# Machine Learning & Data Mining
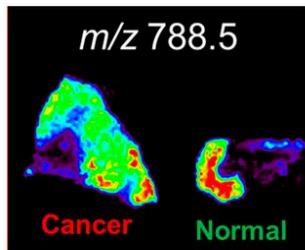## CS/CNS/EE 155

Lecture 16:

Recent Applications

# Announcements

- No lecture next week

- Final:
  - Released at noon on March 14<sup>th</sup>
    - Via Moodle
  - Due at midnight on March 15<sup>th</sup>/16<sup>th</sup>
    - Via Moodle
  - Designed to take 3 hours
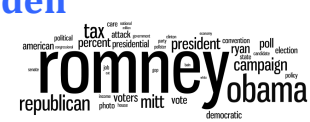
# Today: Three Recent Applications

## Cancer Detection


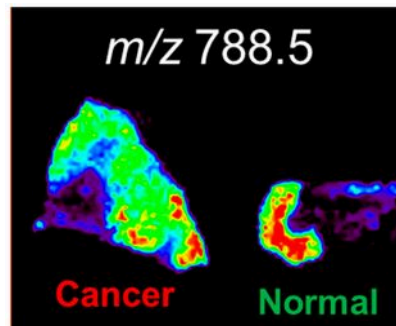
## Personalization via twitter



## Learning Visual Style



Slide material borrowed from Rob Tibshirani, Khalid El-Arini, and Julian McAuley

Image Sources: http://www.pnas.org/content/111/7/2436
https://dl.dropboxusercontent.com/u/16830382/papers/badgepaper-kdd2013.pdf
http://www.cs.cornell.edu/~andreas/iccv15.pdf
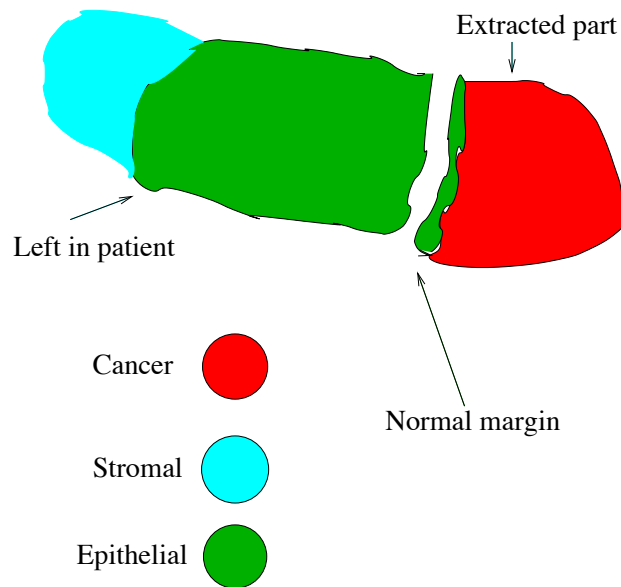
# Cancer Detection

# "Molecular assessment of surgical-resection margins of gastric cancer by mass-spectrometric imaging"

**Proceedings of the National Academy of Sciences (2014)**

Livia S. Eberlin, Robert Tibshirani, Jialing Zhang, Teri Longacre, Gerald Berry, David B. Bingham, Jeffrey Norton, Richard N. Zare, and George A. Poultsides

http://www.pnas.org/content/111/7/2436
http://statweb.stanford.edu/~tibs/ftp/canc.pdf



Extracted part

Left in patient

Cancer

Normal margin

Stromal

Epithelial

## Gastric (Stomach) Cancer

1. Surgeon removes tissue

2. Pathologist examines tissue
   – Under microscope

3. If no margin, GOTO Step 1.

Image Source: http://statweb.stanford.edu/~tibs/ftp/canc.pdf

# Drawbacks

- **Expensive:** requires a pathologist

- **Slow:** examination can take up to an hour

- **Unreliable:** 20%-30% can't predict on the spot



Extracted part
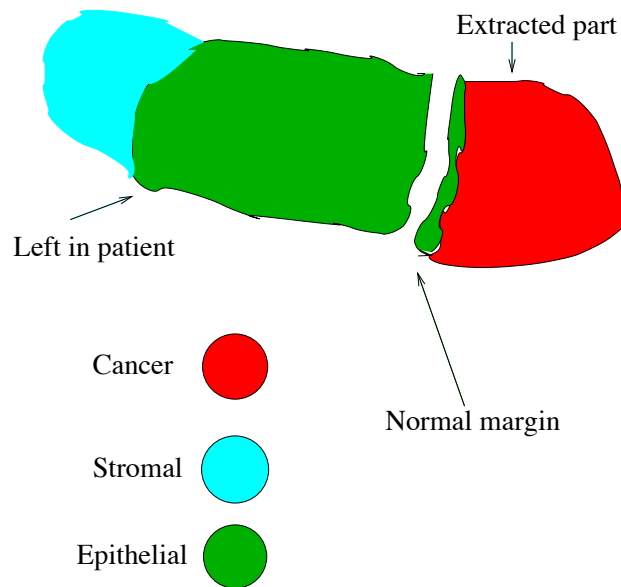
Left in patient

Cancer

Normal margin

Stromal

Epithelial

**Gastric (Stomach) Cancer**

1. Surgeon removes tissue

2. Pathologist examines tissue
   – Under microscope

3. If no margin, GOTO Step 1.

Image Source: http://statweb.stanford.edu/~tibs/ftp/canc.pdf

# Machine Learning to the Rescue!
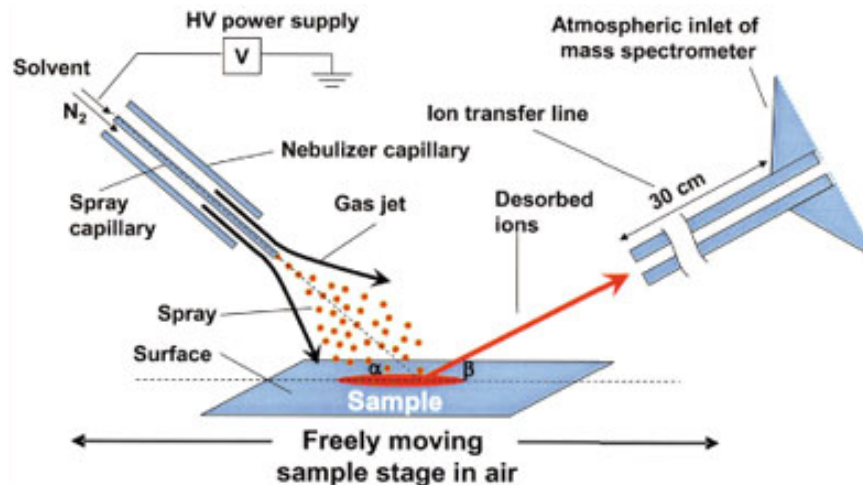## (actually just statistics)

- Lasso originated from statistics community.
  - **But we machine learners love it!**

