mining extraction way take support	5 5 4 4 4 4	0.35% 0.28% 0.28% 0.28%	40.15 23.61 14.78 21.75	enterprise time surveys results	5 H 4 4 L 4 L	0.35% 0.28% 0.28% 0.28%	40.59 20.59 50.39 38.58
mining extraction way take support potential	5 5 4 4 4 4 4 4	0.35% 0.35% 0.28% 0.28% 0.28% 0.28%	40.15 23.61 14.78 21.75 39.97	enterprise time surveys results positive	5 H 4 L 4 L 4 L 4 4	0.35% 0.28% 0.28% 0.28% 0.28%	40.59 20.59 50.39 38.58 56.36

Ē

Keyword Density & Prominen	ce Tool v1.5	b - Mozilla Fi	refox				_
e <u>E</u> dit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ook	marks <u>T</u> ools	s <u>H</u> elp del	.icio.us				
╞ • 🔿 • 💽 😣 👔	8 🖃 🖡	🔒 🔺 http	://www.ranks.nl/co	ji-bin/ranksnl/spider/spider.cgi?lang=	•	G - Google	
Phrase repeats							
Total 2 word phra	ises : 10)2 - Tota	al	Total 3 word phrases : 45	- Tota	Repeat	s:93
Repeats : 246			phrase	repeats	density P	rominence	
phrase	repeats	density	Prominence	customer experience management	Зн	0.63 %	52.99
text analytics	9	1.26 %	58.87	enterprise feedback management	3 н	0.63 %	52.73
of the	6	0.84 %	46.49	of text analytics	3	0.63 %	46.78
and the	4	0.56 %	48.45	analytics can be	2	0.42 %	97.15
e mail	4	0.56 %	62.86	analyze attitudinal information	2	0.42 %	96.66
from sources	4	0.56 %	88.12	and analyze attitudinal	2	0.42 %	96.73
influence networks	4 н	0.56 %	76.00	and survey responses	2	0.42 %	95.54
notes and	4	0.56 %	52,11	applied to extract	2	0.42 %	96.94
of text	4	0.56 %	52.37	articles blog postings	2	0.42 %	96.10
to the	4	0.56 %	60.17	as articles blog	2	0.42 %	96.17
to understand	4	0.56 %	63,55	as varied as	2	0.42 %	96.31
by the	3	0.42 %	34.65	attitudinal information from	2	0.42 %	96.59
call center	3	0.42 %	68.96	be applied to	2	0.42 %	97.01
can be	3	0.42 %	81.68	blog postings e	2	0.42 %	96.03
customer experience	Зн	0.42 %	52.99	call center notes	2	0.42 %	95.75
enterprise feedback	Зн	0.42 %	52.73	can be applied	2	0.42 %	97.08
experience	2 1	0 42 0%	52.02	center notes and	2	0.42 %	95.68
management	5	0.42 70	52.92	ceo of text	2	0.42 %	55.24
feedback management	Зн	0.42 %	52,66	cries for help	2	0.42 %	7.70
in the	3	0.42 %	41.79	e mail call	2	0.42 %	95.89
of opinion	3	0.42 %	69.97	experience management	_		
real time	3	0.42 %	17.01	enterprise	2 н	0.42 %	62.65
seek to	3	0.42 %	28.58	extract and analyze	2	0.42 %	96.80
sentiment analysis	3 L,I	0.42 %	69.52	focus on applications	2	0.42 %	97.96
sentiment extraction	3	0.42 %	37.29	from linguamatics to	2	0.42 %	81.52
the results	3	0.42 %	33,45	from sources as	2	0.42 %	96.45
triggered by	3	0.42 %	26.00	information from sources	2	0.42 %	96.52
a decision	2	0.28 %	20.41	mail call center	2	0.42 %	95.82
a new	2	0.28 %	65.21	management enterprise feedback	2 н	0.42 %	62.58
analytics can	2	0.28 %	97.15	notes and survey	2	0.42 %	95.61
analytics vendor	2	0.28 %	55.02	of opinion leadership	2	0.42 %	80.43
analyze attitudinal	2	0.28 %	96.66	online consumer forums	2	0.42 %	55.90
and analyze	2	0.28 %	96.73	postings e mail	2	0.42 %	95.96
and other	2	0.28.0/n	37.70	roal time two	0	0.40.04	10.50



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 cooccurrence of terms.

Yet something is missing. What? (Hint: It's defined on p. 36.) 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 Corporation, 2008





When we understand, for instance, parts of speech – <subject> <verb> <object> – we're in a position to discern facts and relationships.

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













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.



