

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

#### Announcements

- Kaggle Miniproject is closed
   Report due Thursday
- Public Leaderboard
  - How well you think you did
- Private Leaderboard now viewable
   How well you actually did

∆5d	Team Name	Score 🚱	Entries
<b>↑5</b>	OYGBYNGB 🖈	0.94975	14
<b>↑5</b>	3 Shades of Brown 🕫	0.94750	24
↓ <b>2</b>	TurboBoost 🖈	0.94700	45
<b>†1</b>	JP #	0.94600	11
<b>↑7</b>	Eman 🕫	0.94575	93
new	^thisAlgorithmIsBetterThanOurs #	0.94575	15
ţ5	kanyeblessed 🖈	0.94425	15
<b>↓5</b>	TopSauce 💶	0.94325	15
<b>↑14</b>	The Usual Suspects 💵	0.94300	29
<b>↓6</b>	Kobe Wan Kenobi 🍂	0.94275	47
_	CHM 💵	0.94275	45
<b>†1</b>	SD_YC 🗚	0.94275	50
↑ <b>12</b>	zed 📲	0.94225	20
new	WhiSK	0.94200	12
<b>↓7</b>	tyj518	0.94175	7
<b>†1</b>	The Red Brothers 🕫	0.94175	20
<b>↓8</b>	fboemer	0.94150	3
new	IMSA 📭	0.94150	9
new	ANND	0.94150	10
† <b>2</b>	I Just Kaggled 🏨	0.94100	25
	Δ5d           ↑5           ↓2           ↑1           ↑7           ↓5           ↓5           ↓1           ↓5           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓1           ↓2	A5dTeam Name15OYGBYNGB15OYGBYNGB16Shades of Brown17Shades of Brown18TurboBoost11JP17Eman17Eman18AthisAlgorithmIsBetterThanOurs19TopSauce14The Usual Suspects16Kobe Wan Kenobi11SD_YC11SD_YC12zed13The Red Brothers14The Red Brothers15MiSA16MND17Just Kaggled12Just Kaggled	A5dTeam NameScore Q15OYGBYNGB0.94975153 Shades of Brown0.9470012TurboBoost0.9470011JP0.9460011JP0.9450512Eman0.9457514Anyeblessed0.9442515TopSauce0.9432514The Usual Suspects0.9432515SD_YC0.9427516SD_YC0.9427517SD_YC0.9427518SD_YC0.9427519thisK0.9427511The Red Brothers0.9417512MSA0.9415013MSA0.9415014MND0.94150

#	∆rank	Team Name	Score @	Entries
1	<b>†1</b>	3 Shades of Brown 🏨	0.93984	24
2	<b>↑7</b>	The Usual Suspects 💶	0.93916	29
3	<b>↑15</b>	IMSA 🗈	0.93916	9
4	<b>∱3</b>	kanyeblessed 🔎	0.93865	15
5	<b>↑6</b>	CHM 🕫	0.93831	45
6	<b>↓2</b>	JP #	0.93814	11
7	<b>↓4</b>	TurboBoost 🖈	0.93712	45
8	<b>↓7</b>	OYGBYNGB 🕫	0.93695	14
9	<b>↑6</b>	tyj518	0.93678	7
10	<b>↓2</b>	TopSauce 💶	0.93626	15
11	<b>↓5</b>	^thisAlgorithmIsBetterThanOurs 🎿	0.93626	15
12	<b>↓2</b>	Kobe Wan Kenobi 🏨	0.93592	47
13	<b>↑10</b>	Ball Is Life 🖈	0.93575	7
14		zed 🕫	0.93490	20
15	<b>↑4</b>	ANND	0.93473	10
16	<b>↓2</b>	WhiSK	0.93422	12
17	<b>∱9</b>	Peaches	0.93405	17
18	<b>↑11</b>	AdaBreaker 📌	0.93388	13
19	↓ <b>3</b>	The Red Brothers 📭	0.93388	20
20	<b>↑5</b>	Team Rocket 🕫	0.93371	13