Basic Lasso: $$\underset{w,b}{\operatorname{argmin}} \lambda |w| + \sum_{i=1}^{N} L\left(y_i, w^T x_i - b\right)^2$$

- Train a model to predict cancerous regions!
  - Y = {C,E,S}    (3 classes)
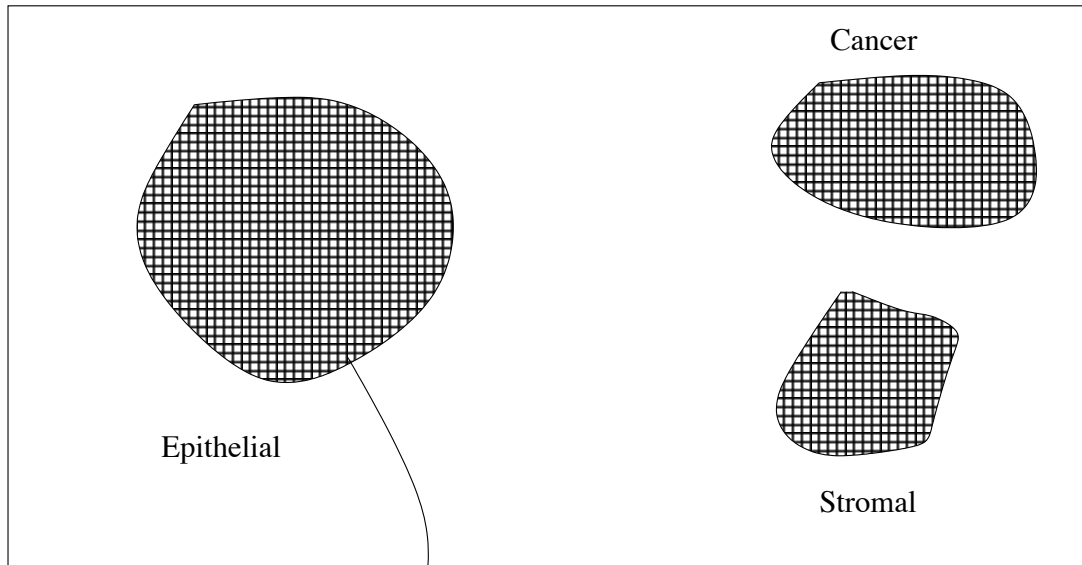  - What is X?
  - What is loss function?

# Mass Spectrometry Imaging

- DESI-MSI (Desorption Electrospray Ionization)



- Effectively runs in real-time  (used to generate x)

http://en.wikipedia.org/wiki/Desorption_electrospray_ionization
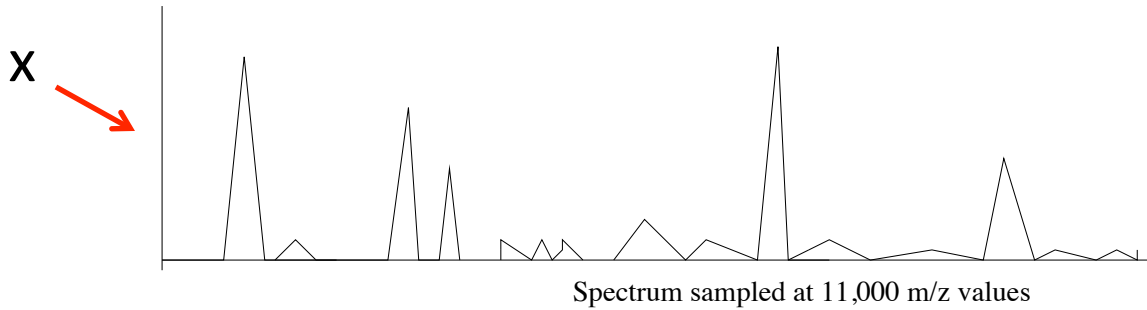
Cancer

Epithelial

Stromal

Spectrum for each pixel

Spectrum sampled at 11,000 m/z values

X

Each pixel is data point

x via spectroscopy
y via cell-type label

Image Source: http://statweb.stanford.edu/~tibs/ftp/canc.pdf

9

m/z 788.5    m/z 885.5    m/z 775.3    m/z 773.5

Cancer    Normal

0 ▭ 100%

m/z 723.5    m/z 812.5    m/z 861.5    m/z 333.5

m/z 215.3    m/z 303.5    m/z 836.5

H&E stained
Cancer    Normal stroma
Normal epithelium
1 mm

X →

Spectrum sampled at 11,000 m/z values

Each pixel has 11K features. Visualizing a few features.

# Multiclass Logistic Regression

**Binary LR:**   $P(y \mid x, w, b) = \dfrac{e^{y\left(w^T x - b\right)}}{e^{y\left(w^T x - b\right)} + e^{-y\left(w^T x - b\right)}}$     $y \in \{-1, +1\}$

**"Log Linear" Property:**   $P(y \mid x, w, b) \propto e^{y\left(w^T x - b\right)}$     $(w_1, b_1) = (-w_{-1}, -b_{-1})$

**Extension to Multiclass:**   $P(y = k \mid x, w, b) \propto e^{w_k^T x - b_k}$     Keep a $(w_k, b_k)$ for each class

**Multiclass LR:**   $P(y = k \mid x, w, b) = \dfrac{e^{w_k^T x - b_k}}{\sum_m e^{w_m^T x - b_m}}$

Referred to as Multinomial Log-Likelihood by Tibshirani

http://statweb.stanford.edu/~tibs/ftp/canc.pdf

# Lasso Multiclass Logistic Regression

$$\underset{w,b}{\operatorname{argmin}} \, \lambda |w| + \sum_i -\ln P(y_i \mid x_i, w, b)$$

$$x \in R^D$$
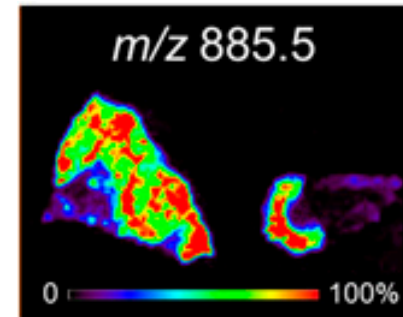$$y \in \{1, 2, \ldots, K\}$$

$$|w| = \sum_k |w_k| = \sum_k \sum_d |w_{kd}|$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_K \end{bmatrix}$$
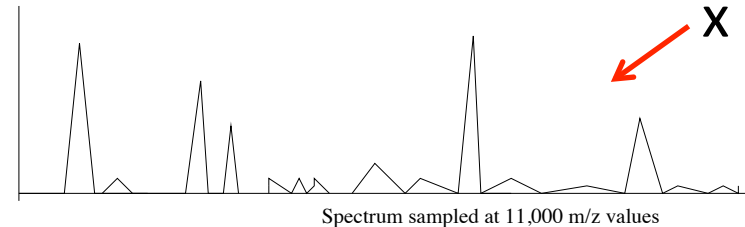
- Probabilistic model
- Sparse weights

