# Using Big Data to Improve Product Quality

Thomas Debeauvais
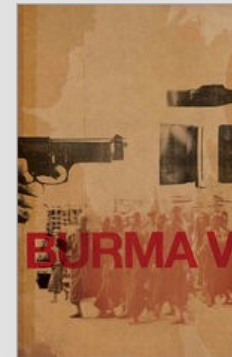
[tdebeauv@uci.edu](mailto:tdebeauv@uci.edu)

# Outline

- Netflix recommendations
- What is big data?
- Industry examples
- A technical example: singular value decomposition for recommender systems
- Thought exercises

"Connecting people to the movies they love"

## Critically-acclaimed Fight-the-System Documentaries

Based on your interest in…



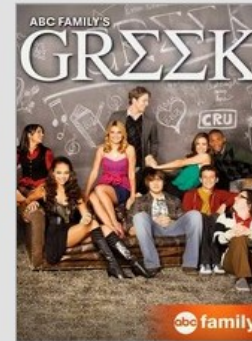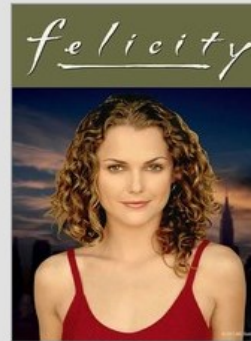## Teen TV Dramas Featuring a Strong Female Lead

Your **taste preferences** created this row.