### Last Week

- Dimensionality Reduction
- Clustering

• Latent Factor Models

- Learn low-dimensional representation of data

#### This Lecture

• Embeddings

Alternative form of dimensionality reduction

• Locally Linear Embeddings

• Markov Embeddings

# Embedding

- Learn a representation U
  - Each column u corresponds to data point
- Semantics encoded via d(u,u')

Distance between points



http://www.sciencemag.org/content/290/5500/2323.full.pdf

• Given: 
$$S = \{x_i\}_{i=1}^N$$

Unsupervised Learning

- Learn U such that local linearity is preserved
  - Lower dimensional than x
  - "Manifold Learning"

Any neighborhood looks like a linear plane



x's



u's

https://www.cs.nyu.edu/~roweis/lle/

• Create B(i)

$$S = \left\{ x_i \right\}_{i=1}^N$$

- B nearest neighbors of x<sub>i</sub>
- Assumption: B(i) is approximately linear
- $-x_i$  can be written as a convex combination of  $x_i$  in B(i)

$$x_i \approx \sum_{j \in B(i)} W_{ij} x_j$$
$$\sum_{j \in B(i)} W_{ij} = 1$$

https://www.cs.nyu.edu/~roweis/lle/



Given Neighbors B(i), solve local linear approximation W:

$$\underset{W}{\operatorname{argmin}} \sum_{i} \left\| x_{i} - \sum_{j \in B(i)} W_{ij} x_{j} \right\|^{2} = \underset{W}{\operatorname{argmin}} \sum_{i} W_{i,*}^{T} C^{i} W_{i,*} \qquad \qquad \sum_{j \in B(i)} W_{ij} = 1$$

$$\begin{aligned} \left\| x_{i} - \sum_{j \in B(i)} W_{ij} x_{j} \right\|^{2} &= \left\| \sum_{j \in B(i)} W_{ij} (x_{i} - x_{j}) \right\|^{2} \\ &= \left( \sum_{j \in B(i)} W_{ij} (x_{i} - x_{j}) \right)^{T} \left( \sum_{j \in B(i)} W_{ij} (x_{i} - x_{j}) \right) \\ &= \sum_{j \in B(i)} \sum_{k \in B(i)} W_{ij} W_{ik} C_{jk}^{i} \\ &= W_{i,*}^{T} C^{i} W_{i,*} \end{aligned}$$

#### https://www.cs.nyu.edu/~roweis/lle/

Given Neighbors B(i), solve local linear approximation W:

$$\underset{W}{\operatorname{argmin}} \sum_{i} \left\| x_{i} - \sum_{j \in B(i)} W_{ij} x_{j} \right\|^{2} = \underset{W}{\operatorname{argmin}} \sum_{i} W_{i,*}^{T} C^{i} W_{i,*}$$

$$\sum_{j \in B(i)} W_{ij} = 1$$

$$C_{jk}^{i} = (x_{i} - x_{j})^{T} (x_{i} - x_{j})$$

- Every x<sub>i</sub> is approximated as a convex combination of neighbors
  - How to solve?



#### Lagrange Multipliers



http://en.wikipedia.org/wiki/Lagrange\_multiplier 11

#### Solving Locally Linear Approximation

Lagrangian:

$$L(W,\lambda) = \sum_{i} \left( W_{i,*}^{T} C^{i} W_{i,*} - \lambda_{i} \left( \vec{1}^{T} W_{i,*} - 1 \right) \right) \qquad \sum_{j} W_{ij} = \vec{1}^{T} W_{i,*}$$

$$\partial_{W_{i,*}} L(W,\lambda) = 2C^i W_{i,*} - \lambda_i \vec{1}$$

### Locally Linear Approximation

• Invariant to:

- Rotation 
$$Ax_i \approx \sum_{j \in B(i)} AW_{ij}x_j$$

$$x_i \approx \sum_{j \in B(i)} W_{ij} x_j$$
$$\sum_{j \in B(i)} W_{ij} = 1$$

-Scaling 
$$5x_i \approx \sum_{j \in B(i)} 5W_{ij}x_j$$

- Translation 
$$x_i + x' \approx \sum_{j \in B(i)} W_{ij} (x_j + x')$$

#### Story So Far: Locally Linear Embeddings

Given Neighbors B(i), solve local linear approximation W:

$$\underset{W}{\operatorname{argmin}} \sum_{i} \left\| x_{i} - \sum_{j \in B(i)} W_{ij} x_{j} \right\|^{2} = \underset{W}{\operatorname{argmin}} \sum_{i} W_{i,*}^{T} C^{i} W_{i,*} \qquad \qquad \sum_{j \in B(i)} W_{ij} = 1$$

Solution via Lagrange Multipliers:

$$C_{jk}^i = (x_i - x_j)^T (x_i - x_k)$$



#### Locally Linear Approximation

https://www.cs.nyu.edu/~roweis/lle/



## **Recall:** Locally Linear Embedding

• Given: 
$$S = \{x_i\}_{i=1}^N$$

- Learn U such that local linearity is preserved
  - Lower dimensional than x
  - "Manifold Learning"



https://www.cs.nyu.edu/~roweis/lle/

# **Dimensionality Reduction**

Given local approximation W, learn lower dimensional representation:

$$\underset{U}{\operatorname{argmin}} \sum_{i} \left\| u_{i} - \sum_{j \in B(i)} W_{ij} u_{j} \right\|^{2}$$

- Find low dimensional U
  - Preserves approximate local linearity



https://www.cs.nyu.edu/~roweis/lle/

Given local approximation W, learn lower dimensional representation:

$$\underset{U}{\operatorname{argmin}} \sum_{i} \left\| u_{i} - \sum_{j \in B(i)} W_{ij} u_{j} \right\|^{2}$$

$$UU^{T} = I_{K}$$
$$\sum_{i} u_{i} = \vec{0}$$

• Rewrite as:

$$\operatorname{argmin}_{U} \sum_{ij} M_{ij} \left( u_i^T u_j \right) = \operatorname{trace} \left( UMU^T u_j \right)$$
$$M_{ij} = 1_{[i=j]} - W_{ij} - W_{ji} + \sum_k W_{ki} W_{kj}$$
$$M = (I_N - W)^T (I_N - W)$$

Symmetric positive semidefinite

https://www.cs.nyu.edu/~roweis/lle/

Given local approximation W, learn lower dimensional representation:

$$\underset{U}{\operatorname{argmin}} \sum_{ij} M_{ij} \left( u_i^T u_j \right) = \operatorname{trace} \left( U M U^T \right) \qquad U U^T = I_K$$

$$\sum_{i} u_{i} = \vec{0}$$

 $uu^T = 1$ 

$$\operatorname{argmin}_{u} \sum_{ij} M_{ij} \left( u_i^T u_j \right) = \operatorname{trace} \left( u M u^T \right)$$
$$= \operatorname{argmax}_{u} \operatorname{trace} \left( u M^+ u^T \right)$$
pseudoinverse

http://en.wikipedia.org/wiki/Min-max\_theorem

#### **Recap:** Principal Component Analysis

0

 $M = V\Lambda V^T$ 

- Each column of V is an Eigenvector
- Each  $\lambda$  is an Eigenvalue ( $\lambda_1 \ge \lambda_2 \ge ...$ )

$$M^+ = V\Lambda^+ V^T$$

$$MM^{+} = V\Lambda\Lambda^{+}V^{T} = V_{1:2}V_{1:2}^{T} = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 0 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & 0 & \\ & & & 0 \end{bmatrix}$$

г

Г

$$\Lambda^{+} = \begin{bmatrix} 1/\lambda_{1} & & \\ & 1/\lambda_{2} & \\ & & 0 & \\ & & & 0 \end{bmatrix}$$

