

# Blazing a Trail for the Benford's Law of Words

December 15, 2016

---

RICHARD B. LANZA, CFE, CPA, CGMA

[WWW.RICHLANZA.COM](http://WWW.RICHLANZA.COM)

[WWW.AUDITSOFTWAREVIDEOS.COM](http://WWW.AUDITSOFTWAREVIDEOS.COM)

Page 1



## Richard B. Lanza, CFE, CGMA



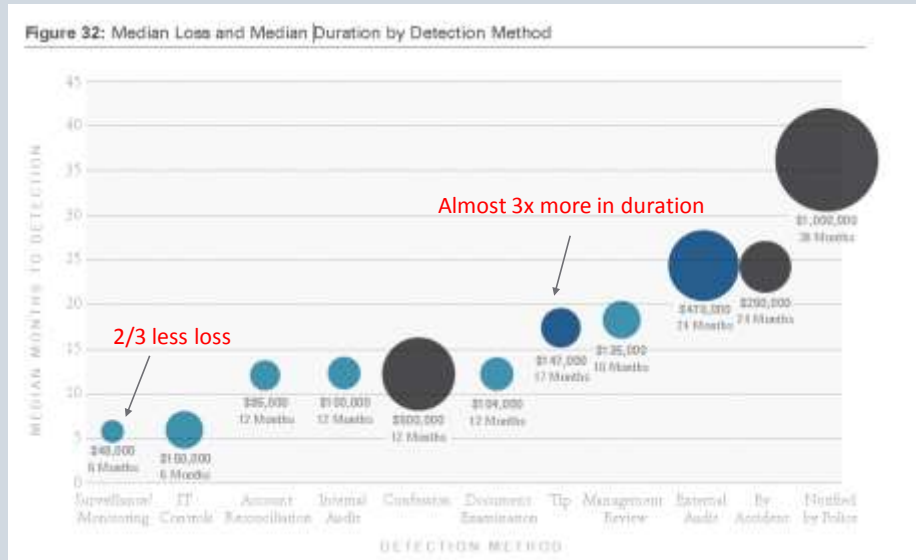
- Assists clients daily in the use of analytic software
- Nearly 25 years of ACL and excel software usage
- Has written and spoken on the use of audit data analytics for over two decades
- Received the outstanding achievement in business award by the Association of Certified Fraud Examiners for developing the publication *Proactively Detecting Fraud Using Computer Audit Reports* as a research project for the IIA
- Recently was a contributing author of:
  - Global Technology Audit Guide (GTAG #13) Fraud In An Automated World – Institute Of Internal Auditors.
  - Data Analytics – A Practical Approach - research whitepaper for the information system accountability control association.
  - Cost Recovery – Turning Your Accounts Payable Department Into A Profit Center – Wiley And Sons.
- In 2015, discovered a new textual analytic technique using letters called the Lanza Approach to Letter Analytics (LALA)<sup>TM</sup>



Please see full bio at [www.richlanza.com](http://www.richlanza.com)

2

## Surveillance is the quickest; lowers fraud impacts



2016 Report to the Nations – Association of Certified Fraud Examiners

3

## PredPol <http://www.predpol.com/>

Predictive Modeling To Improve Police Detection

Santa Cruz experienced:

- 27% decrease in burglary
- 11% decrease in robbery
- 56% increase in arrests

<http://bit.ly/1VyQPQY>



“PredPol **does not** replace the **experience and intuition of our great officers**, but is rather an invaluable added tool that allows our police force to use their patrol time more efficiently and helps stop crime before it happens.” **Chief Mark Yokoyama**

4

<http://bit.ly/1gP3meq>

## EY Global Forensic Data Analytics Survey 2014



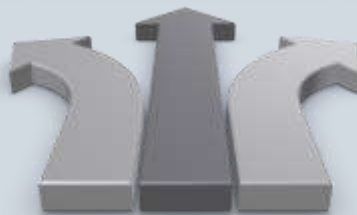
450 executives surveyed

- 72% of respondents believe that emerging big data technologies can play a key role in fraud
- **Only 7% of respondents are aware of any specific big data / Only 2% are using them**
- 12% utilize visualization / 26% apply key word searches
- 62% of respondents indicate that they need to improve management's awareness of the benefits of analytics