**TV Dramas**
**Strong Women.**

As well as your interest in…



## Movies Featuring an Epic Nicolas Cage Meltdown



*http://genresofnetflix.tumblr.com/*

4

# Intuitions

1- People rate movies high

*They **generally** don't bother rating bad movies*

# People rate movies high



**Average Movie Ratings by Users of Netflix Training set**

*http://www.netflixprize.com/community/viewtopic.php?pid=5941#p5941*

# Intuitions

2- Older movies are rated higher

*Nostalgia*

*Netflix keeps only the best of old movies*

# Older movies are rated higher



Rating by movie age

From Paddraic Smyth

# Intuitions

3- Movies have genres

[Koren et al. 2009] 11

# Take-away

- Confirm or reject intuitions with data

# WHAT IS BIG DATA?

BIG DATA

# Data mining

# Statistics

# Machine learning

## classification

kernel approximation

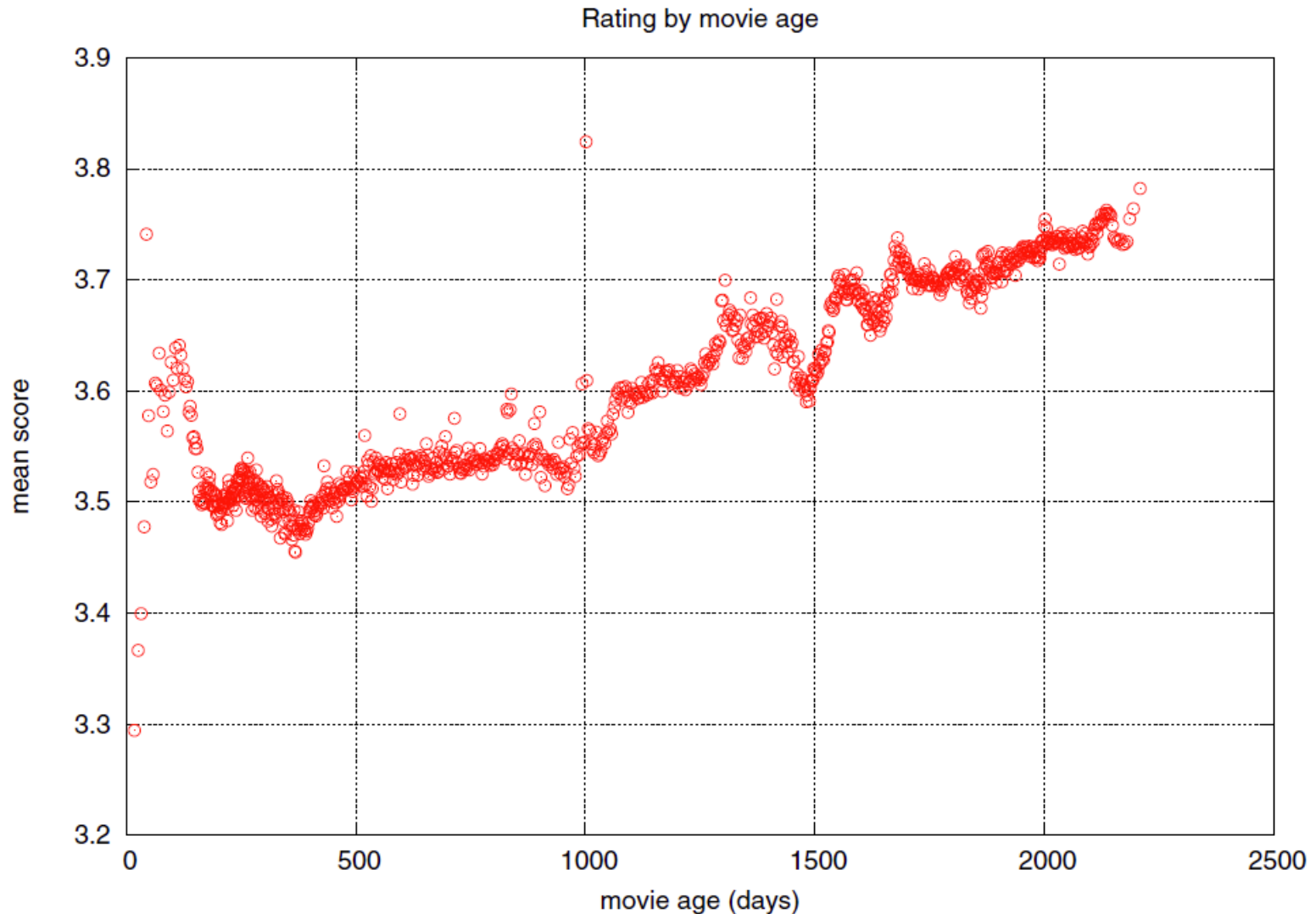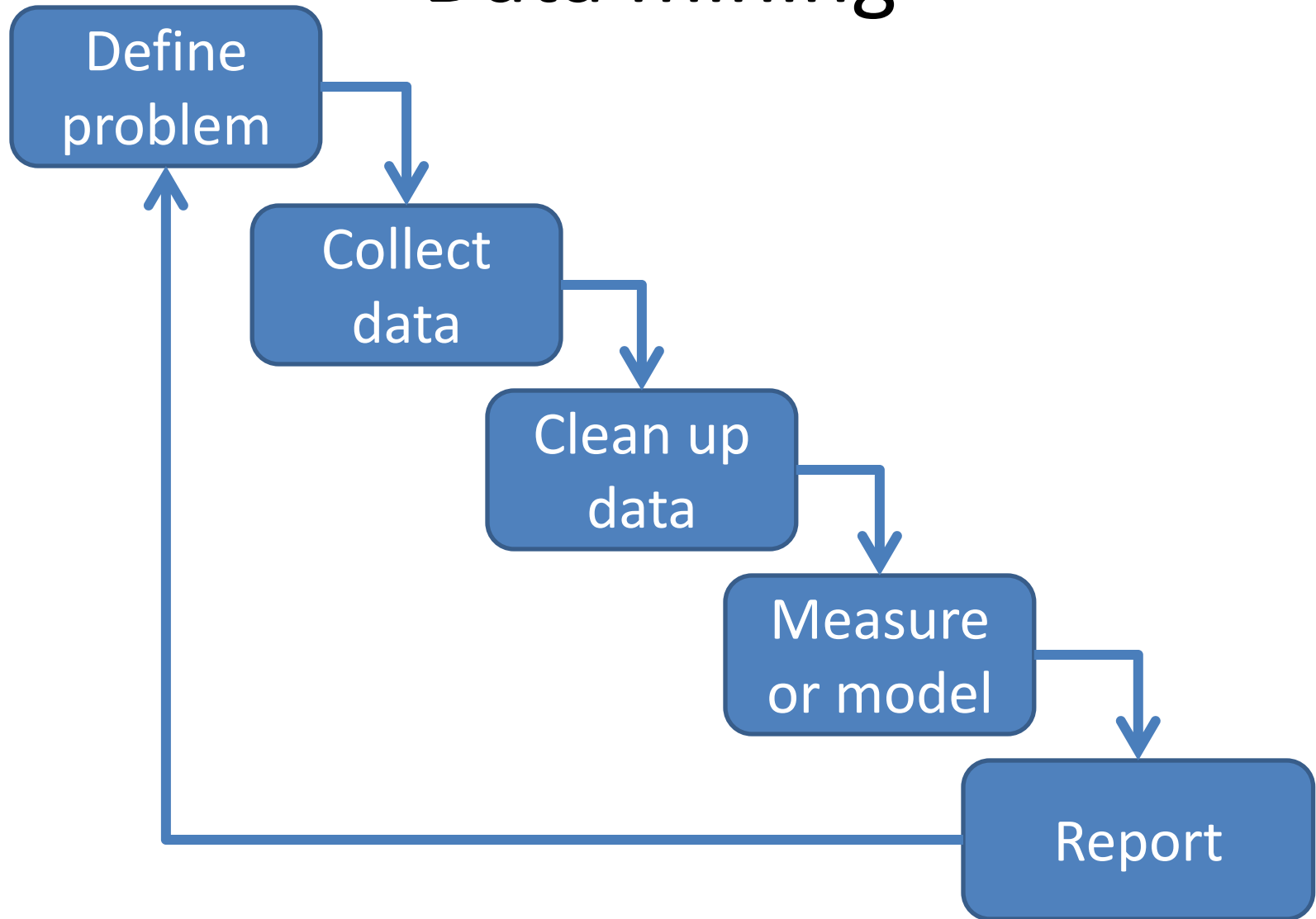SVC Ensemble Classifiers