# Back to the Problem

- Image Tissue Samples

- Each pixel is an x
  - 11K features via Mass Spec
  - Computable in real time
  - 1 prediction per pixel

- y via lab results
  - ~2 weeks turn-around



Visualization of all pixels for one feature



x

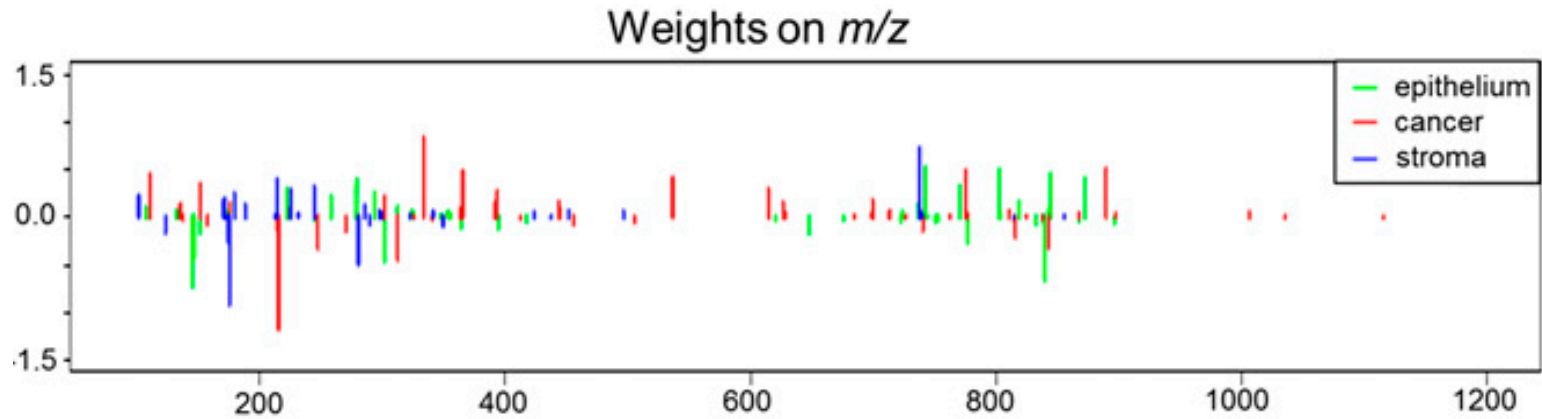Spectrum sampled at 11,000 m/z values

# Learn a Predictive Model

- Training set: 28 tissue samples from 14 patients
  - Cross validation to select λ

- Test set: 21 tissue samples from 9 patients
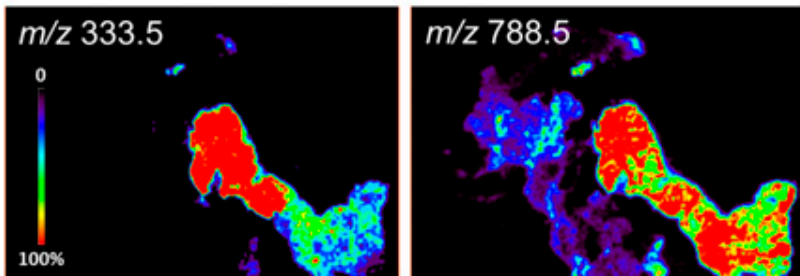
- Test Performance:

≥0.2 margin in probability

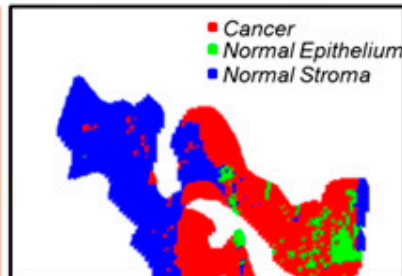| Pathology | Predicted | | | Don't know | Agreement, % | Overall agreement, % |
|---|---|---|---|---|---|---|
| | Cancer | Epithelium | Stroma | | | |
| Cancer | 5,809 | 114 | 2 | 230 | 97.0 | 97.2 |
| Epithelium | 134 | 3,566 | 118 | 122 | 96.8 | |
| Stroma | 25 | 82 | 2,630 | 143 | 96.1 | |
| | Cancer | Normal | | | Agreement, % | Overall agreement, % |
| Cancer | 5,809 | 116 | | 230 | 97.0 | 98.4 |
| Normal | 159 | 6,396 | | 265 | 99.7 | |

Weights on *m/z*

- **Lasso yields sparse weights! (Manual Inspection Feasible!)**
- Many correlated features
  - Lasso tends to focus on one


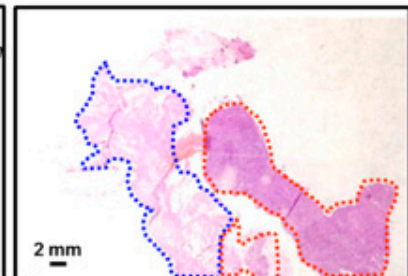A  DESI-MS Ion images — *m/z* 333.5, *m/z* 788.5
B  Lasso Prediction — Cancer, Normal Epithelium, Normal Stroma
C  Pathological Diagnosis

http://cshprotocols.cshlp.org/content/2008/5/pdb.prot4986

15

# Extension: Local Linearity

$$P(y \mid x, w, b) = \frac{e^{w_y^T x - b_y}}{\sum_m e^{w_m^T x - b_m}}$$

- Assumes probability shifts along straight line
  - Often not true

- **Approach:** cluster based on x
  - Train customized model for each cluster

| Patient | 1 | 2 | 3 | 4 | 5 | 6 | Overall |
|---|---|---|---|---|---|---|---|
| Standard training | 0.29% | 4.56% | 6.78% | 0.00% | 13.76% | 2.77% | 3.58% |
| Customized training | 0.71% | 1.89% | 0.82% | 0.40% | 9.43% | 0.92% | 1.89% |

http://statweb.stanford.edu/~tibs/ftp/canc.pdf

# Recap: Cancer Detection



- Seems Awesome!  What's the catch?
  - Small sample size
    - Tested on 9 patients
  - Machine Learning only part of the solution
    - Need infrastructure investment, etc.
    - Analyze the scientific legitimacy
  - Social/Political/Legal
    - If there is mis-prediction, who is at fault?

