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

Lecture 13:

Latent Factor Models &

Non-Negative Matrix Factorization

# Announcements

- Homework 6 Released
  - Due Tuesday March 1$^{st}$

- Miniproject 2 to be released next week
  - Due March 8th

# Today

- Some useful matrix properties
  - Useful for homework

- Latent Factor Models
  - Low-rank models with missing values

- Non-negative matrix factorization

# Recap: Orthogonal Matrix

- ## A matrix U is orthogonal if $UU^T = U^TU = I$

    - For any column u:  $u^Tu = 1$

    - For any two columns u, u':  $u^Tu' = 0$

    - U is a rotation matrix, and $U^T$ is the inverse rotation

    - If $x' = U^Tx$, then $x = Ux'$

# Recap: Orthogonal Matrix

- Any subset of columns of U defines a subspace

$$x' = U_{1:K}^T x$$

Transform into new coordinates
Treat $U_{1:K}$ as new axes

$$proj_{U_{1:K}}(x) = U_{1:K} U_{1:K}^T x$$

Project x onto $U_{1:K}$ in original space
"Low Rank" Subspace

$u_1^T x$

$u_1 u_1^T x$

# Recap: Singular Value Decomposition

$$X = \left[ x_1, ..., x_N \right] \in \mathrm{Re}^{D \times N}$$

$$X = U \Sigma V^T$$

SVD

Orthogonal

Orthogonal    Diagonal

$$\sum_{i=1}^{N} \left\| x_i - U_{1:K} U_{1:K}^T x_i \right\|^2$$

"Residual"

**$U_{1:K}$ is the K-dim subspace with smallest residual**

# Recap: SVD & PCA

$$XX^T = U\Lambda U^T \qquad \qquad \textbf{PCA}$$

Orthogonal   Diagonal

$$X = U\Sigma V^T \qquad \qquad \textbf{SVD}$$

Orthogonal

Orthogonal   Diagonal

$$XX^T = \left(U\Sigma V^T\right)\left(U\Sigma V^T\right)^T = U\Sigma V^T V\Sigma U^T = U\Sigma^2 U^T$$

# Recap: Eigenfaces

- Each col of U' is an "Eigenface"
- Each col of V'ᵀ = coefficients of a student



**225000-dimensional!**

Avg Face

# Matrix Norms

- Frobenius Norm

$$\|X\|_{Fro} = \sqrt{\sum_{ij} X_{ij}^2} = \sqrt{\sum_d \sigma_d^2}$$

- Trace Norm

$$\|X\|_* = \sum_d \sigma_d = \text{trace}\left(\sqrt{X^T X}\right)$$

$$X = U\Sigma V^T$$

Each $\sigma_d$ is guaranteed to be non-negative
By convention: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_D \geq 0$

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_D \end{bmatrix}$$

# Properties of Matrix Norms

$$\|X\|_{Fro}^2 = \mathrm{trace}\left(X^T X\right) = \mathrm{trace}\left(\left(U\Sigma V^T\right)^T U\Sigma V^T\right)$$

$$= \mathrm{trace}\left(V\Sigma^2 V^T\right) = \mathrm{trace}\left(\Sigma^2 V^T V\right)$$

$$= \mathrm{trace}\left(\Sigma^2\right) = \sum_d \sigma_d^2$$

$$X = U\Sigma V^T$$

Each $\sigma_d$ is guaranteed to be non-negative
By convention: $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_D \geq 0$

$$\mathrm{trace}(ABC) = \mathrm{trace}(BCA) = \mathrm{trace}(CAB)$$

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_D \end{bmatrix}$$

# Properties of Matrix Norms

$$\left\| X \right\|_* = \text{trace}\left( \sqrt{\left( U\Sigma V^T \right)^T U\Sigma V^T} \right) = \text{trace}\left( \sqrt{V\Sigma U^T U\Sigma V^T} \right)$$

$$= \text{trace}\left( \sqrt{V\Sigma\Sigma V^T} \right) = \text{trace}\left( \sqrt{V\Sigma^2 V^T} \right) = \text{trace}\left( V\Sigma V^T \right)$$

$$= \text{trace}\left( \Sigma V^T V \right) = \text{trace}\left( \Sigma \right) = \sum_d \sigma_d$$

$$X = U\Sigma V^T$$

Each $\sigma_d$ is guaranteed to be non-negative
By convention: $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_D \geq 0$