NOT WORKING

KNeighbors Classifier

NOT WORKING

SGD Classifier

NOT WORKING

NO

Naive Bayes

YES

Text Data

NO

NOT WORKING

YES

<100K samples

NO

Linear SVC

START

get more data

NO

>50 samples

YES

## scikit-learn algorithm cheat-sheet

## regression

SGD Regressor

NO

ElasticNet Lasso

YES

SVR(kernel='rbf') EnsembleRegressors

NOT WORKING

predicting a category

YES

NO

<100K samples

YES

few features should be important

NO

RidgeRegression SVR (kernel='linear')

YES

do you have labeled data

NO

predicting a quantity

YES

NO

## clustering

Spectral Clustering GMM

NOT WORKING

KMeans

YES

number of categories known

YES

NO

<10K samples

NO

<10K samples

YES

NO

MiniBatch KMeans

MeanShift VBGMM

just looking

YES

NO

tough luck

predicting structure

## dimensionality reduction

Randomized PCA

NOT WORKING

YES

<10K samples

NO

Isomap Spectral Embedding

NOT WORKING

LLE

kernel approximation

*http://n-chandra.blogspot.com/2013/01/picking-machine-learning-algorithm.html*

17

# Visualization

# Data engineering

protobuf objects
through RabbitMQ

map-reduce
filtering and
aggregating

non-game
SQL data

protobuf
objects
store

GREENPLUM

legacy
ETL jobs

Tableau
for viz

*http://bowling-bash.blogspot.com/2013/10/blizzards-hadoop-platform.html*

# Take-aways

- Many tools and processes
- Some old, some recent
- Pick the right one for the problem at hand

# INDUSTRY EXAMPLES

# Take-aways

- Useful
- Everywhere
- Creepy?

# SINGULAR VALUE DECOMPOSITION

Factor vector 2 (y-axis), Factor vector 1 (x-axis)

Movies plotted by factor vectors:
- Julien Donkey-Boy
- Kill Bill: Vol. 1
- Natural Born Killers
- I Heart Huckabees
- Punch-Drunk Love
- The Royal Tenenbaums
- Being John Malkovich
- Lost in Translation
- Belle de Jour
- Annie Hall
- Freddy Got Fingered
- Half Baked
- Citizen Kane
- Scarface
- Freddy vs. Jason
- Road Trip
- Sophie's Choice
- Moonstruck
- The Wizard of Oz
- The Way We Were
- The Sound of Music
- The Waltons: Season 1
- The Longest Yard
- The Fast and the Furious
- Armageddon
- Catwoman
- Coyote Ugly
- Maid in Manhattan
- Runaway Bride
- Stepmom
- Sister Act

[Koren et al. 2009] 27

# Ratings data

| | The Longest Yard | Citizen Kane | Fast & Furious | Kill Bill | Sound of Music |
|---|---|---|---|---|---|
| Alice | 1 | 4 | 1 | | 5 |
| Bob | 5 | | 4 | 4 | |
| Claire | | 4 | 2 | | 4 |
| Dan | 4 | 1 | | 3 | |
| Eve | | 3 | | | 5 |
| Frank | 3 | 3 | 5 | 2 | |

(usually 99% sparse)