Page 5

## The Population of Data Type

Structured Data	Unstructured Data	External Data
<ul style="list-style-type: none"> <li>▪ Accounting records</li> <li>▪ Sub ledger details</li> <li>▪ Monthly performance measures</li> </ul>	<ul style="list-style-type: none"> <li>▪ Documents (Excel, PDF, Word)</li> <li>▪ Emails</li> <li>▪ Network Logs</li> </ul>	<ul style="list-style-type: none"> <li>▪ Geomap Service</li> <li>▪ OFAC, SAM.Gov Watch Lists</li> <li>▪ IRS Tax ID Match</li> </ul>



6

## 90% of Data is Text Based When Did You Last Investigate Text?

---

### **It works fast to quickly gain a perspective of the business process data:**

- Can work in real-time with the data while talking to the client – no prep needed...meaningful questions in seconds
- Look for deviations over a 3-year moving average to the current period

### **If digital analysis/Benford's Law is latitude, letter analytics is longitude**

- Text is far richer in business value and providing a picture than simple digit theory
- The unique word chart provides a more normalized view of activity while total word occurrences by letter provides a more dynamic chart
- The trends can be seen quickly to ask relevant questions and to highlight fraud

### **Why not use another approach, such as Benford's Law, to look at ALL data?**

7

## Red Flag (Key) Word Phrases

---

## Red Flag Word Phrases/Words

- One could build a database of suspicious words and then search the entire data file for these words, looping back to the table to get the next word:
  - bribe
  - fraud
  - plug
  - etc.
- Summaries can be done by person and collectively for additional collusion reviews

## Lessons from WorldCom/ MCI

The fraud was accomplished primarily in two ways:

1. Booking **"line costs"** (interconnection expenses with other telecommunication companies) as capital expenditures on the balance sheet instead of expenses.
2. Inflating revenues with bogus accounting entries from **"corporate unallocated revenue accounts"**.

*In 2002, a small team of internal auditors at WorldCom worked together, often at night and secretly, to investigate and reveal \$3.8 billion worth of fraud....*

Per Wikipedia – MCI Inc.



# Key Words/Phrase Survey Summary Results – Upd. 2015

## Unique Responses

- 4,320 response phrases / 2,153 unique phrases
- Average of 17 phrases per response

## Phrase Occurrences

- Unique phrases 1,424 (66%)
- 2 to 4 occurrences 574 (27%)
- 5 to 19 occurrences 144 (7%)
- 20 and Over 11 (1%)
- 2,153

Per AuditNet® Key  
Words Survey

<http://bit.ly/1XyMwch>

## Phrase Letter Length

- Average of 10, Max of 75 and Min of 2 letters

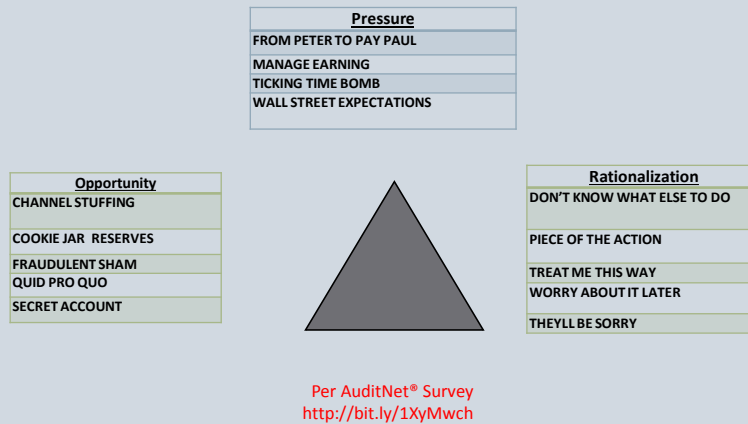
11

# Survey Phrase Summary Results Top Words - 2014

GIFT(S)	52	GREY AREA	20	KICKBACK	16
BRIBE	34	ILLEGAL	19	RESERVE	16
FRAUD	28	MISCELLANEOU			
CORRECT	26	S	35	ADJUSTMENT	15
WRITE OFF	25	PLUG	18	DISCOUNT	15
ERROR	24	WRITEOFF	18		
		CASH	17	OFF THE BOOKS	15
FACILITATION	21	REVERSE	17	PROBLEM	15
		ADJUST	16	OTHER	14
COMMISSION	20	COVER UP	16	OVERRIDE	14
ENTERTAINME				SPREAD	14
NT	20				