# Personalization via twitter



18

# "Representing Documents Through Their Readers"

## Proceedings of the ACM Conference on Knowledge Discovery and Data Mining (2013)

Khalid El-Arini, Min Xu, Emily Fox, Carlos Guestrin

https://dl.dropboxusercontent.com/u/16830382/papers/badgepaper-kdd2013.pdf

## overloaded by news

≥ 1 million news articles & blog posts generated every hour*

* [www.spinn3r.com]

# News Recommendation Engine



corpus

user

**Vector representation:**
- Bag of words
- LDA topics
- etc.

# News Recommendation Engine



corpus

user

**Vector representation:**
- Bag of words
- LDA topics
- etc.

# News Recommendation Engine



**Vector representation:**
- Bag of words
- LDA topics
- etc.

corpus

user

# Challenge

Most common representations don't naturally line up with user interests

**Fine-grained representations (bag of words)** too specific

Haqqani network is considered most ruthless branch of Afghan insurgency
Group that started as part of anti-Soviet jihad has moved into mafia-like violence, intimidation and extortion

**High-level topics (e.g., from a topic model)**
- too fuzzy and/or vague
- can be inconsistent over time

# Goal

Improve recommendation performance through a more natural document representation

# An Opportunity: News is Now Social

- In 2012, Guardian announced more readers visit site via Facebook than via Google search

# Substandard Nerd

@substandardnerd

*Gig Going, Festival Attending, Music Loving, Linux Fettling, Perl Hacking, Cycling, Vegan*

The Gdansk of Oxfordshire ·

https://www.youtube.com/user/apusskidu/featured

**badges**

---

**Substandard Nerd** @substandardnerd                    13 Jan
Stevie Nicks: the return of Fleetwood Mac
guardian.co.uk/music/2013/jan…
📄 View summary

# Approach

Learn a document representation based on how readers publicly describe themselves

# Substandard Nerd

@substandardnerd

*Gig Going, Festival Attending, Music Loving, Linux Fettling, Perl Hacking, Cycling, Vegan*

The Gdansk of Oxfordshire ·

https://www.youtube.com/user/apusskidu/featured

---

**Substandard Nerd** @substandardnerd                    13 Jan

Stevie Nicks: the return of Fleetwood Mac

guardian.co.uk

View summary

Culture 〉 Music 〉 Stevie Nicks

## Stevie Nicks: the return of Fleetwood Mac

Stevie Nicks's tumultuous life as a rock queen led her to addiction, heartbreak and "insanity". As Fleetwood Mac reunite, she tells Caspar Llewellyn Smith why she's going back for more

Using **many** tweets, can we learn that someone who identifies with

via profile badges → **music**

reads articles with these words:



?

**Given:** training set of tweeted news articles from a specific period of time

> **3 million** articles

1. Learn a **badge dictionary** from training set



music

words

badges

concert
sound
vma
rock record stage single **song**
download act west debut track kanye gaga
message listen **band**
tour mtv lady singer
performance chart pop
**album** **music** artist
billboard

2. Use badge dictionary to **encode new articles**



Haqqani network is considered most ruthless branch of Afghan insurgency
Group that started as part of anti-Soviet jihad has moved into mafia-like violence, intimidation and extortion

islam security
**afghanistan**
updated east conflict arab disabled
guardian divorced international
adult **pakistan**

# Advantages

- Interpretable
  - Clear labels
  - Correspond to user interests

- Higher

> **Haqqani network is considered most ruthless branch of Afghan insurgency**
>
> Group that started as part of anti-Soviet jihad has moved into mafia-like violence, intimidation and extortion

# Advantages

- Interpretable
  - Clear labels
  - Correspond to user interests
- Higher-level than words
- Semantically consistent over time

**politics**

**Given:** training set of tweeted news articles from a specific period of time

**3 million** articles

1. Learn a **badge dictionary** from training set

music

words

band
song
album
music
concert
sound
vma
rock record stage single
act
download message mtv
tour performance
chart
track
lady
singer
kanye gaga
pop
billboard
artist

badges

2. Use badge dictionary to **encode new articles**

Haqqani network is considered most ruthless branch of Afghan insurgency
Group that started as part of anti-Soviet jihad has moved into mafia-like violence, intimidation and extortion

afghanistan
pakistan
islam
security
east
conflict
guardian divorced
adult
arab
disabled
international
updated

# Dictionary Learning

- Training data :

$$S = \left\{ \left( z_i, y_i \right) \right\}_{i=1}^{N}$$

Identifies badges
in Twitter profile
of tweeter

Bag-of-words
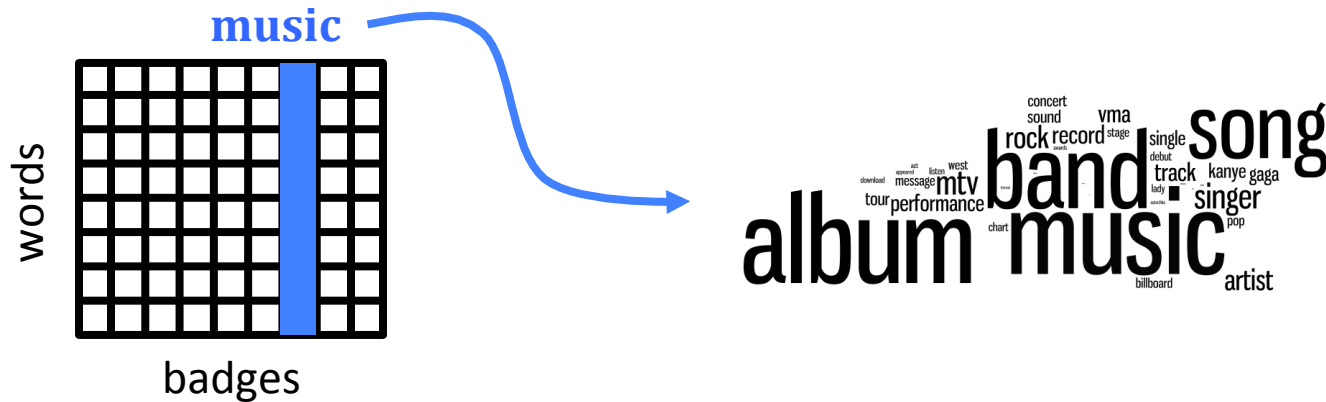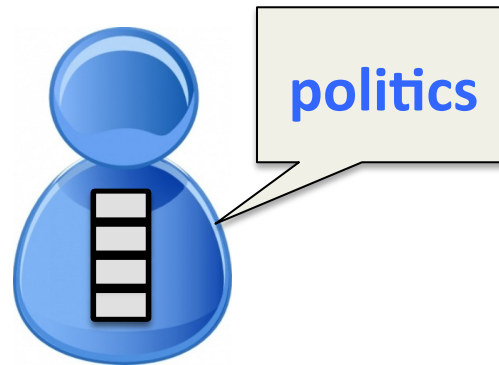representation of
document

Culture 〉 Music 〉 Stevie Nicks

## Stevie Nicks: the return of Fleetwood Mac

Stevie Nicks's tumultuous life as a rock queen led her to
addiction, heartbreak and "insanity". As Fleetwood Mac reunite,
she tells Caspar Llewellyn Smith why she's going back for more

$y$

album
Fleetwood Mac
love

Nicks

*Normalized!*

## Substandard Nerd

@substandardnerd

*Gig Going, Festival Attending, Music Loving, Linux Fettling, Perl
Hacking, Cycling, Vegan*

