

Practical Learnings from Real-World Data Mining

Dr. Shailesh Kumar

Google, Inc.

In theory, theory and practice are same.

In practice, they are not.

-- Lawrence Peter

Berra

Welcome to the Information Age drowning in data and starving for Knowledge

ATATTAGGTTTTTACCTACC
CAGGAAAAGCCAACCAACC
TCGATCTCTTGATAGATCTG
TTCTCTAAACGAACCTTAAA
ATCTGTGTAGCTGTCGCTC
GGCTGCATGCCTAGTGCA
CCTACGCAGTATAACAAT
AATAAATTTTACTGTCGTTG
ACAAGAAACGAGTAACTCG
TCCCTCTTCTGCAGACTGC
TTATTACGCGACCGTAAGC
TAC



This data explosion is enabled by...

- ▶ **Better “Sensors”** – Higher Resolution, More Spectral Bands, Quick Experimental Turnaround, Crowd Sourcing...
 - ▶ **Higher Bandwidth Communication** – Faster Networks and Routers, Better Compression technologies...
 - ▶ **Larger Warehouses** – Cheaper Storage, Multi-Level Caching, Scalable Database/Data warehousing technologies...
 - ▶ **Massive Crunching Power** – Faster Multi-core processors, Parallel Distributed Computing, MapReduce paradigms...
-
- ▶ **Advances in Machine Learning and Data Mining** –

Lessons from Real-World Data Mining

- ▶ **Retail: Discovering Needles in a Haystack?**
 - ▶ **Is this the right needle to look for?**
 - ▶ **Do you understand the nature of your haystack?**
- ▶ **Text: Its all in the “Representation”?**
 - ▶ What is a term?
 - ▶ What does it mean?
 - ▶ How much does it matter?
- ▶ **Computer Vision: Bridging the “Semantic Gap”**
 - ▶ Build semantically deeper features that “mean something”!



The Traditional Market Basket Analysis

Wrong needle in a mysterious haystack!