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_D \end{bmatrix}$$

$$\text{trace}(ABC) = \text{trace}(BCA) = \text{trace}(CAB)$$

# Frobenius Norm = Squared Norm

- Matrix version of L2 Norm:

$$\|X\|_{Fro}^2 = \sum_{ij} X_{ij}^2 = \sum_d \sigma_d^2$$

- Useful for regularizing matrix models

$$X = U\Sigma V^T \qquad \Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_D \end{bmatrix}$$

# Recall: L1 & Sparsity

- w is sparse if mostly 0's:
  - Small L0 Norm

$$\|w\|_0 = \sum_d 1_{[w_d \neq 0]}$$

- Why not L0 Regularization?
  - **Not continuous!**

$$\underset{w}{\arg\min} \; \lambda \|w\|_0 + \sum_{i=1}^{N} \left( y_i - w^T x_i \right)^2$$

- L1 induces sparsity
  - And is continuous!

$$\underset{w}{\arg\min} \; \lambda |w| + \sum_{i=1}^{N} \left( y_i - w^T x_i \right)^2$$

Omitting b &
for simplicity

# Trace Norm = L1 of Eigenvalues

- A matrix X is considered low rank if it has few non-zero singular values:

$$\|X\|_{Rank} = \sum_d 1_{[\sigma_d > 0]}$$

**Not continuous!**

$$\|X\|_* = \sum_d \sigma_d = \text{trace}\left(\sqrt{X^T X}\right)$$

aka "spectral sparsity"

$$X = U\Sigma V^T \qquad \Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_D \end{bmatrix}$$

# Other Useful Properties

- Cauchy Schwarz:

$$\langle A, B \rangle^2 = \operatorname{trace}(A^T B)^2 \le \langle A, A \rangle \langle B, B \rangle = \operatorname{trace}(A^T A)\operatorname{trace}(B^T B) = \|A\|_F^2 \|B\|_F^2$$

- AM-GM Inequality:

$$\|A\|\|B\| = \sqrt{\|A\|^2 \|B\|^2} \le \frac{1}{2}\left(\|A\|^2 + \|B\|^2\right) \qquad \text{True for any norm}$$

- Orthogonal Transformation Invariance of Norms:

$$\|UA\|_F = \|A\|_F \qquad \|UA\|_* = \|A\|_* \qquad \text{If U is a full-rank orthogonal matrix}$$

- Trace Norm of Diagonals

$$\|A\|_* = \sum_i |A_{ii}| \qquad \text{If A is a square diagonal matrix}$$

# Recap: SVD & PCA

- SVD: $$X = U\Sigma V^T$$

- PCA: $$XX^T = U\Sigma^2 U^T$$

- The first K columns of U are the best rank-K subspace that minimizes the Frobenius norm residual:

$$\left\| X - U_{1:K}U_{1:K}^T X \right\|_{Fro}^2$$

# Latent Factor Models

# Netflix Problem



- $Y_{ij}$ = rating user i gives to movie j

- Solve using SVD!

$$y_{ij} \approx u_i^T v_j$$

# Example



$$y_{ij} \approx u_i^T v_j$$

**Miniproject 2: create your own.**

http://www2.research.att.com/~volinsky/papers/ieeecomputer.pdf

# Actual Netflix Problem



- ## Many missing values!

# Collaborative Filtering

- M Users, N Items
- Small subset of user/item pairs have ratings
- Most are missing

- Applicable to any user/item rating problem
  - Amazon, Pandora, etc.

- **Goal:** Predict the missing values.

# Latent Factor Formulation

- Only labels, no features $\qquad S = \left\{ y_{ij} \right\}$

- Learn a **latent** representation over users U and movies V such that:

$$\underset{U,V}{\mathrm{argmin}} \frac{\lambda}{2} \left( \|U\|_{Fro}^2 + \|V\|_{Fro}^2 \right) + \sum_{ij} \left( y_{ij} - u_i^T v_j \right)^2$$

# Connection to Trace Norm

- Suppose we consider all U,V that achieve perfect reconstruction: $Y = UV^T$

- Find U,V with lowest complexity:

$$\underset{Y=UV^T}{\mathrm{argmin}} \frac{1}{2} \left( \left\| U \right\|_{Fro}^2 + \left\| V \right\|_{Fro}^2 \right)$$

- Complexity equivalent to trace norm:

$$\left\| Y \right\|_* = \min_{Y=UV^T} \frac{1}{2} \left( \left\| U \right\|_{Fro}^2 + \left\| V \right\|_{Fro}^2 \right)$$   **Prove in homework!**