The Gdansk of Oxfordshire ·

https://www.youtube.com/user/apusskidu/featured

$z$

gig

music
cycling
linux

# Dictionary Learning

$$S = \left\{ \left( z_i, y_i \right) \right\}_{i=1}^{N}$$

Identifies badges
in Twitter profile
of tweeter

Bag-of-words
representation of
document

- Training Objective:

$$\operatorname*{argmin}_{B,W} \lambda_B \left| B \right| + \lambda_W \left| W \right| + \sum_{i=1}^{N} \left\| y_i - BW_i \right\|^2$$

**"Dictionary"**

words

music

badges

album music band song
record single track
singer artist

**"Encoding"**

Haqqani network is considered most
ruthless branch of Afghan insurgency
Group that started as part of anti-Soviet jihad has moved into
mafia-like violence, intimidation and extortion

badges

articles

islam security
afghanistan
guardian divorced disabled
adult pakistan

$$\underset{B,W}{\text{argmin}} \; \lambda_B |B| + \lambda_W |W| + \sum_{i=1}^{N} \|y_i - BW_i\|^2$$

**"Dictionary"**

music

words

badges

**"Encoding"**

Haqqani network is considered most ruthless branch of Afghan insurgency
Group that started as part of anti-Soviet jihad has moved into mafia-like violence, intimidation and extortion

badges

articles

- Not convex! (because of BW term)

- Convex if only optimize B or W (but not both)

- Alternating Optimization (between B and W)

- How to initialize?

Initialize:

$$W_i = \frac{z_i}{|z_i|}$$

**Use:** $S = \left\{ (z_i, y_i) \right\}_{i=}^{N}$

$z$

gig
music
cycling
linux

$$\operatorname*{argmin}_{B,W} \lambda_B |B| + \lambda_W |W| + \sum_{i=1}^{N} \|y_i - BW_i\|^2$$

- Suppose Badge s often co-occurs with Badge t
  - $B_s$ & $B_t$ are correlated

- From perspective of W, B's are features.
  - Lasso tends to focus on one correlated feature



Many articles might be about Gig, Festival & Music simultaneously.

$$\underset{B,W}{\arg\min} \, \lambda_B |B| + \lambda_W |W| + \sum_{i=1}^{N} \|y_i - BW_i\|^2$$

- Suppose Badge s often co-occurs with Badge t
  - $B_s$ & $B_t$ are correlated

- From perspective of W, B's are features.
  - Lasso tends to focus on one correlated feature

- Graph Guided Fused Lasso:

$$\underset{B,W}{\arg\min} \, \lambda_B |B| + \lambda_W |W| + \lambda_G \sum_{i=1}^{N} \underbrace{\sum_{(s,t)\in E(G)}}_{\text{Graph G of related Badges}} \omega_{st} |W_{is} - W_{it}| + \sum_{i=1}^{N} \|y_i - BW_i\|^2$$

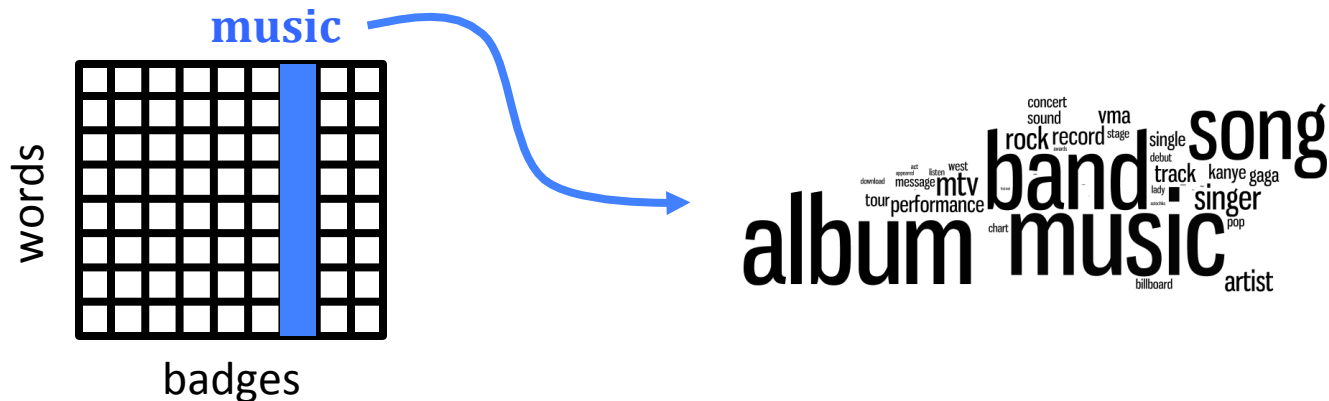Co-occurance Rate
On Twitter Profiles

# Encoding New Articles

- Badge Dictionary B is already learned

- Given a new document j with word vector $y_j$
  - Learn Badge Encoding $W_j$:

$$\operatorname*{argmin}_{W_j} \lambda_W \left| W_j \right| + \lambda_G \sum_{(s,t) \in G} \left| W_{js} - W_{jt} \right| + \left\| y_j - B W_j \right\|^2$$

# Recap: Badge Dictionary Learning

1. Learn a **badge dictionary** from training set



music

words

badges

album band music song rock record concert sound vma stage single debut track kanye gaga singer pop artist billboard west listen message mtv tour performance download act appeared chart lady

2. Use badge dictionary to **encode new articles**

Haqqani network is considered most ruthless branch of Afghan insurgency
Group that started as part of anti-Soviet jihad has moved into mafia-like violence, intimidation and extortion



afghanistan pakistan islam security adult guardian divorced east conflict arab disabled international updated
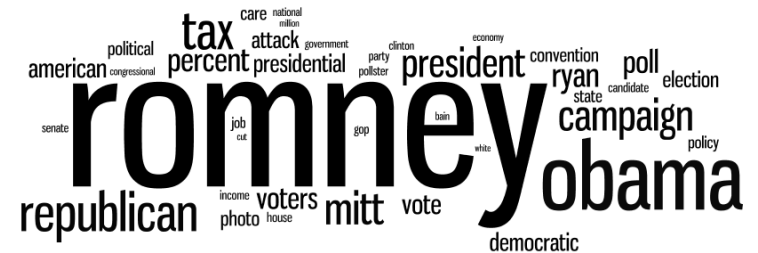
# Examining **B**

**music**

**Biden**

**soccer**

**Labour**

# Badges Over Time



September 2012

music     Biden

September 2010
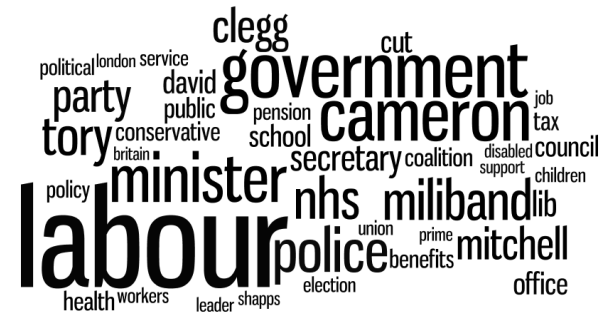
# A Spectrum of Pundits

"top conservatives on Twitter"