Т

Given local approximation W, learn lower dimensional representation:

$$\underset{U}{\operatorname{argmin}} \sum_{ij} M_{ij} \left( u_i^T u_j \right) = \operatorname{trace} \left( U M U^T \right) \qquad U U^T = I_K$$

• K=1:

 $\sum_{i} u_{i} = \vec{0}$ 

- u = principal eigenvector of M<sup>+</sup>
- u = smallest non-trivial eigenvector of M
  - Corresponds to smallest non-zero eigenvalue

#### • General K

- U = top K principal eigenvectors of M<sup>+</sup>
- U = bottom K non-trivial eigenvectors of M
  - Corresponds to bottom K non-zero eigenvalues

https://www.cs.nyu.edu/~roweis/lle/

http://en.wikipedia.org/wiki/Min-max\_theorem

### **Recap:** Locally Linear Embedding

- Generate nearest neighbors of each x<sub>i</sub>, B(i)
- Compute Local Linear Approximation:

$$\underset{W}{\operatorname{argmin}} \sum_{i} \left\| x_{i} - \sum_{j \in B(i)} W_{ij} x_{j} \right\|^{2} \qquad \qquad \sum_{j \in B(i)} W_{ij} = 1$$

Compute low dimensional embedding

$$\underset{U}{\operatorname{argmin}} \sum_{i} \left\| u_{i} - \sum_{j \in B(i)} W_{ij} u_{j} \right\|^{2} \qquad \qquad UU^{T} = I_{K}$$
$$\sum_{i} u_{i} = \vec{0}$$

#### **Results for Different Neighborhoods**



#### https://www.cs.nyu.edu/~roweis/lle/gallery.html

#### **Embeddings vs Latent Factor Models**

- Both define low-dimensional representation
- Embeddings preserve distance:

$$\left\|\boldsymbol{u}_{i}-\boldsymbol{u}_{j}\right\|^{2} \approx \left\|\boldsymbol{x}_{i}-\boldsymbol{x}_{j}\right\|^{2}$$

• Latent Factor preserve inner product:

$$\boldsymbol{u}_i^T \boldsymbol{u}_j \approx \boldsymbol{x}_i^T \boldsymbol{x}_j$$

• Relationship:

$$\left\|u_{i} - u_{j}\right\|^{2} = \left\|u_{i}\right\|^{2} + \left\|u_{j}\right\|^{2} - 2u_{i}^{T}u_{j}$$

#### **Visualization Semantics**



#### **Latent Factor Model**

Similarity measured via dot product Rotational semantics Can interpret axes Can only visualize 2 axes at a time



Similarity measured via distance Clustering/locality semantics Cannot interpret axes Can visualize many clusters simultaneously



#### Latent Markov Embeddings

## Latent Markov Embeddings

- Locally Linear Embedding is conventional unsupervised learning
  - Given raw features x<sub>i</sub>
  - I.e., find low-dimensional U that preserves approximate local linearity
- Latent Markov Embedding is a feature learning problem
  - E.g., learn low-dimensional U that captures user-generated feedback

## **Playlist Embedding**



- Users generate song playlists
  - Treat as training data
- Can we learn a probabilistic model of playlists?

## **Probabilistic Markov Modeling**

• Training set:

$$S = \{s_1, \dots, s_{|S|}\} \qquad D = \{p_i\}_{i=1}^N \qquad p_i = \left\langle p_i^1, \dots, p_i^{N_i} \right\rangle$$
  
Songs Playlists Playlist Definition

- Goal: Learn a probabilistic Markov model of playlists:  $P(p_i^j \,|\, p_i^{j-1})$
- What is the form of P?