Page 12

## Fraud Triangle Phrase Examples



## 2015 Updates to Key Words

### AuditNet® LLC (Jim Kaplan) added:

- More Key Words
- Spam Words
- SEC and Terrorist Words
- Social Media Terms

### Stephen Valance enhanced:

- The classification of the 2014 survey key words

### Coney B.V., Amsterdam (Joris Joppe & Pieter de Kok)

- Translated the 2014 survey to Dutch

### Rich Lanza

- Finalized the list and organized all data for use

## Moving Beyond “Bad” Words to Other Word Lists

### THE CORPUS OF CONTEMPORARY AMERICAN ENGLISH (COCA)

450 MILLION WORDS, 1990-2012

The screenshot shows the 'Full-text corpus data' page from the COCA website. It features a navigation menu on the left and a main content area with a table of license options. The table lists various license types (AUP-1, AUP-2, NEW-COCA) and their corresponding prices for different user categories (Individual, Academic, Commercial).

License	Explanation	COCA	COCA-1	GENE	COCA-1000	COCA-1000-1000	COCA-1000-1000-1000
AUP-1	For use by students on college personal (professor, teacher, or student). For use by post-grad research (not allowed with colleagues or students).	\$240	\$240	\$240	\$175	\$175	\$175
AUP-2	For use by teachers on college personal (professors, teachers, students). Can be used for multiple classes in the same semester, or distributed (within a password) on a campus internal network.	\$300	\$300	\$300	\$225	\$225	\$225
NEW-COCA	Any other use*, including commercial.	\$750	\$750	\$750	\$1,275	\$1,275	\$1,275

\*To purchase the data.





# Secret Life of Pronouns LIWC

**Linguistic Inquiry and Word Count**

Home | Dictionaries | How to Use | Try Online | Contact Us

**Click here to buy LIWC**

LIWC2007  
LIWC2007

Compare 2007 and 2001 dictionaries

**What is LIWC?**

Linguistic Inquiry and Word Count (LIWC) is a text analysis software program designed by James W. Pennebaker, Roger L. Bond, and Martha E. Francis. LIWC calculates the degree to which people use different categories of words across a wide range of tasks, including journal, speeches, poems, or transcribed daily speech. With a click of a button, you can determine the degree you feel any number of negative emotions, self-references, social words, and 70 other linguistic dimensions.

The LIWC program can analyze hundreds of thousands of words (and files or Microsoft Word documents) in seconds. The LIWC2007 program also allows you to build your own dictionaries to analyze dimensions of language specifically relevant to your interests. The previous version of LIWC2007 has a feature that will highlight in color all the words found in a particular file when it is analyzed. With the Microsoft version, users can also create dictionaries that include lexical phrases (e.g., "you know") as well as individual words and word stems.

The student version of LIWC, LIWC000, only analyzes plain text files using the LIWC2007 and earlier LIWC2001 dictionaries. LIWC000 is the student version that is ideal for people with limited text analysis needs.

**LIWC license**

A single user license for LIWC2007 or LIWC000 enables you to install the software on as many as ten computers, however discounts available for multi-user licenses (see [Licensing and Academic Use](#)).