- Limit badges to **progressive** and **TCOT**

- Predict political alignments of likely readers?



more conservative

Dowd · Krugman · Parker · Rich · Kristof · Kristol · Klein · Friedman · Zakaria · Ignatius · Goldberg · Brooks · Krauthammer · Coulter
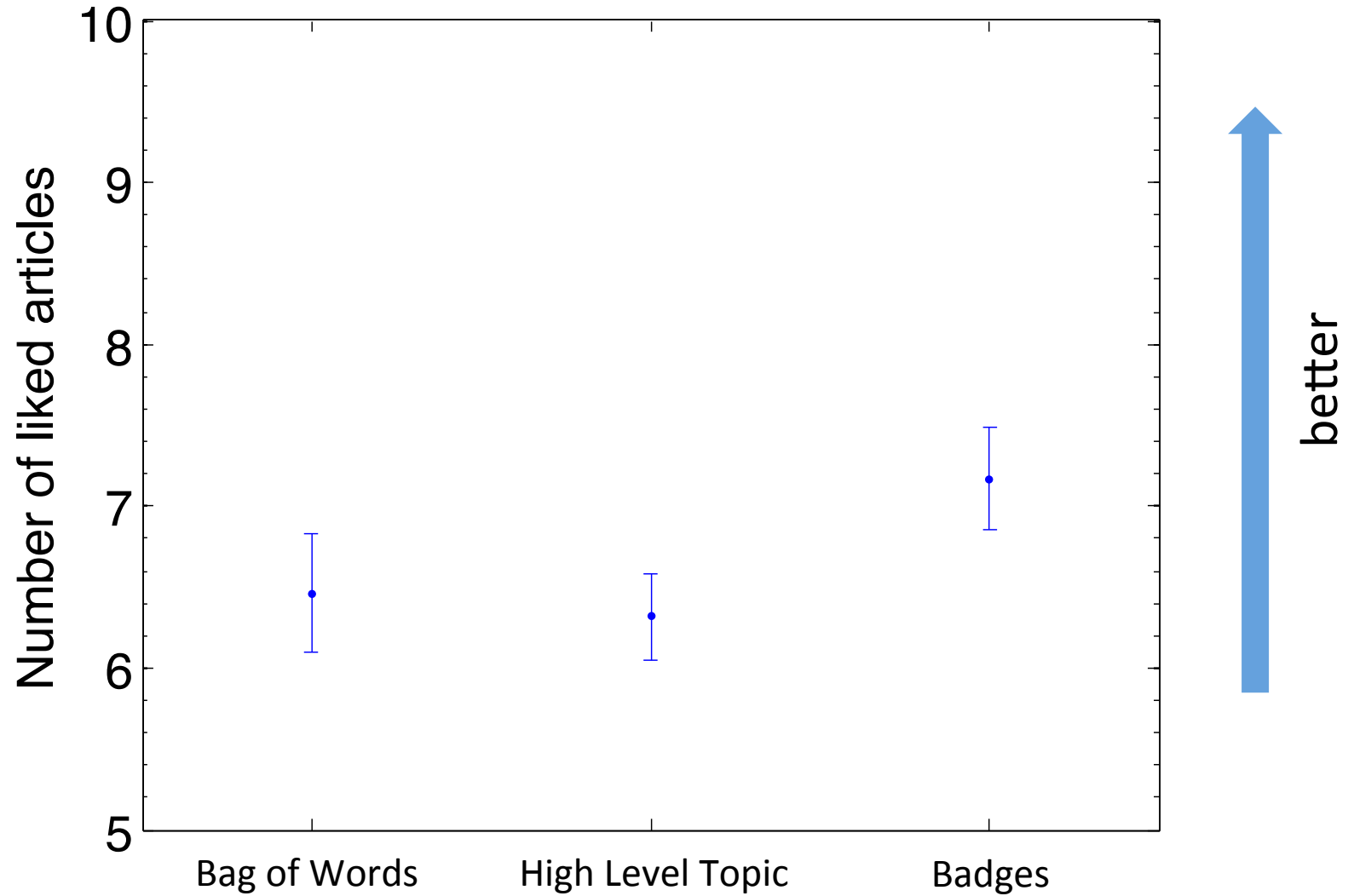
0   0.2   0.4   0.6   0.8   1

- Took all articles by columnist
- Looked at encoding score
  - progressive vs TCOT
- Average

# User Study

- Which representation best captures user preferences over time?

- Study on Amazon Mechanical Turk with 112 users

  1. Show users random 20 articles from Guardian, from time period 1, and obtain ratings

  2. Pick random representation

     - bag of words, high level topic, Badges

  3. Represent user preferences as mean of liked articles

  4. GOTO next time period

     - Recommend according to preferences

     - GOTO STEP 2

# User Study

# Recap: Personalization via twitter

- Sparse Dictionary Learning
  - Learn a new representation of articles
  - Encode articles using dictionary
  - Better than Bag of Words
  - Better than High Level Topics

- Based on social data
  - Badges on twitter profile & tweeting
  - Semantics not directly evident from text alone
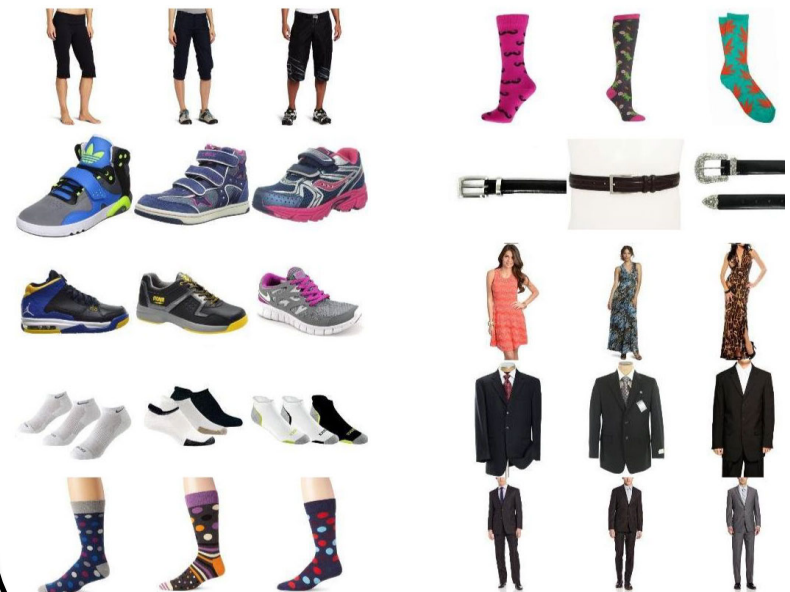
# Learning Visual Style

# Learning Visual Clothing Style with Heterogeneous Dyadic Co-occurrences

Andreas Veit, Balazs Kovacs, Sean Bell, Julian McAuley, Kavita Bala, Serge Belongie, ICCV 2015



Visually Compatible

Visually Incompatible

http://vision.cornell.edu/se3/projects/clothing-style/

# Training Data

- ## Ground set of items
  - ~1M items
  - Image of item x
  - Category of item c
    - Coat, belt, pants, socks, etc.