# Proof (One Direction)

We will prove: $\quad \|Y\|_* \geq \min_{Y=AB^T} \frac{1}{2}\left(\|A\|_{Fro}^2 + \|B\|_{Fro}^2\right) \qquad Y = U\Sigma V^T$

Choose: $\quad A = U\sqrt{\Sigma}, \qquad B = V\sqrt{\Sigma}$

Then:

$$\min_{Y=AB^T} \frac{1}{2}\left(\|A\|_{Fro}^2 + \|B\|_{Fro}^2\right) \leq \frac{1}{2}\left(\left\|U\sqrt{\Sigma}\right\|_{Fro}^2 + \left\|V\sqrt{\Sigma}\right\|_{Fro}^2\right)$$

$$= \frac{1}{2}\left(\mathrm{trace}\left(\left(U\sqrt{\Sigma}\right)^T\left(U\sqrt{\Sigma}\right)\right) + \mathrm{trace}\left(\left(V\sqrt{\Sigma}\right)^T\left(V\sqrt{\Sigma}\right)\right)\right)$$

$$= \frac{1}{2}\left(\mathrm{trace}\left(\sqrt{\Sigma}U^T U\sqrt{\Sigma}\right) + \mathrm{trace}\left(\sqrt{\Sigma}V^T V\sqrt{\Sigma}\right)\right)$$

$$= \frac{1}{2}\left(\mathrm{trace}\left(\sqrt{\Sigma}\sqrt{\Sigma}\right) + \mathrm{trace}\left(\sqrt{\Sigma}\sqrt{\Sigma}\right)\right)$$

$$= \frac{1}{2}\left(\mathrm{trace}\left(\Sigma\right) + \mathrm{trace}\left(\Sigma\right)\right) = \mathrm{trace}\left(\Sigma\right) = \|Y\|_*$$

# Interpreting Model

- Latent-Factor Model Objective

$$\operatorname*{argmin}_{U,V} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \sum_{ij}\left(y_{ij} - u_i^T v_j\right)^2$$

- Related to:

$$\operatorname*{argmin}_{W} \lambda\|W\|_* + \sum_{ij}\left(y_{ij} - w_{ij}\right)^2$$

**Find the best low-rank approximation to Y!**

$$\|W\|_* = \min_{W=UV^T} \frac{1}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right)$$  Equivalent when U,V = rank of W

# User/Movie Symmetry

$$\underset{U,V}{\text{argmin}} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \sum_{ij}\left(y_{ij} - u_i^T v_j\right)^2$$

- **If we knew V, then linear regression to learn U**
  - Treat V as features

- **If we knew U, then linear regression to learn V**
  - Treat U as features

# Optimization

$$\operatorname*{argmin}_{U,V} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \sum_{ij} \omega_{ij}\left(y_{ij} - u_i^T v_j\right)^2 \qquad \omega_{ij} \in \{0,1\}$$

- Only train over observed $y_{ij}$

- Two ways to Optimize
  - Gradient Descent
  - Alternating optimization
    - Closed Form (for each sub-problem)
  - Homework question

# Gradient Calculation

$$\underset{U,V}{\operatorname{argmin}} \frac{\lambda}{2} \left( \|U\|_{Fro}^2 + \|V\|_{Fro}^2 \right) + \frac{1}{2} \sum_{ij} \omega_{ij} \left( y_{ij} - u_i^T v_j \right)^2$$

$$\partial_{u_i} = \lambda u_i - \sum_j \omega_{ij} v_j \left( y_{ij} - u_i^T v_j \right)^T$$

Closed Form Solution (assuming V fixed):

$$u_i = \left( \lambda I_K + \sum_j \omega_{ij} v_j v_j^T \right)^{-1} \left( \sum_j \omega_{ij} y_{ij} v_j \right)$$

# Gradient Descent Options

- ## Stochastic Gradient Descent
  - Update all model parameters for single data point

- ## Alternating SGD:
  - Update a single column of parameters at a time

$$u_i = u_i - \eta \partial_{u_i}$$

$$\partial_{u_i} = \lambda u_i - v_j \left( y_{ij} - u_i^T v_j \right)$$

# Alternating Optimization

- Initialize U & V randomly

- Loop
  - Choose next $u_i$ or $v_j$
  - Solve optimally:

$$u_i = \left( \lambda I_K + \sum_j \omega_{ij} v_j v_j^T \right)^{-1} \left( \sum_j \omega_{ij} y_{ij} v_j \right)$$