http://www.cs.cornell.edu/People/tj/publications/chen\_etal\_12a.pdf

## First Try: Probability Tables

P(s s')	s <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>	S <sub>4</sub>	<b>S</b> <sub>5</sub>	s <sub>6</sub>	S <sub>7</sub>	S <sub>start</sub>
S <sub>1</sub>	0.01	0.03	0.01	0.11	0.04	0.04	0.01	0.05
s <sub>2</sub>	0.03	0.01	0.04	0.03	0.02	0.01	0.02	0.02
S <sub>3</sub>	0.01	0.01	0.01	0.07	0.02	0.02	0.05	0.09
s <sub>4</sub>	0.02	0.11	0.07	0.01	0.07	0.04	0.01	0.01
<b>S</b> <sub>5</sub>	0.04	0.01	0.02	0.17	0.01	0.01	0.10	0.02
s <sub>6</sub>	0.01	0.02	0.03	0.01	0.01	0.01	0.01	0.08
<b>S</b> <sub>7</sub>	0.07	0.02	0.01	0.01	0.03	0.09	0.03	0.01

## First Try: Probability Tables

$s_1$ 0.01       0.03       0.01       0.11       0.04       0.04       0.01       0.05 $s_2$ 0.03       0.01       0.04       0.02       0.01       0.02       0.02 $s_3$ 0.01       0.01       0.01       0.07       0.02       0.02       0.09 $s_4$ 0.02       0.11       0.07       0.07       0.04       0.01       0.01 $s_4$ 0.02       0.11       0.07       0.07       0.04       0.01       0.01 $s_5$ $s_5$ $s_6$ $s_7$ $\#Parameters = O( S ^2)  !!$ $\#Parameters = O( S ^2)  !!       \#Parameters = O( S ^2)  !       \#Parameters = O( S ^2)  ! $	P(s s')	s <sub>1</sub>	s <sub>2</sub>	S <sub>3</sub>	S <sub>4</sub>	<b>S</b> <sub>5</sub>	s <sub>6</sub>	S <sub>7</sub>	S <sub>start</sub>	
$s_2$ 0.03       0.01       0.04       0.03       0.02       0.01       0.02       0.02 $s_3$ 0.01       0.01       0.01       0.07       0.02       0.02       0.09 $s_4$ 0.02       0.11       0.07       0.01       0.04       0.01       0.01 $s_5$ $s_5$ $s_5$ $s_7$ $#Parameters = O( S ^2) !!!$ $#Parameters = O( S ^2) !!!$	s <sub>1</sub>	0.01	0.03	0.01	0.11	0.04	0.04	0.01	0.05	
$s_3$ 0.01       0.01       0.07       0.02       0.05       0.09 $s_4$ 0.02       0.11       0.07       0.07       0.04       0.01       0.01 $s_5$ $s_6$ $s_7$ $#Parameters = O( S ^2) !!!$ $#Parameters = O( S ^2) !!!$	s <sub>2</sub>	0.03	0.01	0.04	0.03	0.02	0.01	0.02	0.02	
$s_4$ 0.02       0.11       0.07       0.01       0.07       0.04       0.01       0.01 $s_5$ $s_6$ $s_7$ #Parameters = O( S ^2) !!! $s_1$ $s_1$ $s_2$ $s_1$ $s_2$ <td>S<sub>3</sub></td> <td>0.01</td> <td>0.01</td> <td>0.01</td> <td>0.07</td> <td>0.02</td> <td>0.02</td> <td>0.05</td> <td>0.09</td>	S <sub>3</sub>	0.01	0.01	0.01	0.07	0.02	0.02	0.05	0.09	
$s_{5}$ $s_{6}$ $s_{7}$ #Parameters = O( S  <sup>2</sup> ) !!!	s <sub>4</sub>	0.02	0.11	0.07	0.01	0.07	0.04	0.01	0.01	
$s_{6}$ $s_{7}$ #Parameters = O( S  <sup>2</sup> ) !!!	s <sub>5</sub>									
#Parameters = $O( S ^2)$ !!!	s <sub>6</sub>									
	<b>S</b> <sub>7</sub>	$#Parameters = O( S ^2)$ []]								