- ## Pairwise relationships
  - "frequently bought together"
  - Interpret as visually compatible

# Training Goal
(ignoring regularization)

Penalizes too far

Penalizes too close

Embedding of image

Embedding of image

$$\operatorname*{argmin}_{\Theta} \sum_{(i,j)\in D} L^{+}\big(\Phi(x_i),\Phi(x_j)\big) + \sum_{(i,j)\in \tilde{D}} L^{-}\big(\Phi(x_i),\Phi(x_j)\big)$$
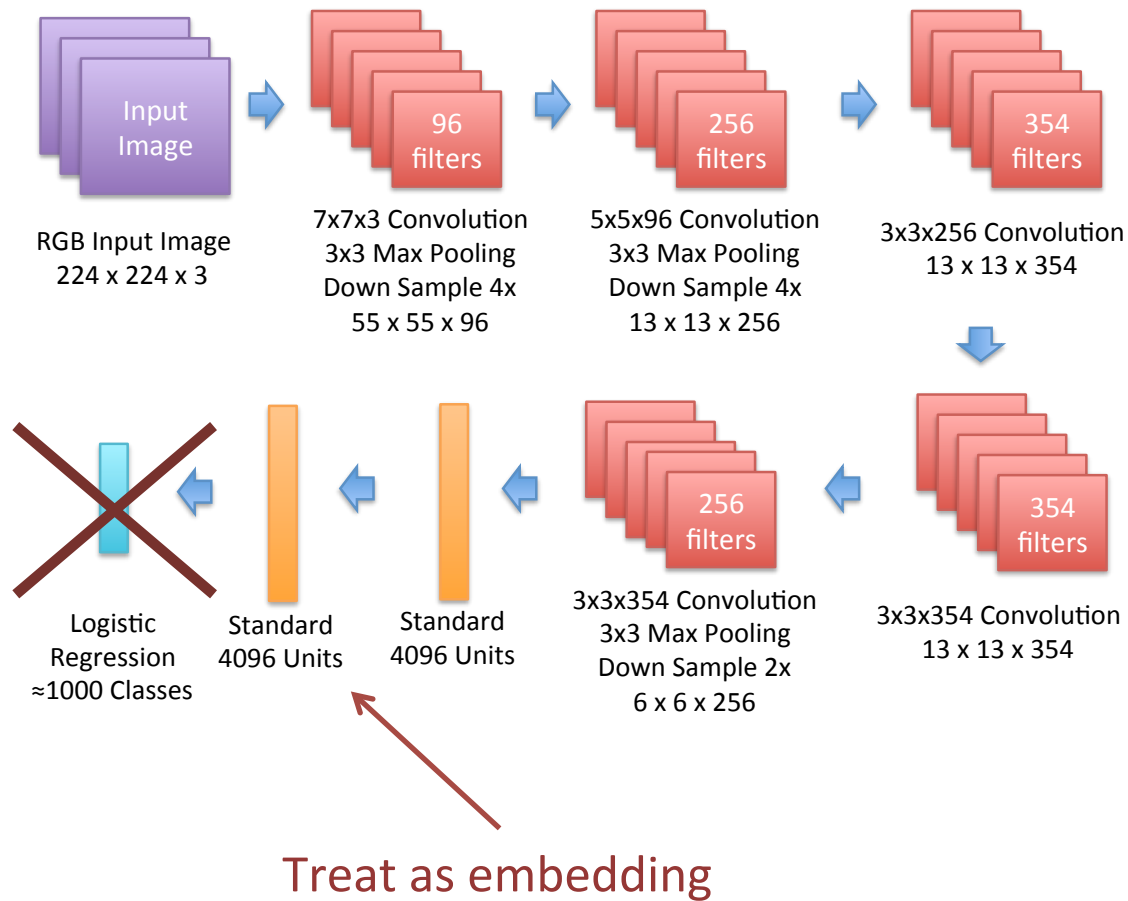
All Model
Parameters

Compatible
Pairs

Incompatible
Pairs

Only pairs in different categories.

# Recall: Convolutional Neural Networks



RGB Input Image
224 x 224 x 3

7x7x3 Convolution
3x3 Max Pooling
Down Sample 4x
55 x 55 x 96

5x5x96 Convolution
3x3 Max Pooling
Down Sample 4x
13 x 13 x 256

3x3x256 Convolution
13 x 13 x 354

3x3x354 Convolution
13 x 13 x 354

3x3x354 Convolution
3x3 Max Pooling
Down Sample 2x
6 x 6 x 256

Standard
4096 Units

Standard
4096 Units

Logistic
Regression
≈1000 Classes

Treat as embedding

# Siamese Convolutional Neural Networks



More details: http://www.cs.cornell.edu/~kb/publications/SIG15ProductNet.pdf

# Recap: Training Goal

Penalizes too far

Penalizes too close

Embedding of image

Embedding of image

$$\operatorname*{arg\,min}_{\Theta} \sum_{(i,j)\in D} L^{+}\big(\Phi(x_i),\Phi(x_j)\big) + \sum_{(i,j)\in \tilde{D}} L^{-}\big(\Phi(x_i),\Phi(x_j)\big)$$

All Model
Parameters

Compatible
Pairs

Incompatible
Pairs

Only pairs in different categories.

**Model Embedding via Siamese Convolutional Neural Network!**

# Training Details

- Want embedding dimension smaller
  - E.g., 128 rather than 4096

- Need to subsample negative pairs
  - Most items are not frequently bought together
  - Negative component can overwhelm objective

http://www.cs.cornell.edu/~andreas/iccv15.pdf

# Suggesting Outfits

Upper
Garment

Lower
Garment

Footware



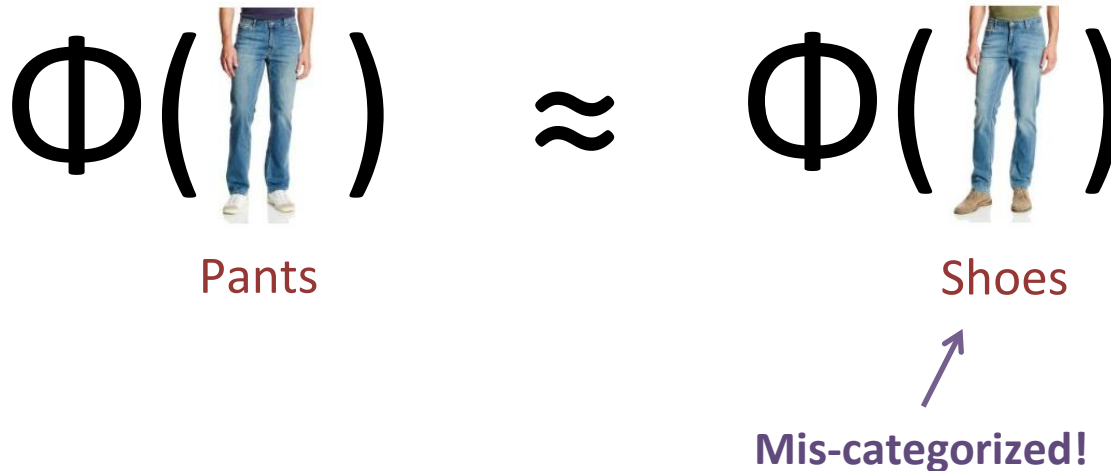http://www.cs.cornell.edu/~andreas/iccv15.pdf