  - (assuming all other variables fixed)

# Tradeoffs

- Alternating optimization much faster in terms of #iterations

  – But requires inverting a matrix:

$$u_i = \left( \lambda I_K + \sum_j \omega_{ij} v_j v_j^T \right)^{-1} \left( \sum_j \omega_{ij} y_{ij} v_j \right)$$

- Gradient descent faster for high-dim problems

  – Also allows for streaming data

$$u_i = u_i - \eta \partial_{u_i}$$

A scatter plot titled with axes "Factor vector 1" (horizontal) and "Factor vector 2" (vertical), both ranging from -1.5 to 1.5. Movies plotted: Freddy Got Fingered, Half Baked, Freddy vs. Jason, Road Trip, Kill Bill: Vol. 1, Natural Born Killers, Scarface, Julien Donkey-Boy, I Heart Huckabees, Punch-Drunk Love, The Royal Tenenbaums, Being John Malkovich, Lost in Translation, Belle de Jour, Citizen Kane, Annie Hall, 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.

http://www2.research.att.com/~volinsky/papers/ieeecomputer.pdf

# Recap: Collaborative Filtering

- **Goal:** predict every user/item rating

- **Challenge:** only a small subset observed

- **Assumption:** there exists a low-rank subspace that captures all the variability in describing different users and items

# Aside: Multitask Learning

- M Tasks: $S^m = \left\{ (x_i, y_i^m) \right\}_{i=1}^N$

$$\underset{W}{\text{argmin}} \frac{\lambda}{2} R(W) + \frac{1}{2} \sum_m \sum_i \left( y_i - w_m^T x_i \right)^2$$

↑
Regularizer

- Example: personalized recommender system
  – One task per user:

# How to Regularize?

$$\underset{W}{\operatorname{argmin}} \frac{\lambda}{2} R(W) + \frac{1}{2} \sum_{m} \sum_{i} \left( y_i - w_m^T x_i \right)^2 \qquad S^m = \left\{ (x_i, y_i^m) \right\}_{i=1}^{N}$$

- Standard L2 Norm:

$$\underset{W}{\operatorname{argmin}} \frac{\lambda}{2} \|W\|^2 + \sum_{m} \sum_{i} \left( y_i - w_m^T x_i \right)^2 = \sum_{m} \left[ \frac{\lambda}{2} \|w_m\|^2 + \sum_{i} \left( y_i - w_m^T x_i \right)^2 \right]$$

- Decomposes to independent tasks
  - For each task, learn D parameters

# How to Regularize?

$$\underset{W}{\operatorname{argmin}} \frac{\lambda}{2} R(W) + \frac{1}{2} \sum_m \sum_i \left( y_i - w_m^T x_i \right)^2 \qquad S^m = \left\{ (x_i, y_i^m) \right\}_{i=1}^N$$

- Trace Norm:

$$\underset{W}{\operatorname{argmin}} \frac{\lambda}{2} \|W\|_* + \sum_m \sum_i \left( y_i - w_m^T x_i \right)^2$$

- Induces W to have low rank across all task

# Recall: Trace Norm & Latent Factor Models

- Suppose we consider all U,V that achieve perfect reconstruction: W=UV$^T$

- Find U,V with lowest complexity:

$$\underset{W=UV^T}{\operatorname{argmin}} \frac{1}{2}\left( \|U\|_{Fro}^2 + \|V\|_{Fro}^2 \right)$$

- Claim: complexity equivalent to trace norm:

$$\|W\|_* = \min_{W=UV^T} \frac{1}{2}\left( \|U\|_{Fro}^2 + \|V\|_{Fro}^2 \right)$$

# How to Regularize?

$$\underset{W}{\text{argmin}} \, \frac{\lambda}{2} R(W) + \frac{1}{2} \sum_m \sum_i \left( y_i - w_m^T x_i \right)^2 \qquad S^m = \left\{ (x_i, y_i^m) \right\}_{i=1}^{N}$$

- Latent Factor Approach

$$\underset{U,V}{\text{argmin}} \, \frac{\lambda}{2} \left( \|U\|_{Fro}^2 + \|V\|_{Fro}^2 \right) + \frac{1}{2} \sum_m \sum_i \left( y_i - u_m^T V x_i \right)^2$$