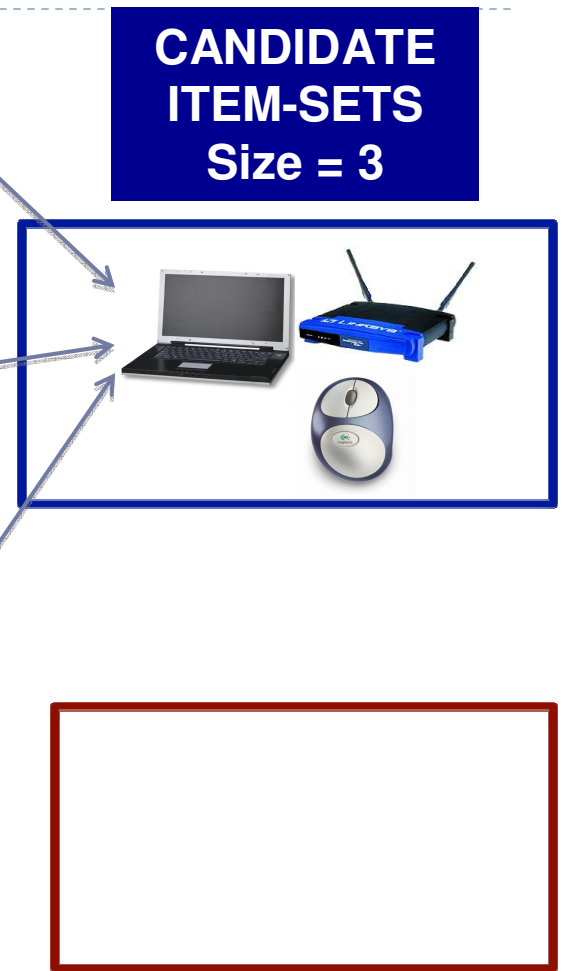
**FREQUENT
ITEM-SETS**
Size = 1



**CANDIDATE
ITEM-SETS**
Size = 2



**FREQUENT
ITEM-SETS**
Size = 2



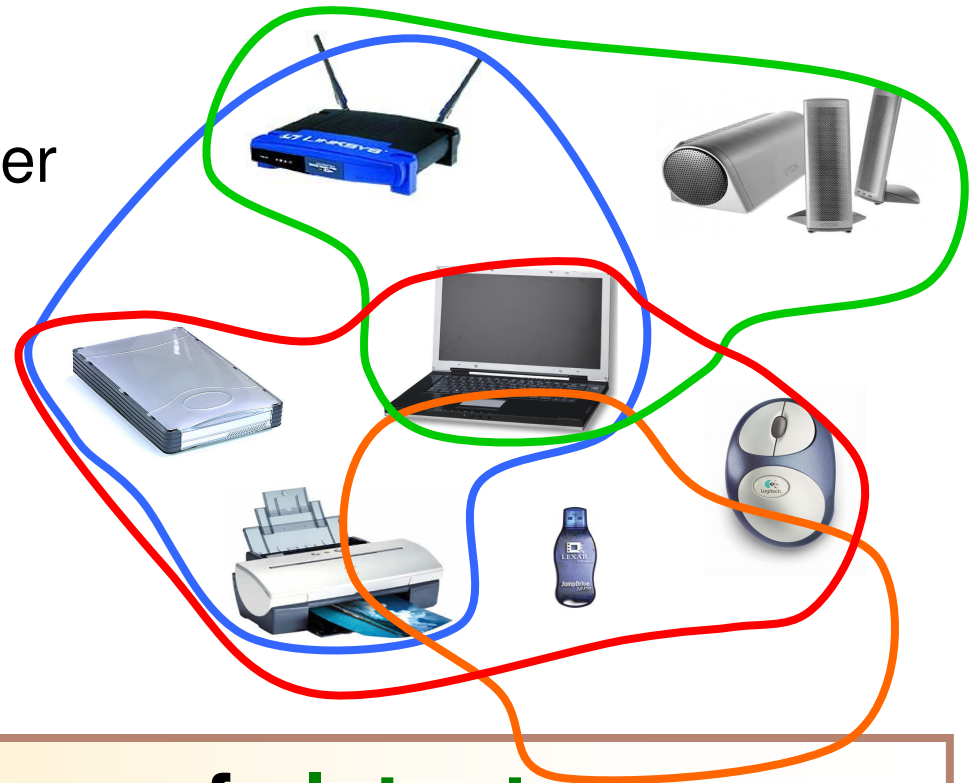
**FREQUENT
ITEM-SETS**
Size = 3

Lesson: **Know your data (Haystack)**

What process generated the data?

Few buy a complete “logical” product group in the same basket

- ❑ already have other products
- ❑ buy them from another retailer
- ❑ buy them at a different time
- ❑ got them as gifts
- ❑