Data integration via 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>

Alta Plana

©Alta Plana Corporation, 2008

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>
```

</GateDocument>

©Alta Plana Corporation, 2008

Other Integration Models

Application integration works in some instances.



Integration models

Another model of application/component integration.



Example: E-mail

For Text-BI/DW, we're most interested in IE. 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
----------	--------

	History B	ookmarks <u>T</u> ools <u>H</u> elp del <u>.</u> icio.us						
• 🧼 •	ی چ	🏠 📑 🔒 🗋 http://www.calepa.ca	a.gov/Customer	/CSForm.a	isp 🔻 🖡	> G • I		
		Who was the service provider?						
		Board, Department, Select Board, De	partment. or (Office			-	
		or Office:	t with up?				_	
		C Ganaral Information C Brahlam	Pocolution	От	chnical Acc	ictoreo		
		C Permitting/Licensing Assistance	C Other	- 16	chinical Ass			
			o unon j					
		Check as Appropriate						
		Statements	Agree	Agree	Disagree	Disagree	Comment	
		Staff was courteous and helpful.	0	0	0	C	0	
		Staff provided complete, accurate information to you.	0	0	0	C	0	
		A timely response was provided.	•	0	0	0	0	
		My overall experience was positive.	0	0	0	С	0	
		Please complete the section below permitting/licensing/registration as	r if your cont sistance.	tact with	us involve	d		
		The regulations were understandable.	0	0	0	C	0	
		The application instructions were understandable.	0	0	0	0	0	
		The terms and conditions of the permit, license, or registration were understandable.	0	0	0	0	0	
		Comments:	rour service	expecta	tions, pleas	se describe t	the situation,	
		As a result of your experience with	nvolved and us, what ser	i the dat	e the incide ated improv	ent occurred	: you	



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. C C C Please complete the section below if your contact with us involved permitting/licensing/registration assistance. The regulations were understandable. O C C C C \mathbf{O} \mathbf{O} O. 0 \mathbf{C} The application instructions were understandable. The terms and conditions of the Ō. C O. O. C. 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:

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?

Example: Attitudinal Data

Exercise: Identify the attitudinal information in this excerpt from Dell's IdeaStorm.com –

"Dell really... REALLY need to stop overcharging... and when i say overcharing... i mean atleast double what you would pay to pick up the ram yourself."

What Sentiment is expressed?

Subject?

Polarity?

Intensity?

Opinion?

Applications

Text Analytics is applied in many domains – Life sciences.

- Intelligence and law enforcement.
- E-discovery (legal) and compliance.
- Publishing and information services.
- Insurance and financial services.
- Voice of the Customer (sales, marketing, and product) applications.


Voice of the Customer

You see the value of the voice of the customer, the need to understand –

- the totality of stakeholder needs and opinions, whether explicitly stated or indirectly implied.
- individual views and collective, market thinking.
- customers wherever and however they express themselves.
- Our thesis is that customer voices are most frequently expressed in text.

Voice of the Customer

I. Quantitative



www.andersonanalytics.com/index.php?mact=News,cntnt01,getfile,1&cntnt01filename=SCIP0208TomAndersonArticle.pdf&cntnt01returnid=46&page=46

Voice of the Customer

Additional concepts and tools apply...

"Net Promoter is a discipline by which companies profitably grow by focusing on their customers."

"One simple question - **Would you recommend us to a friend or colleague?** - allows companies to track promoters and detractors and produces a clear measure of an organization's performance through its customers' eyes."

-- http://www.netpromoter.com/netpromoter/index.php



Applications

Take law enforcement as an example-

- Sources: case files, crime reports, incident and victimization databases, legal documents
- Targets: crime patterns, criminal investigation, networks



An Attensity lawenforcement example – NLP to identify roles and relationships.

Alta Plana



Applications



Alta Plana

©Alta Plana Corporation, 2008

Applications