- Learns a feature projection x' = Vx

- Learns a K dimensional model per task

# Tradeoff

- D*N parameters:

$$\underset{W}{\operatorname{argmin}} \sum_m \left[ \frac{\lambda}{2} \|w_m\|^2 + \frac{1}{2} \sum_i \left( y_i - w_m^T x_i \right)^2 \right]$$

- D*K + N*K parameters:

$$\underset{U,V}{\operatorname{argmin}} \frac{\lambda}{2} \left( \|U\|_{Fro}^2 + \|V\|_{Fro}^2 \right) + \frac{1}{2} \sum_m \sum_i \left( y_i - u_m^T V x_i \right)^2$$

  – Statistically more efficient
  – Great if low-rank assumption is a good one

# Multitask Learning

- M Tasks: $$S^m = \left\{ (x_i, y_i^m) \right\}_{i=1}^N$$

$$\operatorname*{argmin}_{U,V} \frac{\lambda}{2}\left( \|U\|_{Fro}^2 + \|V\|_{Fro}^2 \right) + \frac{1}{2}\sum_m \sum_i \left( y_i^m - u_m^T V x_i \right)^2$$

- Example: personalized recommender system
  - One task per user:
  - If x is topic feature representation
    - V is subspace of correlated topics
    - Projects multiple topics together

# Reduction to Collaborative Filtering

$$\underset{U,V}{\mathrm{argmin}}\ \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \frac{1}{2}\sum_m\sum_i\left(y_i^m - u_m^T V x_i\right)^2 \qquad S^m = \left\{(x_i, y_i^m)\right\}_{i=1}^N$$

- Suppose each $x_i$ is single indicator $x_i = e_i$

- Then: $\quad V x_i = v_i$

- Exactly Collaborative Filtering!

$$x_i = \begin{bmatrix} \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix}$$

$$\underset{U,V}{\mathrm{argmin}}\ \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \frac{1}{2}\sum_m\sum_i\left(y_i^m - u_m^T v_i\right)^2$$

# Latent Factor Multitask Learning vs Collaborative Filtering

$$\underset{U,V}{\arg\min} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \frac{1}{2}\sum_m\sum_i\left(y_i^m - u_m^T V x_i\right)^2$$

- Projects x into low-dimensional subspace Vx

- Learns low-dimensional model per task

$$\underset{U,V}{\arg\min} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \frac{1}{2}\sum_m\sum_i\left(y_i^m - u_m^T v_i\right)^2$$

- Creates low dimensional feature for each movie

- Learns low-dimensional model per user

# General Bilinear Models

$$\underset{U,V}{\arg\min} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \sum_i \left(y_i - z_i^T U^T V x_i\right)^2 \quad S = \left\{(x_i, z_{i,} y_i)\right\}$$

- Users described by features z

- Items described by features x

- Learn a projection of z and x into common low-dimensional space
  - Linear model in low dimensional space

# Why are Bilinear Models Useful?

$$\operatorname*{argmin}_{U,V} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \frac{1}{2}\sum_m\sum_i\left(y_i - u_m^T v_i\right)^2$$

U: MxK
V: NxK

$$\operatorname*{argmin}_{U,V} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \frac{1}{2}\sum_m\sum_i\left(y_i - u_m^T V x_i\right)^2$$

U: MxK
V: DxK

$$\operatorname*{argmin}_{U,V} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \frac{1}{2}\sum_i\left(y_i - z_i^T U^T V x_i\right)^2$$

U: FxK
V: DxK

$$S = \left\{(x_i, z_i, y_i)\right\}$$

# Story So Far: Latent Factor Models

$$\operatorname*{argmin}_{U,V} \frac{\lambda}{2}\left(\|U\|_{Fro}^2 + \|V\|_{Fro}^2\right) + \frac{1}{2}\sum_i \left(y_i - z_i^T U^T V x_i\right)^2 \quad S = \left\{(x_i, z_{i,}y_i)\right\}$$

- **Simplest Case:** reduces to SVD of matrix Y
  - No missing values
  - (z,x) indicator features

- **General Case:** projects high-dimensional feature representation into low-dimensional linear model