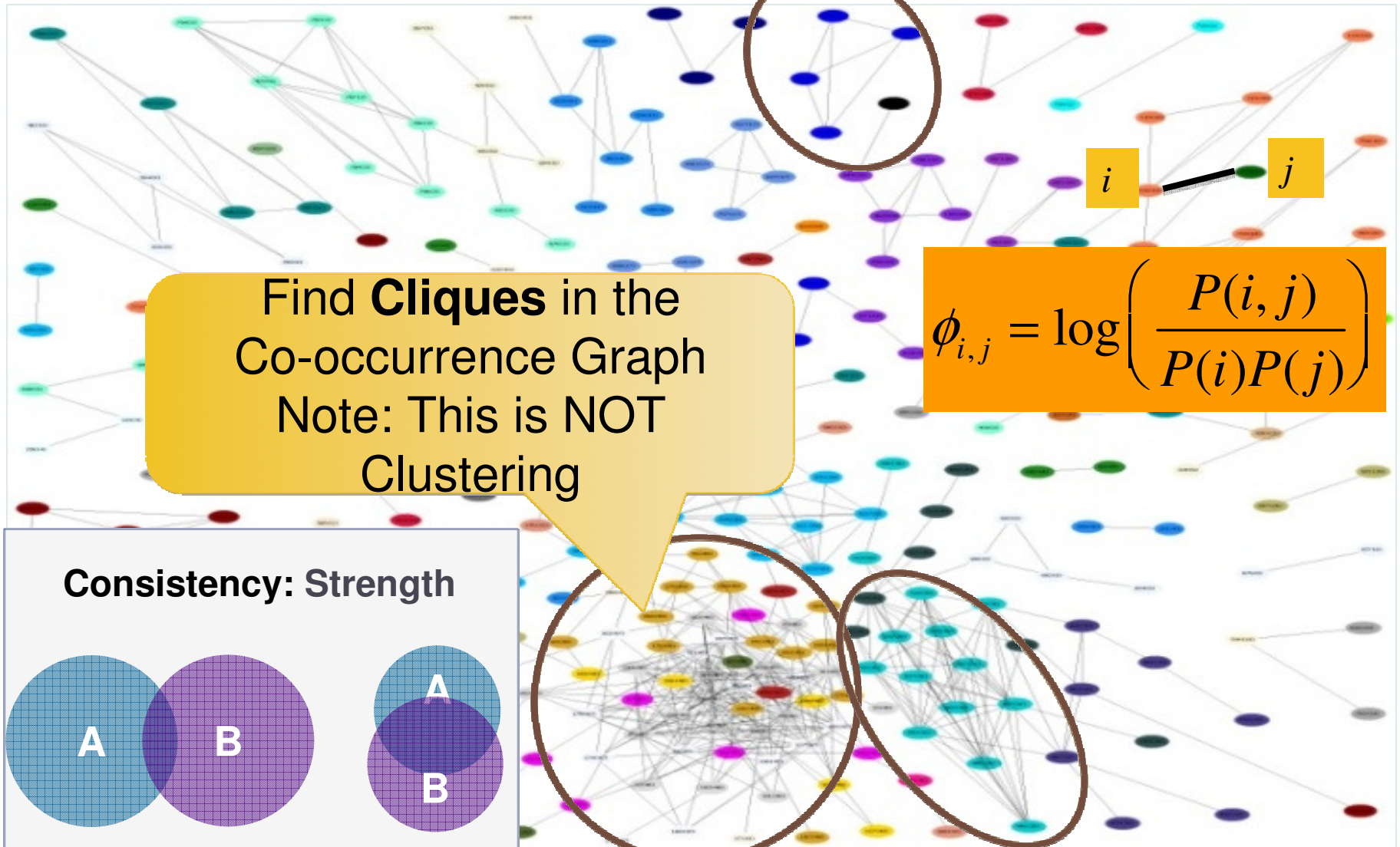


mixture of, **projections** of, **latent**

▶ **intentions**

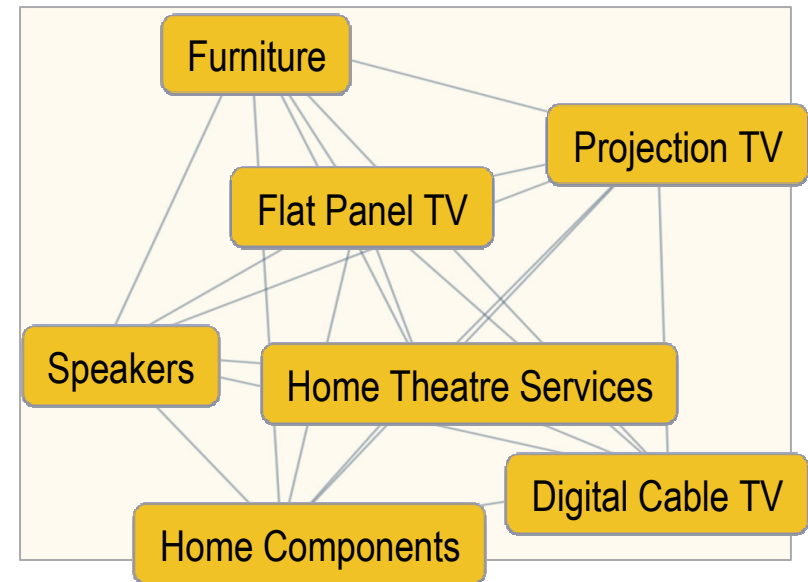
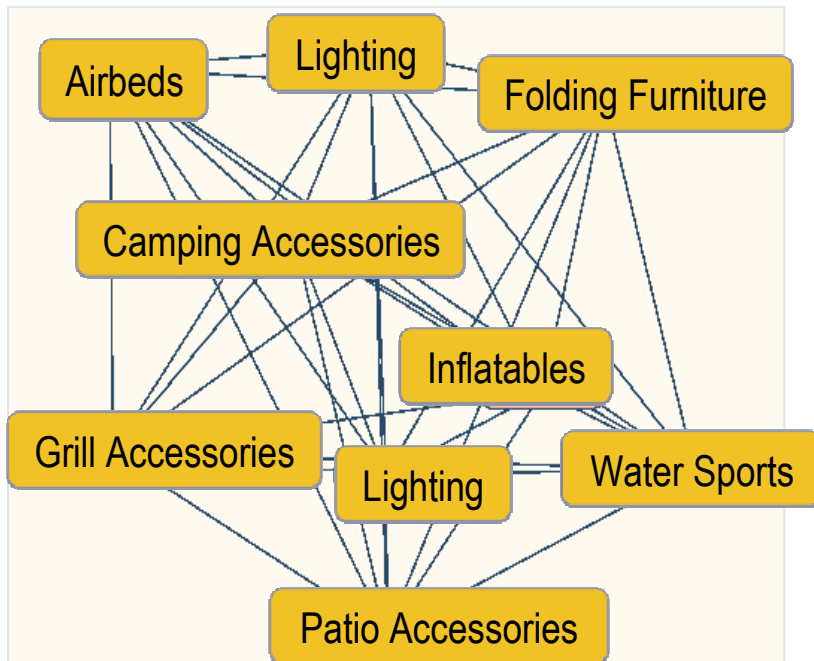
LESSON. EXTRACT THE ESSENCE, LET GO OF data