- Customer Relationship Management (CRM) Sources: customer e-mail, letters, call centers
 - Targets: product and service quality issues, product management, contact routing and CRM automation

Finance and compliance

- Sources: financial & news reports, corporate filings & documents, trading records
- Targets: insider trading, reporting irregularities, money laundering and illegal transactions, pricing anomalies

Health Care Case Management

- Sources: clinical research databases, patient records, insurance filings, regulations
- Targets: enhance diagnosis and treatment, promote quality of service, increase utilization, control costs

Intelligence and counter-terrorism

- Sources: news and investigative reports, communications intercepts, documents
- Targets: organization associations and networks, behavioral/attack patterns, strategy development

Case study: IBM's MedTAKMI

MEDLINE from the National Center for Biotechnology Information hosts links to many widely used information sources such as the PubMed database of 15 million biomedical journal abstracts. Visit www.ncbi.nlm.nih.gov.

Alta Plana



Case study: IBM's MedTAKMI

MedTAKMI = Text Analysis and Knowledge MIning for Biomedical Documents (*ibm.com*):

- An extension of the TAKMI (Text Analysis and Knowledge MIning) system originally developed for text mining in CRM applications.
- Goal is to extract relationships among biomedical entities (e.g. proteins and genes), from patterns such as "A inhibits B" and "A activates B," where A and B represent specific entities.
- Work starts with a "syntactic parser" that identifies entities and basic binary (a noun and a verb) and ternary (two nouns and a verb) relationships.

Alta Plana

Text Analytics for BI/DW Practitioners

Case study: IBM's MedTAKMI

Figure 1 MedTAKMI architecture





Case study: IBM's MedTAKMI

2D Map								X
							Save Copy Pr	int
cross category Mus musculus		protein	tumor suppressor	breast cancer	Bc12-associated X protein	myelocytomatosis oncogene	transcription factor	
Cross Ade: By Frequecy	protein	632 (100.0%)	57 (9.01%)	42 (6.64%)	60 (9.49%)	24 (3.79%)	27 (4.27%)	
C By Alphabet	tumor suppressor	57 (38.51%)	148 (100.0%)	4 (2.7%)	8 (5.4%)	6 (4.05%)	10 (6.75%)	1
Amino Acid	breast cancer	42 (35.9%)	4 (3.41%)	117 (100.0%)	2 (1.7%)	8(6.83%)	3 (2.56%)	1
+Dry Lab Methods	Bcl2-associated X-prot	60 (59.41%)	8 (7.92%)	2 (1.98%)	101 (100.0%)	4 (3.96%)	3 (2.97%)	
+cellular component	myelocytomatosis onco	24 (38.1%)	6 (9.52%)	8 (12.7%)	4 (6.34%)	63 (100.0%)	2 (3.17%)	
+Root of LocusLink phenotype	transcription factor	27 (45.0%)	10 (16.66%)	3 (5.0%)	3 (5.0%)	2 (3.33%)	60 (100.0%)	
+MeSH Major (Tree) Minor MeSH	G elongation factor	28 (53.85%)	3 (5.76%)	1 (1.92%)	19 (36.54%)	1 (1.92%)	1 (1.92%)	1
Major MeSH +Protein By Species	proliferative cell nuclea	19 (37.25%)	1 (1.96%)	3 (5.88%)	2 (3.92%)	3 (5.88%)	1 (1.96%)	1
Age Tags Country	period	18 (37.5%)	2 (4.16%)	4 (8.33%)	2 (4.16%)	3 (6.25%)	0 (0.0%)	
Data Completed Publication Data	enhancer of rudimentar	17 (42.5%)	2 (5.0%)	23 (57.5%)	0 (0.0%)	1 (2.5%)	3 (7.5%)	
URL Full-Text Last Revision Date	epiregulin	17 (42.5%)	2 (5.0%)	23 (57.5%)	0 (0.0%)	1 (2.5%)	3 (7.5%)	
Organization Name Prime Name of Substance	progesterone receptor	13 (32.5%)	0 (0.0%)	20 (50.0%)	0 (0.0%)	1 (2.5%)	0 (0.0%)	1
Personal Name as Subject Publication Type	Harvey rat sarcoma viru	12 (32.42%)	5 (13.51%)	0 (0.0%)	1 (2.7%)	9 (24.32%)	1 (2.7%)	1
Subset Secondary Source Identifier	tumor necrosis factor	11 (33.33%)	3 (9.09%)	0 (0.0%)	4 (12.12%)	1 (3.03%)	2 (6.06%)	
Source Summary For Patient In	epidermal growth factor	11 (33.33%)	2 (6.06%)	6 (18.18%)	0 (0.0%)	3 (9.09%)	0 (0.0%)	
Journal Title Abbreviation	vascular endothelial gro	9 (29.03%)	0 (0.0%)	1 (3.22%)	1 (3.22%)	1 (3.22%)	1 (3.22%)	
Transliterated/Vernacular Title	vitamin D receptor	10 (33.33%)	0 (0.0%)	2 (6.66%)	1 (3.33%)	1 (3.33%)	1 (3.33%)	
Update Of	estrogen receptor	8 (28.57%)	0 (0.0%)	16 (57.14%)	0 (0.0%)	1 (3.57%)	0 (0.0%)	
+PROPERTY Minor Qualifier	protein-L-isoaspartate	7 (25.92%)	0 (0.0%)	2 (7.4%)	0 (0.0%)	3 (11.11%)	1 (3.7%)	1
Major Qualifier +SACCHARIDE	MAS1 oncogene	9 (33.33%)	0 (0.0%)	11 (40.74%)	0 (0.0%)	1 (3.7%)	0 (0.0%)	1
+SMART Domain	protein kinase	7 (25.92%)	1 (3.7%)	0 (0.0%)	4 (14.81%)	1 (3.7%)	1 (3.7%)	-
Author Data Created	cornichon	8 (30,76%)	1 (3.84%)	3 (11.54%)	0 (0.0%)	0 (0.0%)	1 (3.84%)	-