# Singular Value Decomposition

**Users** **Tastes** **Movies**

| | The Longest Yard | Citizen Kane | Fast & Furious | Kill Bill | Sound of Music |
|---|---|---|---|---|---|
| Alice | 1 | 4 | 1 | M | 5 |
| Bob | 5 | M | 4 | 4 | M |
| Claire | M | 4 | 2 | M | 4 |
| Dan | 4 | 1 | M | 3 | M |
| Eve | M | 3 | M | M | 5 |
| Frank | 3 | 3 | 5 | 2 | M |

**=**

| *0.1 | Factor 1 | Factor 2 |
|---|---|---|
| Alice | -4 | 7 |
| Bob | -5 | -4 |
| Claire | -4 | 2 |
| Dan | -4 | -5 |
| Eve | -4 | 0 |
| Frank | -4 | 0 |

**X**

| | F1 | F2 |
|---|---|---|
| F1 | 18 | 0 |
| F2 | 0 | 4 |

**X**

| *0.1 | The Longest Yard | Citizen Kane | Fast & Furious | Kill Bill | Sound of Music |
|---|---|---|---|---|---|
| F1 | -4 | -4 | -4 | -4 | -5 |
| F2 | -5 | 4 | -6 | 1 | 4 |

M = global average = 3.3

[Koren et al. 2009] 30

[Koren et al. 2009] 31
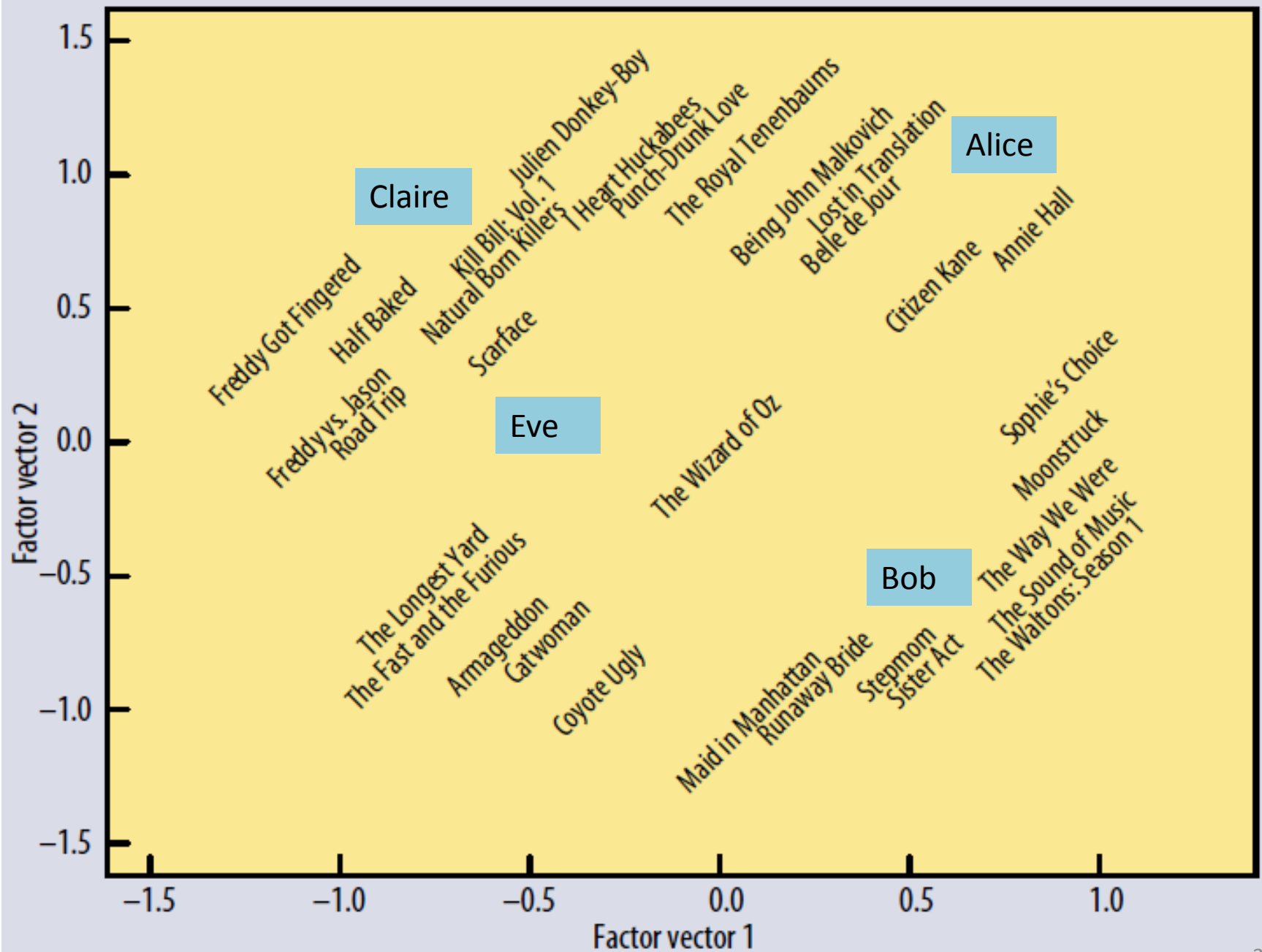
# Densifying sparse matrices