# Suggesting Outfits

- Given query item i
  - Embedding $\varphi_i = \Phi(x_i | \Theta)$
  - Category $c_i$

- For other categories
  - Recommend item with closest embedding $\varphi$

- **Not robust to label noise!**

http://www.cs.cornell.edu/~andreas/iccv15.pdf

# Label Noise

- Amazon category labels are noisy
  - Eg., some pants mis-categorized as shoes

- Pants are visually very similar

$$\Phi(\quad) \quad \approx \quad \Phi(\quad)$$

Pants                               Shoes

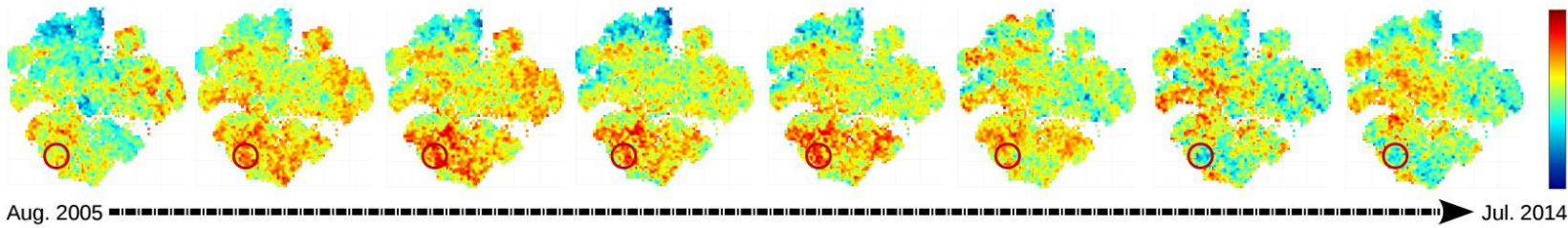**Mis-categorized!**
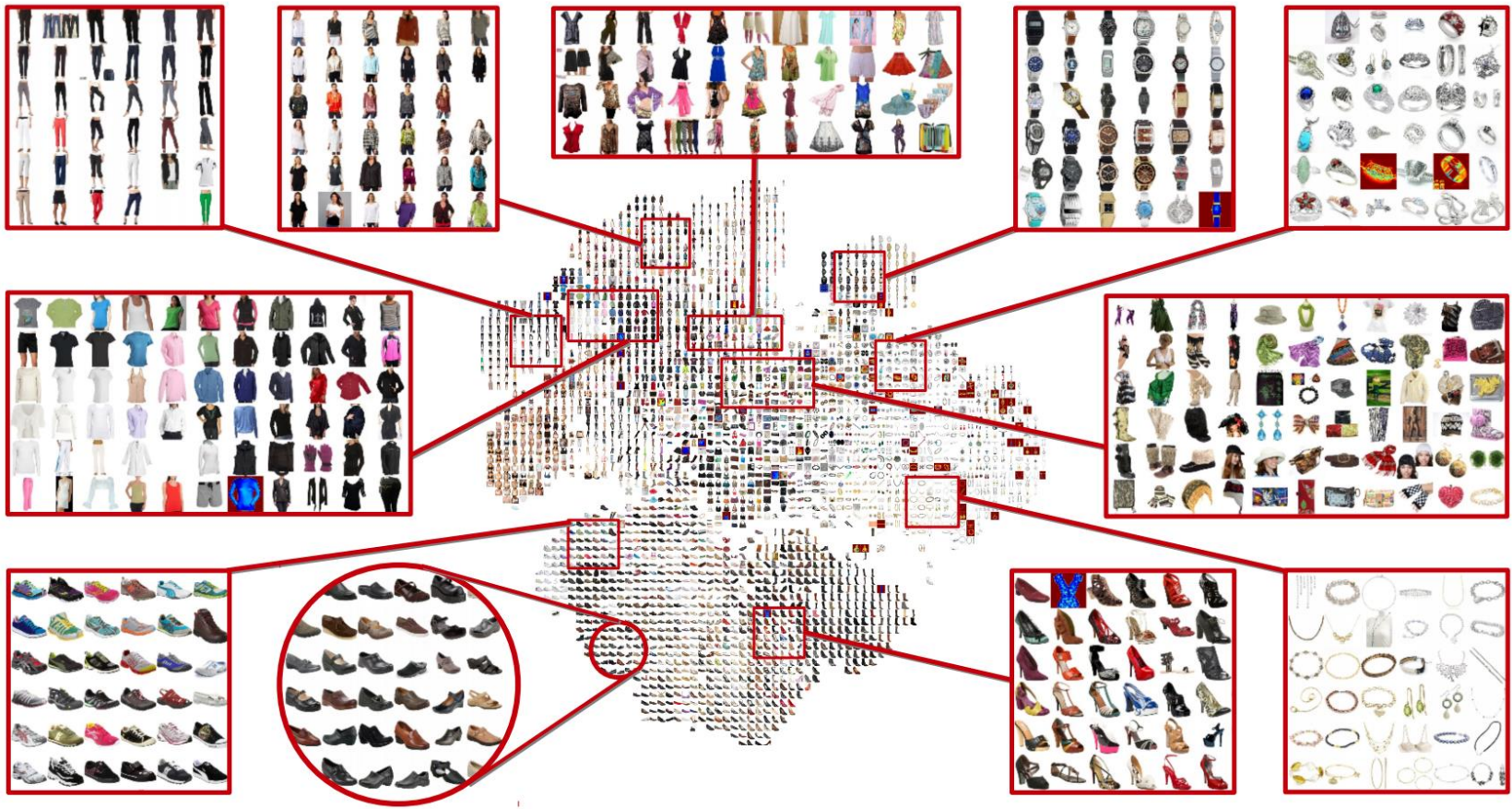
# Making Robust Suggestions

- ## Mis-categorizations are rare
  - Instead of predicting closest shoe…
  - Predict closest cluster of shoes!

- ## Preprocessing: cluster every category

- ## Given input query (category=pants)
  - Find closest cluster center (category=shoes)
  - Output shoes item close to cluster center

http://www.cs.cornell.edu/~andreas/iccv15.pdf

# Compute Coherence of Outfit

## Least coordinated



## Most coordinated

Aug. 2005 ←———————————————————————————————————————————→ Jul. 2014

http://cseweb.ucsd.edu/~jmcauley/pdfs/www16a.pdf

# Next Lecture

- Survey of Advanced Topics
  - Last lecture of the course!

- Next Thursday: Miniproject 2 due