Case study: IBM's MedTAKMI

Entity extraction here is recognition of gene, protein, and chemical names from biomedical text based on a domain dictionary with two million entities.

Categories are constructed from public ontologies.

	S_TKc	STYKc	TyrKc	HATPase_c	HisKA	HR1	SAM	ITAM
Leukemia	9 (5.23%)	9 (5.23%)	9 (5.23%)	3 (1.74%)	2 (1.16%)	2 (1.16%)	2 (1.16%)	1 (0.58%
HMG-CoA lyase deficiency	7 (10.29%)	7 (10.29%)	7 (10.29%)	3 (4.41%)	2 (2.94%)	3 (4.41%)	2 (2.94%)	0 (0.0%)
Hepatic lipase deficiency	7 (10.29%)	7 (10.29%)	7 (10.29%)	3 (4.41%)	2 (2.94%)	3 (4.41%)	2 (2.94%)	0 (0.0%)
Miller-Dieker lissencephaly syndrome	2 (5.4%)	2 (5.4%)	2 (5.4%)	1 (2.7%)	1 (2.7%)	0 (0.0%)	1 (2.7%)	0 (0.0%)
Colorectal cancer	0 (0.0%)	0 (0.0%)	0 (0.0%)	2 (6.06%)	1 (3.03%)	0 (0.0%)	0 (0.0%)	1 (3.03%
Lupus erythematosus	1 (3.57%)	1 (3.57%)	1 (3.57%)	3 (10.71%)	3 (10.71%)	0 (0.0%)	0 (0.0%)	0 (0.0%)
Osteosarcoma	0 (0.0%)	0 (0.0%)	0 (0.0%)	1 (3.7%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	1 (3.7%)
Histiocytoma	0 (0.0%)	0 (0.0%)	0 (0.0%)	1 (3.7%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	1 (3.7%)
Li-Fraumeni syndrome	0 (0.0%)	0 (0.0%)	0 (0.0%)	1 (3.7%)	0 (0.0%)	0 (0.0%)	0 (0.0%)	1 (3.7%)

Figure 11 2D map analysis correlating LocusLink phenotypes (vertical axis) and signaling proteins (horizontal axis) in 1051 papers

www.research.ibm.com/journal/sj/433/uramoto.html

Alta Plana

©Alta Plana Corporation, 2008

Getting Started: Options

- Deployment Options -
 - Traditional software installation.
 - SaaS / Hosted.
 - Outsourced to a service bureau.
 - Web Services APIs.
 - VOC text analytics embedded in line-of-business applications.
- Consider evaluation and implementation Best Practices. What did my study discover?

Alta Plana

Getting Started: Analytical Styles

Unified Analytics –

- Approaches build on familiar BI and predictiveanalytics tools and approaches...
- Adding data and text mining...
- Extracting entities, facts, sentiment, etc....
- Relying on semantic integration...
- ...for true, 360° enterprise views.

Getting Started: Vendors

- Let's categorized vendors by analytic style Text-BI.
 - Text-integrated predictive analytics/data mining.Focused on Voice of the Customer/CustomerExperience Management.Toolkits.

Research Study

The study centered on lessons solicited from: individuals with experience applying Voice of the Customer (VoC) text analytics to real-world business problems at their organizations.

a number of vendor representatives and industry analysts.

It looked at Best Practices defined as –

generalized principles, techniques, and methodologies derived from theory, academic and industrial research, direct experience, and customer

stories. Alta Plana

Research Study

For the study,

- I had free-form discussions with five industry analysts, six end users, and a variety of vendor executives.
- I additionally conducted a formal, small-sample survey of VoC text-analytics practitioners – end users and consultants – that attracted 26 valid responses.



Length of Respondent Experience

The plurality of respondents have been using text analytics for Voice of the Customer work for two years or more!

Length of Experience	Response Percent
still evaluating/not using	31%
less than 6 months	15%
6 months to less than one year	8%
one year to less than two years	12%
two years or more	35%



Satisfaction with VOC Text Analytics



There were no responses below Neutral.



©Alta Plana Corporation, 2008

Information Analyzed



Alta Plana

©Alta Plana Corporation, 2008

ROI Measured, Planned & Achieved



Alta Plana

©Alta Plana Corporation, 2008

Solution Providers

What should a prospect look for?

Response	Percent
deep sentiment/opinion extraction	80%
ability to use specialized dictionaries or taxonomies	76%
broad information extraction capability	60%
adaptation for particular sectors, e.g., hospitality, retail, health care, communications	56%
predictive-analytics integration	48%
BI (business intelligence) integration	48%
support for multiple languages	48%
ability to create custom workflows	32%
low cost	32%
hosted or "as a service" option	32%
specialized VoC analysis interface	24%

Alta Plana

©Alta Plana Corporation, 2008

Advice to End Users

Advice responses lend themselves to clustering in four categories:

- 1. Business Goals
- 2. Evaluation/Pilot
- 3. Implementation/Start-up
- 4. Operational Principles

Subsequent slides present the moreinteresting responses classified in these categories...

Business Goals

Start with the always-useful advice –

Clearly define initiatives based on strategic objectives.

- A number of items address making the case
 - Unless there is a taste for innovation in the organization, it is hard to fund a project.
 - You should understand your business. Are they ready?
 - Is your business ready to take action on the data?
 - Understand that you're going to run into opposition.
 - Try to impress on the people you're talking to that there's urgency for change.

You need an in-house team championing the effort.

Alta Plana

Evaluation/Pilot

Pilot with open source tools to learn.

- Study the TA industry. Develop detailed documentation of your needs
- Do the homework. Go through benchmarking 3-4 tools in the domain. Look at extra services: initiation, end-user training.

Do a proof of concept evaluation.

- Understand the problem you're solving and the audience you're solving it for.
- You want to look for a tool that you can use across different constituencies.

Adherence to UIMA – the Unstructured Information Management Architecture – is important.

Alta Plana

Implementation/Start-up

- The largest cluster of best-practices/advice
 responses. Some see a progressive approach –
 Think big, start small, use an incremental approach.
 Start in an area where value can be created quickly and expand from there.
 - Win fast and win often.
 - Those quick wins mean a lot. It took us a while to understand the technology.
 - Don't go after the biggest problem first.
 - Collect requirements from all potential internal customers and ensure they are included when building out reporting.

Provide simple access to data and analysis. Alta Plana Corporation, 2008 TDWI World

Implementation/Start-up

One advisor suggests looking for services while two emphasize proper training and one stresses "Have the **right tool** in place."

Two focus on staff:

Dedicate someone to the job.

I'll emphasize getting a core team together that knows the technology and can support it in the organization.

Two focus on start-up time:

Don't underestimate the time/effort needed to build libraries.

It takes time to develop the screening rules to make text analytics really useful. Alta Plana

Implementation/Start-up

One user advises –

Keep the project off executive radar until you know if it is meaningful in your business.

- While another explains
 - I didn't tell anybody what I did. I just did it.

Operational Principles

There were three (inconsistent) accuracy points – Users should demand sensitivity of analysis.
Don't expect the high levels of accuracy we experience in the quant side.
Accuracy is less critical to VOC than most uses of Text Analytics (80%+ good enough).



Operational Principles

Other points –

- Text analytics is not push button. It takes work.Be open to possibilities. VOC text analytics can be transformative in a whole variety of ways.
- Don't try to mix initiatives, look for synergies and relations but keep scope within reach
- Integrate both structured and unstructured data.
- [Forget ROI.] We need to get a Return on Time [spent on manual processes and on implementing automated text analytics].

Market View

Let's look at key attributes of the text analytics market –

- I estimate the overall software market at \$250 million/year, growing at over 25% annually. IDC and Hurwitz analysts agree.
- There's been significant consolidation, which will continue.
- There's been tremendous emergence of new solution focused companies.
- Community has been picking up steam. Open source shouldn't be far behind. Alta Plana

104

Text Analytics for BI/DW Practitioners

Market View: Vendor Consolidation



Text Analytics for BI/DW Practitioners

Market View: Community

4th Annual Text Analytics Summit













©Alta Plana Corporation, 2008

To Learn More...

Visit my Business Intelligence Network expert channel, *http://www.b-eye-network.com/sethgrimes* You can find my VOC Text Analytics report there. Get in touch:

Seth Grimes Alta Plana Corporation grimes@altaplana.com +1 301-270-0795 Thanks!