Number of ratings x100

– Storage x100

– SVD is O(N$^3$), CPU x1M

Minimize square error on known ratings using stochastic gradient descent

|  | Armageddon | Citizen Kane | Fast & Furious | Kill Bill | Sound of Music |
|---|---|---|---|---|---|
| Alice | 1 | 4 | 1 | M | 5 |
| Bob | 5 | M | 4 | 4 | M |
| Claire | M | 4 | 2 | M | 4 |
| Dan | 4 | 1 | M | 3 | M |
| Eve | M | 3 | M | M | 5 |
| Frank | 3 | 3 | 5 | 2 | M |

$$\min_{q_*, p_*} \sum_{(u,i) \in \kappa} (r_{ui} - q_i^T p_u)^2 + \lambda(\| q_i \|^2 + \| p_u \|^2)$$

# Similar techniques

- Latent Semantic Analysis
  - Movies -> documents, ratings -> tf-idf
- Latent Dirichlet Allocation
  - Cf "A theory of aspects as latent topics"
- Principal Component Analysis
  - Eigen value decomposition
  - Of the covariance matrix

# Take-aways

- Matrices
- Need practical methods
- The math escalates quickly
  - But you may need it ...

# HOW WOULD YOU DO …

# Amazon

- "Frequently bought together"
  - Association rule mining

- Customer clustering
  - SVD


- (Stock management
  - Poisson process?)

# stack**overflow**

Questions | Tags | Users | Badges | Unanswered

Ask Question

**Title** | Is peer to peer possible with python's asyncore?

**Questions that may already have your answer**

| 2 | **Is is possible to code an app that can detect and connect to peers without a central servers? How?** |
| 0 | **is it possible peer to peer connection over wan without any server?** 1 |
| 0 | **Python Error "Connection reset by peer" 10054** |
| 0 | **Web Audio API and WebRTC** 1 |
| 0 | **can peer A connect to peer B if server is on peer B?** 1 |

**B** *I* | 🌐 " { } 🖼 | ≣ ≣ ≣ ≣ | ↶ ↷ | **?**

Links | Images | Styling/Headers | Lists | Blockquotes | Code | HTML | advanced help »

```
"This module provides the basic infrastructure for writing asynchronous socket service
clients and servers."
http://docs.python.org/2/library/asyncore.html#module-asyncore

If I want a non- client-server architecture, such as peer to peer, can asynchat/asyncore
still do the job?
```

draft saved

"This module provides the basic infrastructure for writing asynchronous socket service clients and servers." http://docs.python.org/2/library/asyncore.html#module-asyncore

If I want a non- client-server architecture, such as peer to peer, can asynchat/asyncore still do the job?

**Tags**

python ✕ | asyncore ✕

# Dating website

- Match %
  - SVD on weighted questions
    - "Would you rather be weird or normal?"= 99% weird
    - "Sex before marriage?"
  - SVD on tf-idf from profile essays
  - Cosine distance in SVD space
- Actual information in a profile
  - Everybody loves travelling, but they say it differently
  - "I love travelling!" = "travel" = low tf-idf
  - "Have a passport!" = "passport" = high tf-idf
  - Tf-idf is not always appropriate!

Thomas Debeauvais, tdebeauv@uci.edu

# THANKS

# References

- Padhraic Smyth's slides on Netflix recsys [http://www.ics.uci.edu/~smyth/courses/cs277/public_slides/recommender_systems_part2.pdf](http://www.ics.uci.edu/~smyth/courses/cs277/public_slides/recommender_systems_part2.pdf)

- Koren et al. 2009: Matrix factorization techniques for recommender systems

- [http://www.theatlantic.com/technology/archive/2014/01/how-netflix-reverse-engineered-hollywood/282679/](http://www.theatlantic.com/technology/archive/2014/01/how-netflix-reverse-engineered-hollywood/282679/)

- [http://online.wsj.com/news/articles/SB10001424052702303453004579290632128929194](http://online.wsj.com/news/articles/SB10001424052702303453004579290632128929194)