Page 19

# Secret Life of Pronouns LIWC – Dictionary Page Example

(www.liwc.net)

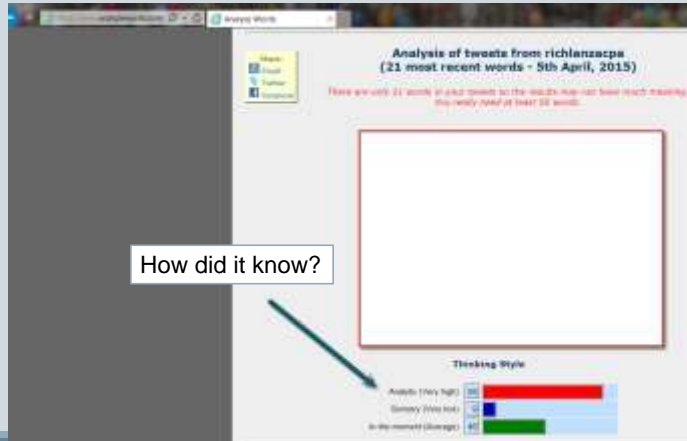
31 Posemo			32 Negemo			
hugs	promising	wealthy	);	dismay*	ignorant	p-
humor*	proud	welcom*	{	disreput*	ignore	p-
humour*	prouder	well	abandon*	diss	ignored	p-
hurra*	proudest	wellbeing	abuse*	dissatisf*	ignoring	p-
ideal*	proudly	welness	abus*	distraught	ignoring	pi
ily*	radian*	win	ache*	distress*	immoral*	pi
importance	readiness	winn*	aching*	distrust*	impatien*	pi
important	ready	wins	advers*	disturb*	impersonal	pi
importantly	reassur*	wisdom	afraid	domina*	impolite*	pi

Page 20

# Secret Life of Pronouns

## Twitter Analysis

AnalyzeWords.com – Twitter Analysis



Page 21

## Word Summarization and Review

---

## Words On The Rise / Words Equal

Page 23

Search_ON_PMI_CODE	COUNT Year 2013	COUNT Year 2014
1 ADVANCE	0	2
2 ALCOHOL	0	1
3 AUDIT	2	18
4 BALANCE	441	481
5 CANCEL	0	2
6 CASH	32	54
7 CHQ	0	1
8 CONSULTING	10	15
9 CONTRACT	4	5
10 CONTROLLER	0	3
11 CORRECT	30	12
12 CORRECTION	214	373
13 CREDIT	1	4
14 DONATIONS	1	1
15 EQUIPMENT	0	0
16 ESTIMATE	0	5
17 EXTRA	0	2
18 FEE	29	46
19 FEES	88	93
20 FSN	15	87
21 GIFT	0	3
22 OTHER	31	42
23 RECLAS	332	496
24 REFUND	12	38
25 RESERVE	84	153
26 RESERVE	24	42
27 RESERVE	5	8
28 REWARD	0	4
29 WORTH OFF	8	14

Page 24

Search_ON_PMI_CODE	COUNT Year 2013	COUNT Year 2014
4 ADJUSTMENT	308	338
54 GIFT CARD	3	2

## Summarize Words Analysis

### What You Need

- Table With Description Fields

### How You Do It

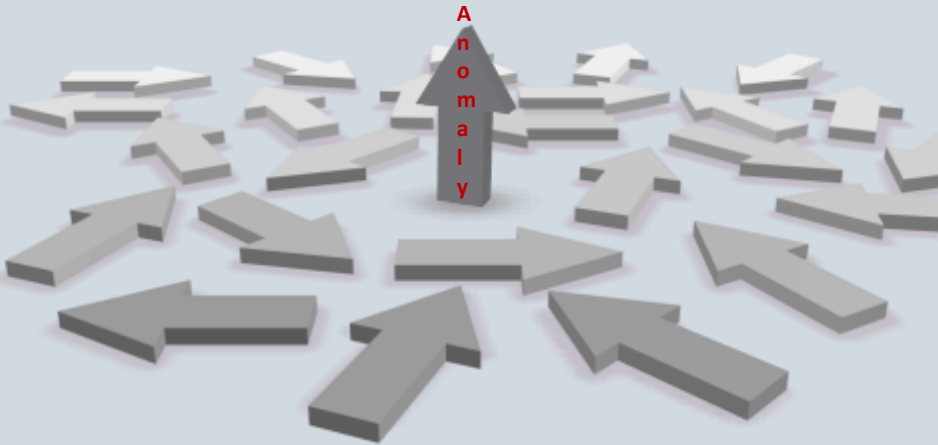
- SPLIT the field to be searched into words
  - 10 to 20 words should work well
- Extract every split field to a new table
- SUMMARIZE on the field to understand usage
  - Sum by month/year as well to trend
  - Sum by enterer





But, Isn't It About Finding the Deviations?