Pair-wise Co-occurrence Statistics



Lesson: **Look for the right Insight**

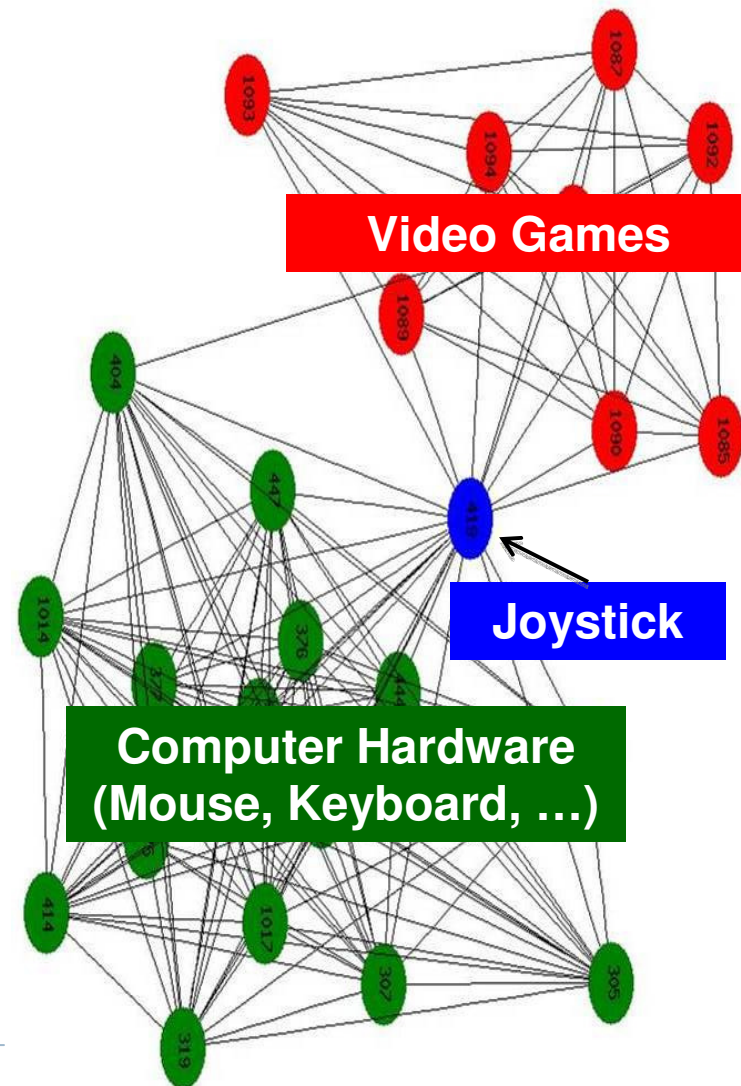
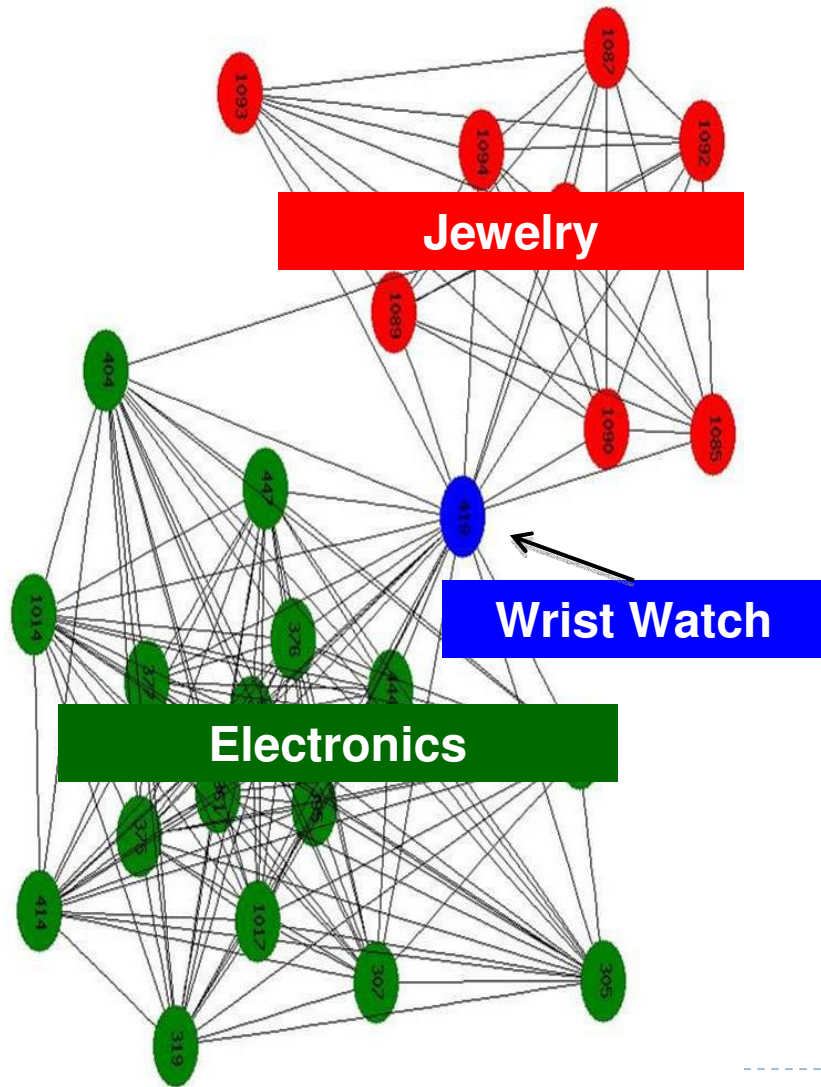
“Frequent” vs. “Logical” Itemset