# Non-Negative Matrix Factorization

# Limitations of PCA & SVD

All features
non-negative

**PCA/SVD
Solution**

**Better Solution?**

# Non-Negative Matrix Factorization



- Assume Y is non-negative
- Find non-negative U & V

# CS 155 Non-Negative Face Basis

# CS 155 Eigenface Basis



-13.1664 , -27.5141

-25.3403 , -42.383

11.728 , 24.1556

16.6788 , 5.5092

3.2845 , -29.9722

8.5987 , 40.4183

-2.4102 , -20.5946

25.5004 , 29.0106

12.6747 , -13.4101

27.3545 , 1.2238

23.686 , -7.2213

-18.3888 , 6.9628

# Aside: Non-Orthogonal Projections

- If columns of A are not orthogonal, $A^TA \neq I$
  - How to reverse transformation $x' = A^Tx$?
  - **Solution: Pseudoinverse!**

$$A = U\Sigma V^T$$

SVD

$$A^+ = V\Sigma^+U^T$$

Pseudoinverse

**Intuition:** use the rank-K orthogonal basis that spans A.

$$A^{+T}A^Tx = U\Sigma^+V^TV\Sigma U^Tx$$

$$= U_{1:K}U_{1:K}^Tx$$

$$\Sigma^+ = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \sigma_D \end{bmatrix} \qquad \sigma^+ = \begin{cases} 1/\sigma & \text{if } \sigma > 0 \\ 0 & \text{otherwise} \end{cases}$$

# Objective Function

$$\underset{U \geq 0, V \geq 0}{\operatorname{argmin}} \sum_{ij} \ell(y_{ij}, u_i^T v_j)$$

- **Squared Loss:**
  - Penalizes squared distance

$$\ell(a, b) = (a - b)^2$$

- **Generalized Relative Entropy**
  - Aka, unnormalized KL divergence
  - Penalizes ratio

$$\ell(a, b) = a \log \frac{a}{b} - a + b$$

- **Train using gradient descent**

http://hebb.mit.edu/people/seung/papers/nmfconverge.pdf

# SVD/PCA vs NNMF

- **SVD/PCA:**
  - Finds the best orthogonal basis faces
    - Basis faces can be neg.
  - Coeffs can be negative
  - Often trickier to visualize
  - Better reconstructions with fewer basis faces
    - Basis faces capture the most variations

- **NNMF:**
  - Finds best set of non-negative basis faces
  - Non-negative coeffs
    - Often non-overlapping
  - Easier to visualize
  - Requires more basis faces for good reconstructions

# Non-Negative Latent Factor Models

$$\underset{U,V}{\operatorname{argmin}} \frac{\lambda}{2} \left( \|U\|_{Fro}^2 + \|V\|_{Fro}^2 \right) + \sum_i \ell \left( y_i, z_i^T U^T V x_i \right) \quad S = \left\{ (x_i, z_i, y_i) \right\}$$

- **Simplest Case:** reduces to NNMF of matrix Y
  - No missing values
  - (z,x) indicator features
- **General Case:** projects high-dimensional non-negative features into low-dimensional non-negative linear model

# Modeling NBA Gameplay Using Non-Negative Spatial Latent Factor Models

# Fine-Grained Spatial Models

- Discretize court
  - 1x1 foot cells
  - 2000 cells

- 1 weight per cell
  - 2000 weights

$$F_s(\mathbf{x}):$$



40 feet

# Fine-Grained Spatial Models

- Discretize court
  - 1x1 foot cells
  - 2000

- 1 weig
  - 2000 weights

**But most players haven't played that much!**

$$F_s(\mathbf{x}):$$

40 feet

Visualizing location factors L

Visualizing players $B_bL$

| Kawhi Leonard | Carmelo Anthony | Dirk Nowitzki | Dion Waiters | John Wall | Tim Duncan | Kyrie Irving | Shawn Marion | Jeremy Lin | David Lee |

http://www.yisongyue.com/publications/icdm2014_bball_predict.pdf

# Training Data



**STATS SportsVU**
2012/2013 Season, 630 Games,
80K Possessions, 380 frames per possession

# Prediction

- ## Game state: **x**
  - Coordinates of all players
  - Who is the ball handler

- ## Event: y
  - Ball handler will shoot
  - Ball handler will pass (to whom?)
  - Ball handler will hold onto the ball
  - 6 possibilities