---



29

Can You Read This?

---

*It deosn't mttar in waht oredr the  
ltteers in a wrod are, the olny  
iprmoetnt tihng is taht the frist and  
lsat ltteer be at the rghit pclae.*

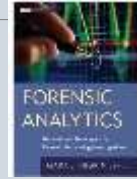
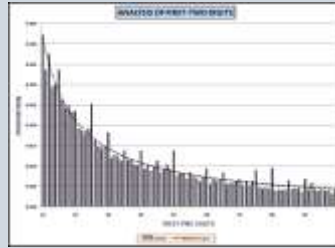
30

# Benford's Law – The Basis of Digital Analysis

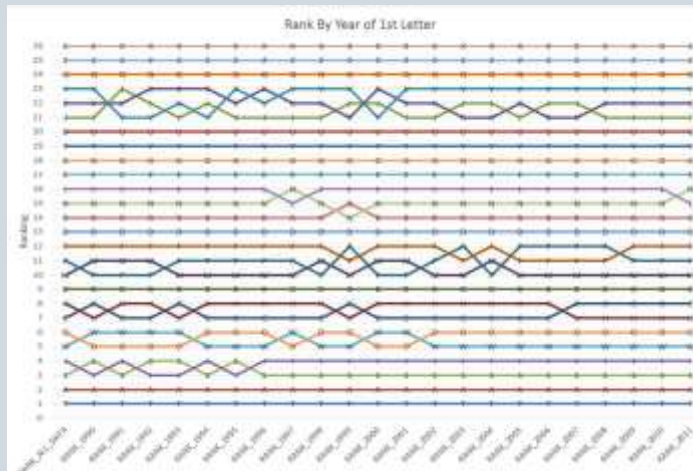
TABLE 5.1 The Expected Digit Frequencies of Benford's Law

Digit	Position Number			
	1st	2nd	3rd	4th
0		.1146	.1013	.1101
1	.2091	.1181	.1028	.1014
2	.1760	.1082	.1007	.1010
3	.1474	.1043	.1007	.1006
4	.1261	.1021	.1006	.1002
5	.1116	.1008	.1005	.1000
6	.1000	.1000	.1000	.1000
7	.0909	.1000	.1000	.1000
8	.0833	.1000	.1000	.1000
9	.0769	.1000	.1000	.1000

Source: Algori, M. J. 1996. "A taxpayer compliance application of Benford's Law." *The Journal of the American Taxation Association* 18 (Spring): 22-91.  
 The table shows that the expected proportion of numbers with a first digit 2 is 0.1760 and the expected proportion of numbers with a fourth digit 4 is 0.1002.