#### Second Try: Hidden Markov Models

$$P(p_i, z) = P(End \mid z^{N_i}) \prod_{j=1}^{N_i} P(z^j \mid z^{j-1}) \prod_{j=1}^{N_j} P(p_i^j \mid z^j)$$

 $P(z^{j} | z^{j-1})$  • #Parameters = O(K<sup>2</sup>)

 $P(p_i^j | z^j)$  • #Parameters = O(|S|K)

#### Problem with Hidden Markov Models

$$P(p_i, z) = P(End \mid z^{N_i}) \prod_{j=1}^{N_i} P(z^j \mid z^{j-1}) \prod_{j=1}^{N_j} P(p_i^j \mid z^j)$$

Need to reliably estimate P(s|z)

$$S = \{s_1, ..., s_{|S|}\} \qquad D = \{p_i\}_{i=1}^N \qquad p_i = \langle p_i^1, ..., p_i^{N_i} \rangle$$

Lots of "missing values" in this training set

## Latent Markov Embedding

$$P(s \mid s') \propto \exp\left\{-\left\|u_s - v_{s'}\right\|^2\right\}$$

u<sub>s</sub>: entry point of song s v<sub>s</sub>: exit point of song s

$$P(s \mid s') = \frac{\exp\{-\|u_s - v_{s'}\|^2\}}{\sum_{s''} \exp\{-\|u_{s''} - v_{s'}\|^2\}}$$

- "Log-Radial" function
  - (my own terminology)

http://www.cs.cornell.edu/People/tj/publications/chen\_etal\_12a.pdf

#### **Log-Radial Functions**

$$\frac{P(s \mid s')}{P(s'' \mid s')} = \frac{\exp\{-\|u_s - v_{s'}\|^2\}}{\exp\{-\|u_{s''} - v_{s'}\|^2\}}$$

2K parameters per song 2|S|K parameters total



Each ring defines an equivalence class of transition probabilities

### Learning Problem

$$S = \left\{ s_1, \dots s_{|S|} \right\}$$
Songs

$$D = \left\{ p_i \right\}_{i=1}^N$$

Playlists

$$p_i = \left\langle p_i^1, \dots, p_i^{N_i} \right\rangle$$

**Playlist Definition** 

• Learning Goal:

$$\operatorname{argmax}_{U,V} \prod_{i} P(p_{i}) = \prod_{i} \prod_{j} P(p_{i}^{j} \mid p_{i}^{j-1})$$
$$P(s \mid s') = \frac{\exp\{-\|u_{s} - v_{s'}\|^{2}\}}{\sum_{s''} \exp\{-\|u_{s''} - v_{s'}\|^{2}\}} = \frac{\exp\{-\|u_{s} - v_{s'}\|^{2}\}}{Z(s')}$$

http://www.cs.cornell.edu/People/tj/publications/chen\_etal\_12a.pdf

## Minimize Neg Log Likelihood

$$\underset{U,V}{\operatorname{argmax}} \prod_{i} \prod_{j} P(p_{i}^{j} \mid p_{i}^{j-1}) = \underset{U,V}{\operatorname{argmin}} \sum_{i} \sum_{j} -\log P(p_{i}^{j} \mid p_{i}^{j-1})$$

- Solve using gradient descent
  - Homework question: derive the gradient formula
  - Random initialization
- Normalization constant hard to compute:
  - Approximation heuristics
    - See paper

http://www.cs.cornell.edu/People/tj/publications/chen\_etal\_12a.pdf

Lecture 14: Embeddings

 $P(s \mid s') = \frac{\exp\{-\|u_s - v_{s'}\|^2\}}{Z(s')}$ 

### **Simpler Version**

• Dual point model:

$$P(s \mid s') = \frac{\exp\{-\|u_s - v_{s'}\|^2\}}{Z(s')}$$

Single point model:

$$P(s \mid s') = \frac{\exp\{-\|u_s - u_{s'}\|^2\}}{Z(s')}$$

- Transitions are symmetric
  - (almost)
- Exact same form of training problem

## Visualization in 2D

#### Simpler version: Single Point Model

$$P(s \mid s') = \frac{\exp\{-\|u_s - u_{s'}\|^2\}}{Z(s')}$$

Single point model is easier to visualize



http://www.cs.cornell.edu/People/tj/publications/chen\_etal\_12a.pdf

## Sampling New Playlists

• Given partial playlist:

$$p = \left\langle p^1, \dots p^j \right\rangle$$

Generate next song for playlist p<sup>j+1</sup>

- Sample according to:

$$P(s \mid p^{j}) = \frac{\exp\left\{-\left\|u_{s} - v_{p^{j}}\right\|^{2}\right\}}{Z(p^{j})} \qquad P(s \mid p^{j}) = \frac{\exp\left\{-\left\|u_{s} - u_{p^{j}}\right\|^{2}\right\}}{Z(p^{j})}$$

**Dual Point Model** 

Single Point Model

http://www.cs.cornell.edu/People/tj/publications/chen\_etal\_12a.pdf

#### Demo

#### http://jimi.ithaca.edu/~dturnbull/research/lme/lmeDemo.html

### What About New Songs?

• Suppose we've trained U:

$$P(s \mid s') = \frac{\exp\{-\|u_s - u_{s'}\|^2\}}{Z(s')}$$

- What if we add a new song s'?
  - No playlists created by users yet...
  - Only options:  $u_{s'} = 0$  or  $u_{s'} = random$ 
    - Both are terrible!

# Song & Tag Embedding

- Songs are usually added with tags
  - E.g., indie rock, country
  - Treat as features or attributes of songs
- How to leverage tags to generate a reasonable embedding of new songs?
  - Learn an embedding of tags as well!

$$S = \left\{ s_1, \dots s_{|S|} \right\}$$

$$D = \left\{ p_i \right\}_{i=1}^N$$

$$p_i = \left\langle p_i^1, \dots, p_i^{N_i} \right\rangle$$

Songs

Playlists

**Playlist Definition** 

$$T = \left\{T_1, \dots T_{|S|}\right\}$$

Tags for Each Song

Learning Objective:  $\underset{U,A}{\operatorname{argmax}} P(D | U) P(U | A, T)$ 

Same term as before: 
$$P(D | U) = \prod_{i} P(p_i | U) = \prod_{i} \prod_{j} P(p_i^j | p_i^{j-1}, U)$$

Song embedding  $\approx$  average of tag embeddings:  $P(U \mid A, T) = \prod_{s} P(u_s \mid A, T_s) \propto \prod_{s} \exp\left\{-\lambda \left\| u_s - \frac{1}{|T_s|} \sum_{t \in T_s} A_t \right\|^2\right\}$ 

Solve using gradient descent:

http://www.cs.cornell.edu/People/tj/publications/moore\_etal\_12a.pdf

## Visualization in 2D



#### http://www.cs.cornell.edu/People/tj/publications/moore\_etal\_12a.pdf

## **Revisited:** What About New Songs?

No user has yet s' added to playlist
 – So no evidence from playlist training data:

s' does not appear in  $D = \{p_i\}_{i=1}^N$ 

- Assume new song has been tagged  $\mathsf{T}_{\mathsf{s}'}$ 

- The  $u_{s'}$  = average of  $A_t$  for tags t in  $T_{s'}$ 

Implication from objective:

 $\underset{U,A}{\operatorname{argmax}} P(D | U) P(U | A, T)$ 

## **Recap: Embeddings**

 Learn a low-dimensional representation of items U

 Capture semantics using distance between items u, u'

Can be easier to visualize than latent factor models

#### Next Lecture

- Recent Applications of Latent Factor Models
- Low-rank Spatial Model for Basketball Play Prediction
- Low-rank Tensor Model for Collaborative Clustering
- Miniproject 1 report due Thursday.