- ▶ Novel – Not obvious from the data (support = 0)
 - ▶ Useful – product bundling, recommendations, layout
 - ▶ Exhaustive – “No insight left behind!” – however
 - ▶ “rare”
-

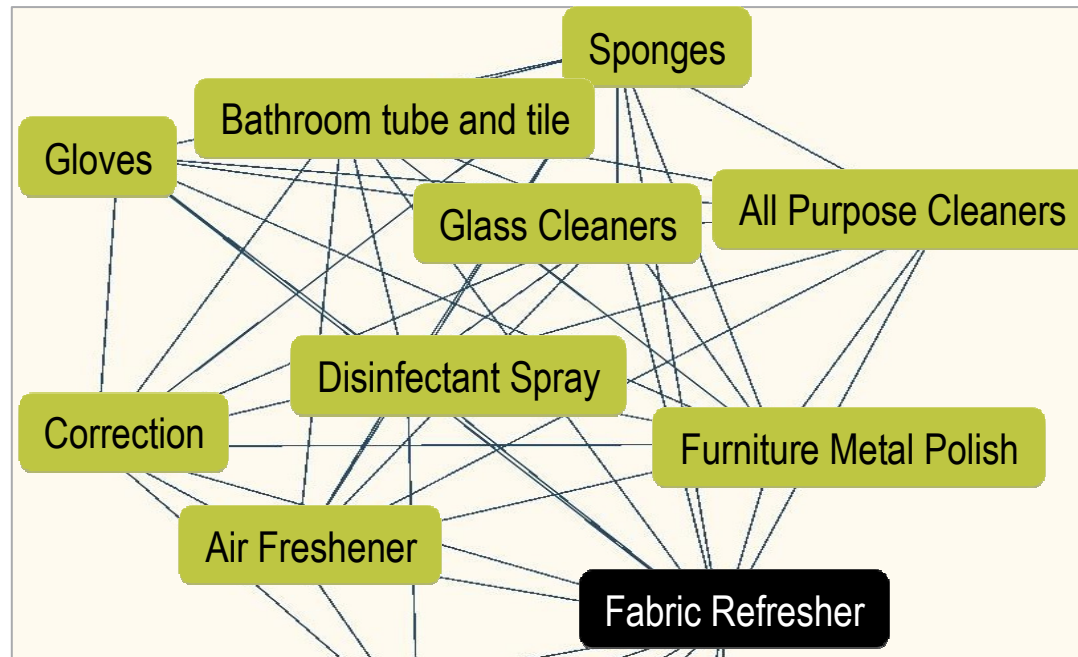
Lesson: **Let the data speak for itself!**

Bridge Products – drive traffic across depts.