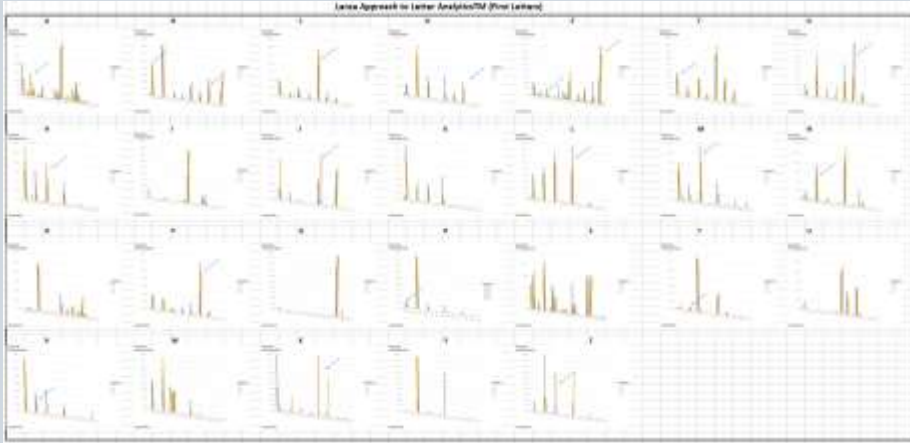


# COCA Ranking – In First Letters The Benford's Law of Words?



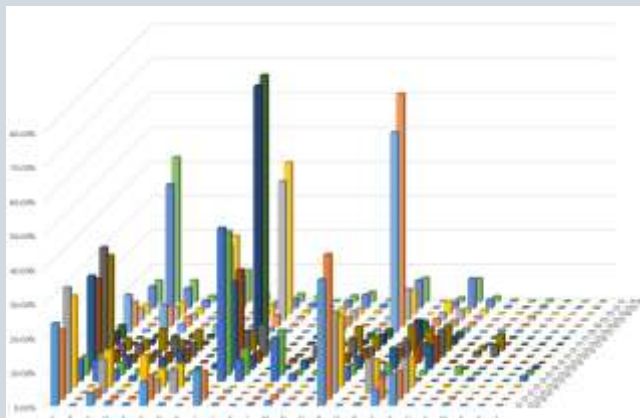


## A Benford's Law For Letters and Words?



Page 33

## General Ledger Fingerprint



## 50% From 1<sup>st</sup> and Last Letters

	A	B	(A) x (B)
# of Letters in a Word	Word Occurrences per COCA	% of Letters Analyzed in First and Last Letters	Final %
1	3.51%	100%	3.51%
2	16.02%	100%	16.02%
3	20.71%	67%	13.88%
4	17.26%	50%	8.63%
5	11.29%	40%	4.52%
6	8.53%	33%	2.81%
7	7.74%	22%	1.70%
8	5.40%	25%	1.35%
	90.46%		52.42%

35

## A Benford's Law For Words The Dashboard



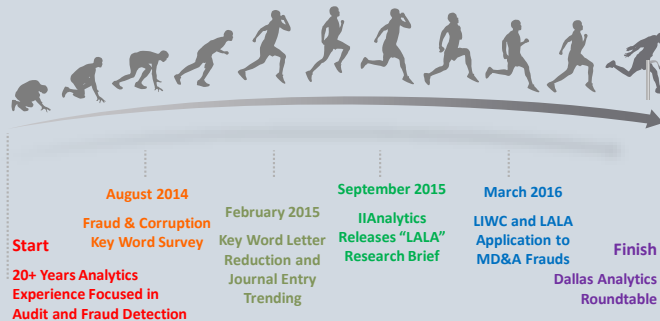
# Lanza Approach to Letter Analytics™ (“LALA”)

Identifies word deviations swiftly by relating letter frequency patterns to benchmarks of the English language and prior period letter occurrences. Focus is placed on:

- ✓ First letter (26 letters)
- ✓ Last letter (26 letters)
- ✓ First two letters (702 letters)
- ✓ Last two letters (702 letters)



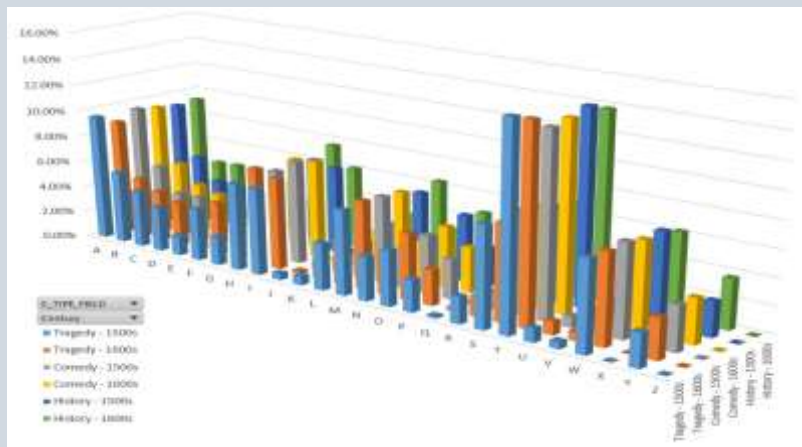
# Letter Analytics Lifecycle 2014 to Present