- ## **Goal:** Learn P(y|**x**)



http://www.yisongyue.com/publications/icdm2014_bball_predict.pdf

# Logistic Regression
## (Simple Version: Just for Shooting)

$$P(y \mid \mathbf{x}) = \frac{\exp\{F(y \mid \mathbf{x})\}}{Z(\mathbf{x} \mid F)} \qquad Z(\mathbf{x} \mid F) = \sum_{y' \in \{s, \perp\}} \exp\{F(y' \mid \mathbf{x})\}$$

$$F(y' \mid \mathbf{x}) = \begin{cases} F_s(\mathbf{x}) & y' = s \quad \text{Shot} \\ F_\perp & y' = \perp \quad \text{Hold on to ball} \end{cases}$$

Offset or bias



D

K        D

L        K

M    $F_s$    =    M    $B^T$

Location Factors

Player Factors

$$P(y = s \mid \mathbf{x}) = \frac{1}{1 + \exp\{-F_s(\boldsymbol{x}) + F_\perp\}}$$

# Learning the Model

- Given training data: $\qquad S = \{(\boldsymbol{x}, y)\}$

Player Configuration

What Happened Next

- Learn parameters of model:

$$\underset{F_s, F_\perp}{\operatorname{argmin}} \frac{\lambda}{2} \|F_s\|^2 + \sum_{(\boldsymbol{x}, y) \in S} \ell\left(y, F_s(x) - F_\perp\right)$$

Log Loss

$$P(y = s \mid \mathbf{x}) = \frac{1}{1 + \exp\left\{-F_s(\boldsymbol{x}) + F_\perp\right\}}$$

# Optimization via Gradient Descent

$$\underset{B \geq 0, L \geq 0, F_\perp}{\mathrm{argmin}} \quad \frac{\lambda}{2}\left(\|B\|^2 + \|L\|^2\right) + \sum_{(\boldsymbol{x},y)} \ell\left(y, B_{b(\boldsymbol{x})}^T L_{l(\boldsymbol{x})} - F_\perp\right)$$

$$\partial_{L_i} = \lambda_1 L_i - \sum_{(\boldsymbol{x},y)} \frac{\partial \log P(y \mid \boldsymbol{x})}{\partial L_i}$$

$$\frac{\partial \log P(y \mid \boldsymbol{x})}{\partial L_i} = \left(1_{[y=s]} - P(s \mid \boldsymbol{x})\right) B_{b(\boldsymbol{x})}$$



http://www.yisongyue.com/publications/icdm2014_bball_predict.pdf

# Where are Players Likely to Receive Passes?



$$F^1_p = P^T \quad M$$

D, K, D, M, K

**Location Factors**

**Player Factors**

Enforce Non-Negativity
(Accuracy Worse)
(More Interpretable)



## Visualizing Location Factors M



| Tony Parker | Dirk Nowitzki | LeBron James | Monta Ellis | Manu Ginobili | David Lee | Jose Calderon | Chandler Parsons | Goran Dragic | Joachim Noah |

http://www.yisongyue.com/publications/icdm2014_bball_predict.pdf

# How do passes tend to flow?



$F_p^2 = Q_1^T Q_2$

Target Location Factors

Source Location Factors

$Q_1$

$Q_2$

http://www.yisongyue.com/publications/icdm2014_bball_predict.pdf

# How do passes tend to flow?



Target Location Factors

Source Location Factors

Passing From "X"

Passing To "X"

http://www.yisongyue.com/publications/icdm2014_bball_predict.pdf

# Tensor Latent Factor Models

# Tensor Factorization



$$Y = U \otimes V \otimes W$$

Y (Missing Values) with dimensions M, N, Q

V with dimensions N, K

W with dimensions K, Q

U with dimensions M

# Tri-Linear Model

$$\underset{U,V,W}{\mathrm{argmin}} \frac{\lambda}{2}\left(\|U\|^2_{Fro} + \|V\|^2_{Fro} + \|W\|^2_{Fro}\right) + \sum_i \ell\left(y_i, \left\langle U^T z_i, V^T x_i, W^T q_i \right\rangle\right)$$

- Prediction via 3-way dot product: $\quad \left\langle a,b,c \right\rangle = \sum_k a_k b_k c_k$
  - Related to Hadamard Product

- **Example:** online advertising
  - User profile z
  - Item description x
  - Query q

Solve using
Gradient Descent

# Summary: Latent Factor Models

- ## Learns a low-rank model of a matrix of observations Y
  - Dimensions of Y can have various semantics

- ## Can tolerate missing values in Y

- ## Can also use features

- ## Widely used in industry

# Next Week

- Embeddings

- Deep Learning

- Next Thursday: Recitation on Advanced Optimization Techniques