Lesson: **Serendipity is good!**

I don't know what I don't know – let data help



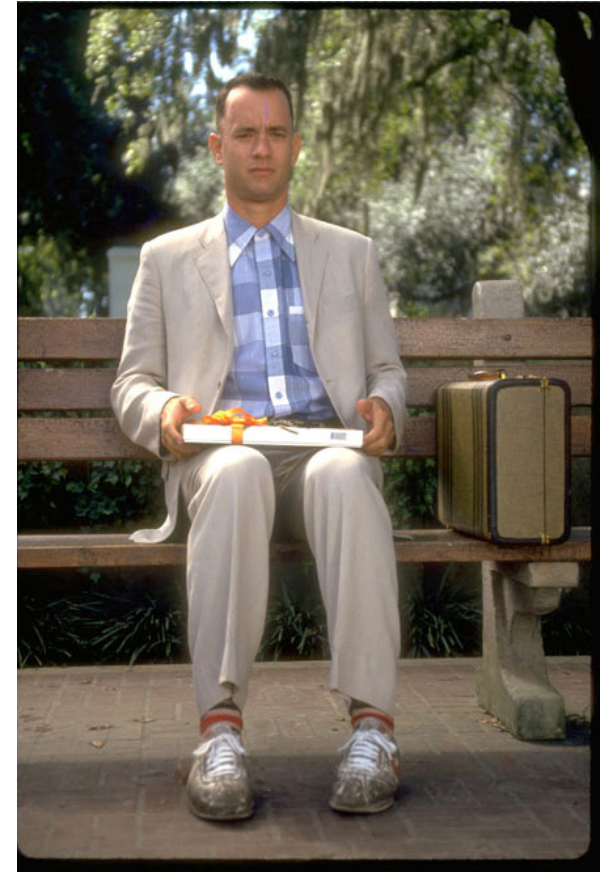
What Products Does

Fabric Refresher

Bridge To?

The Spirit of “Insight Discovery”!

“**Data Mining** is like a box of chocolates. You never know what (**Insights**) you're gonna get.”!!!



Lessons from Real-World Data Mining

- ▶ Retail: Discovering Needles in a Haystack?
 - ▶ Is this the right needle to look for?
 - ▶ Do you understand the nature of your haystack?
- ▶ **Text: Its all in the “Representation”?**
 - ▶ **What is a term?**
 - ▶ **What does it mean?**
 - ▶ **How much does it matter?**
- ▶ Computer Vision: Bridging the “Semantic Gap”
 - ▶ Build semantically deeper features that “mean something”!



Lesson: **Things are not what they appear**

What is a word in “Bag-of-Words”?

▶ **Segmentation:** What is a term?

- ▶ New York Stock Exchange → 4 words?
- ▶ “New York” “Stock Exchange” → 2 phrases?
- ▶ “New York Stock Exchange” → 1 phrase?

▶ **Disambiguation & Equivalencing:** What does it mean?

- ▶ ‘**rock** band’, ‘**rock** climbing’, ‘**rocking** chair’, ‘the **rock**’
- ▶ **orange** juice, **orange** flag, **orange** blog,
- ▶ **apple** store, **apple** pie, the big **apple**

▶ **Weighting**

- ▶ lightening,, thunder, rain, **road**, **umbrella**, **chocolate**
-



Lesson: **Supervision is best delayed**

Unsupervised, Statistical > NLP

High Cohesiveness Low Frequency

russian president boris yeltsin
prime minister viktor chernomyrdin
prime minister ryutaro hashimoto
palestinian leader yasser arafat
serbian president slobodan milosevic
syrian president hafez al-assad
british foreign secretary douglas hurd
iraqi deputy prime minister tariq aziz
secretary general kofi annan
senate majority leader bob dole
lieutenant general raoul cedras
nba commissioner david stern
bulls coach phil jackson
federal reserve policy makers

Medium Cohesiveness Medium Frequency

new york stock exchange
dow jones industrial average
nasdaq composite index
new york mercantile exchange
automated teller machines
international monetary fund
u . n . officials
u . n . peacekeepers
u . n . spokesman
an international tribunal
intensive care unit
highway traffic safety administration
mitsubishi heavy industries ltd
rio de janeiro
freshly ground black pepper