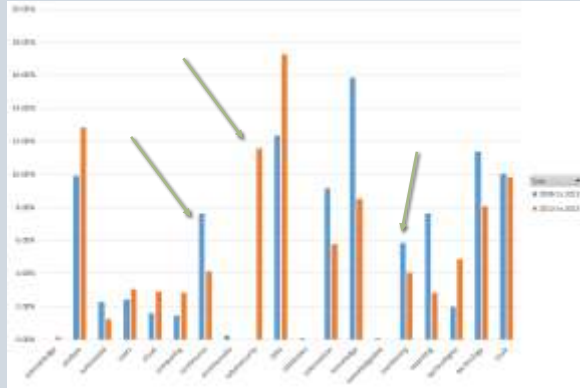
## British Pop Songs - Britburn



## Shakespeare's Plays By Century



## Protiviti Reports on Internal Auditing 2006 to 2012 vs. 2013 to 2015 – Key Words



<http://www.protiviti.com/IASurvey>

## LALA - Where Can It Be Used

Use Case	Specific Analytics
<ul style="list-style-type: none"> <li>Gain insight to a business process and its deviations</li> </ul>	<ul style="list-style-type: none"> <li>Learn new facts about the process through organizing and trending description fields from purchase and sales orders.</li> <li>Trend business process documentation over time for specific items (i.e. travel policy) and holistically across all documents</li> <li>Identify new areas of risk from customer feedback blogs, questionnaires, Emails and social media postings</li> <li>Analyze test result documentation in GRC description fields</li> <li>Pinpoint the common threads between safety and manufacturing shutdown reports</li> </ul>
<ul style="list-style-type: none"> <li>Assess journal entry risks and financial accounting trends</li> </ul>	<ul style="list-style-type: none"> <li>Trend word usage in journal entry names, and line descriptions to better visualize the monthly activity</li> <li>Map ledger-focused unusual key words to identify entries worthy of discussion</li> <li>Determine the letter fingerprint of the monthly journal entry titles, and their rate of change throughout the year</li> </ul>

## LALA - Where Can It Be Used

Use Case	Specific Analytics
<ul style="list-style-type: none"> <li>Profile employees for corruption and collusion</li> </ul>	<ul style="list-style-type: none"> <li>Determine network links between employees by trending the words in employee Emails</li> <li>Assess hourly payroll time descriptions to gain a new perspective of what everyone is working on through their words</li> <li>Perform key word searching of travel expense business descriptions</li> </ul>
<ul style="list-style-type: none"> <li>Pinpoint computer application issues and concerns</li> </ul>	<ul style="list-style-type: none"> <li>Trend employee web page access pages, searches terms and documented posts to social media through company networks</li> <li>Summarize file directory and file names by department</li> <li>Review error log tables over time to identify new error patterns or areas of increased exposure</li> </ul>

43

## Useful Links on LALA

- <http://bit.ly/1jFD87b> - Blog announcing the discovery of letter analytics.
- <http://bit.ly/1RZpolz> - Research Paper #1 – Focused on explaining the letter analytic concept with reference to a benchmark for the English Language and an analysis of British song titles from 1960 to 1999.
- <http://bit.ly/1QebYkL> - Research Paper #2 – Provides a more in-depth analysis of the population of text data and how letters can explain text variations over time more quickly than word summaries. Three examples are provided including Shakespeare's plays, Berkshire Hathaway shareholder reports and my personal Emails.
- <http://bit.ly/1W0CAZO> - Predictive Analytics Times article on how Word clouds analysis could improved with letter analytic visualizations
- ✓ <http://bit.ly/1TGwvPS> and <http://bit.ly/21mEbsU> - ACFE Fraud Magazine articles on "The Benford's Law of Words – Parts 1 and 2"
- ✓ <http://bit.ly/28LVoLd> - A Better Way To Win At Audit Wheel of Fortune Using Letter Analytics

44

# Thank You!

---

Richard B. Lanza, CFE, CGMA  
Cash Recovery Partners, LLC  
Phone: 973-729-3944  
Email: rich@richlanza.com  
[www.AuditSoftwarePros.com](http://www.AuditSoftwarePros.com)