Low Cohesiveness High Frequency

he would not elaborate
sometime next year
few weeks before
nearly two years ago
made no comment
gain market share
tons of cocaine
had no chance
made no comment
materials contained herein
people have been killed
as soon as possible
dollar rose as high as
threatened to pull out
forced to give up
as far as possible
as early as tomorrow
kilometers (<num> miles) north
kilos (<num> pounds)



Lesson: **Stay Language Independent**

Chinese Giga Word Corpus

世界贸易组织	World Trade Organization	10.504958
总统克林顿	President Clinton	9.011669
亚洲金融危机	Asian Financial Crisis (1997)	9.702676
巴塞罗那奥运	Barcelona Olympics (1992)	10.333340
总统布什	President Bush	7.339422
俄罗斯总统叶利钦	Russian President Yeltsin	11.698954
糖尿病	diabetes	11.583491
艾滋病毒	Human immunodeficiency virus	14.140515
络绎不绝	in endless stream	8.637969
紧锣密鼓	an intense publicity campaign	11.915520
兴致勃勃	be highly interested in	11.415492



Lesson: **Does it mean what it says?**

Disambiguation and Equivalence

Equivalencing

SIMILARITY = 0.995

- ▶ we **filed** a **suit** charging **dell** of **illegal** **behavior**
- ▶ they **submitted** a **case** accusing **apple** of **unauthorized** **conduct**

Disambiguation

SIMILARITY = 0.171

- ▶ i was **right** to avoid a **suit** against **apple**
- ▶ on my **right** was a man in a **suit** drinking **apple** juice

**You shall know a word by the company it
keeps**

-- *Firth, J. R. 1957:11*

Lesson: **What Independence assumption?**

How much does a term weigh?

- ▶ IDF Weights: Global weights irrespective of “context”
- ▶ Contextual Weights: Local weights based on other words



idf

10.2765	hinduism
8.6589	hindu
8.6259	finger
7.8524	kerala
7.3432	mother
6.7895	smile
6.6507	child
6.576	women
6.5535	point
6.4512	happy
6.0312	orange
5.2129	india
4.312	family

ctx

8.8989	child
8.8033	smile
8.338	happy
7.982	mother
6.0989	women
4.8763	family
4.208	india
2.9307	hinduism
2.8871	hindu
2.8318	orange
1.4355	kerala
0.2292	point
0	finger

Lessons from Real-World Data Mining

- ▶ Retail: Discovering Needles in a Haystack?
 - ▶ Is this the right needle to look for?
 - ▶ Do you understand the nature of your haystack?
- ▶ Text: Its all in the “Representation”?
 - ▶ What is a term?
 - ▶ What does it mean?
 - ▶ How much does it matter?
- ▶ **Computer Vision: Bridging the “Semantic Gap”**
 - ▶ **Build features that “mean something”!**



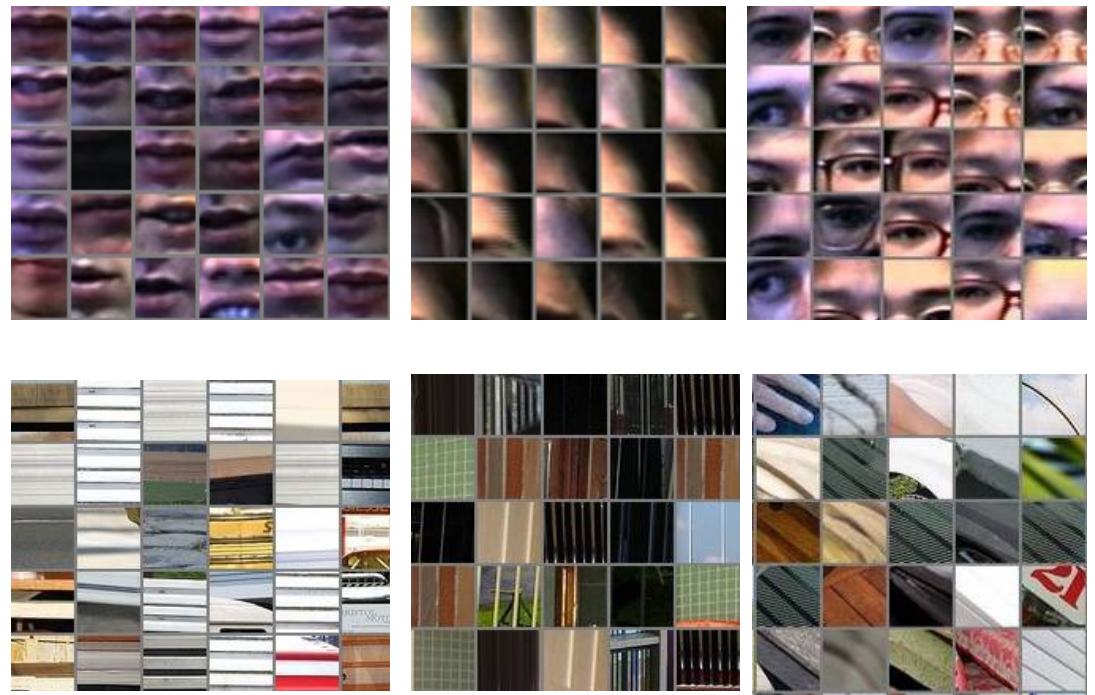
Lesson: Use Semantically deep features

Lines \rightarrow Objects

SIFT Bag-Of-Words



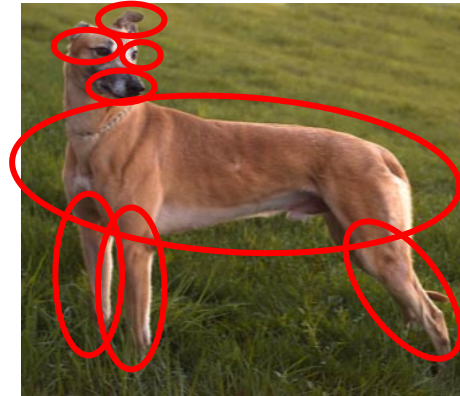
Semantically Deep features



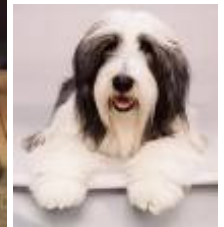
Learning a Hierarchy-of-Objects

The “Building Blocks of Learning”

- ▶ Syntactic Composition/Conjunction – “is a PART of”



- ▶ Semantic Equivalence/Disjunction – “is a TYPE of”



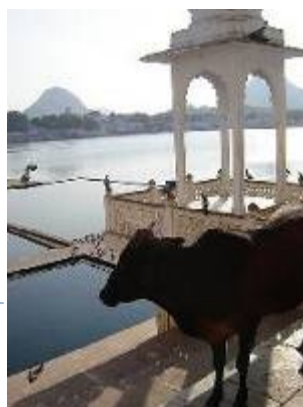
Retrieving Similar Images



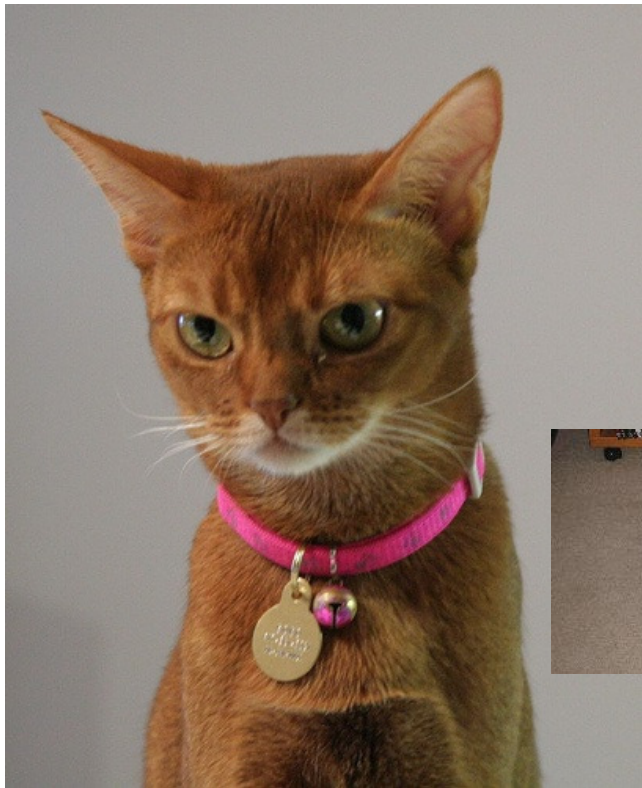
Retrieving Similar Images



Retrieving Similar Images



Content based Image Retrieval



```
ERROR: undefined
OFFENDING COMMAND: f'~
STACK:
```