

Machine Learning & Data Mining CS/CNS/EE 155

Lecture 15: Recent Applications of Latent Factor Models

Today

- Miniproject 1 Reports Due
- Latent Spatial Models for Basketball Play Prediction
- Latent Tensor Models for Collaborative Clustering
- aka: Stuff Yisong Likes







http://projects.yisongyue.com/bballpredict/

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)

Interpretable





Logistic Regression (Simple Version: Just for Shooting)

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

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

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

Example



Tim Duncan



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

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

Fine-Grained Spatial Models

- Discretize court
 - 1x1 foot cells
 - 2000 cells

$$F_{s}(\mathbf{x}) = F_{s}^{T} \phi_{s}(\mathbf{x})$$

2000 dim coeff vector

Indicator feature vector (1 in cell that ball handler is) (0 in all other entries)





Story so Far: Spatial Logistic Regression

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

Probability of shooting log-linear
 In spatial cell feature representation



Training Data



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

Learning the Model

• Given training data:



Learn parameters of model:

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

$$\mathsf{Log Loss}$$

Spatial Regularization

- Self-defined feature vector
 Has spatial structure
- Expect F_s to vary smoothly

 But not enough training data
 Spatial Regularization:



$$\operatorname{argmin}_{F_{s},F_{\perp}} \lambda_{1} \left\| F_{s} \right\|^{2} + \lambda_{2} \sum_{i,j \in N(i)} \kappa_{ij} \left(F_{s,i} - F_{s,j} \right)^{2} + \sum_{(x,y) \in S} \ell \left(y, F_{s}^{T} \phi_{s}(x) - F_{\perp} \right)$$
$$\kappa_{ij} = \exp \left\{ -d(i,j)^{2} / \sigma \right\}$$

Lecture 15: Recent Applications of Latent Factor Models

Combining Several Ideas

- Multiclass prediction We've only seen shots so far
 - Predict different outcomes (e.g., shot, pass, etc)
- Spatial Model
 - Captures spatial structure
- Multitask prediction
 - One prediction task per player
- Interpretable Model
 - Low-dimensional representation
 - Easy to visualize

Discretize into fine-grained cells (regularize neighbors to be similar)

Talk about this next

Latent factor part Related to multitask prediction

Multitask Prediction

- One task per player
 - F_s is actually a matrix of coefficients
 - Each row is a player
- Input data x contains id of ball handler
 - One task per ball handler





Dirk Nowitski





Multitask Prediction

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

$$= \frac{1}{1 + \exp\{-(F_{s,b(x),l(x)} - F_\perp)\}}$$

$$\mathsf{F}_s$$

$$\mathsf{Ballhander ID}$$

$$\mathsf{Location}$$

$$\mathsf{D}=2000$$

Learning Objective:

Not Enough Training Data!

- Most players don't shoot that often
 - 2000 weights per player
- Cell discretization is very fine-grained
 Spatial regularization helps some
- How to share data across players?
 Latent Factor Model!



Latent Factor Model (Shooting)

Assume that players can be represented by K-dimensional vector Learn a common K-dimensional latent feature representation



• Shooting Score: $F_s(x) = B_{b(x)}^T L_{l(x)}$

New Prediction Model

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

= $\frac{1}{1 + \exp\{-(B_{b(x)}^T L_{l(x)} - F_{\perp})\}}$
Ballhander ID Location
$$F_s = M B^T \int_{\text{Location Factors}}^{\text{D}}$$

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



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



Visualizing location factors L (K=10)



Visualizing players

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

Training Objective

Spatial Regularization

 K_{ij}

$$\underset{B \ge 0, L \ge 0, F_{\perp}}{\operatorname{argmin}} \lambda_{1} \left(\left\| B \right\|^{2} + \left\| L \right\|^{2} \right) + \lambda_{2} \sum_{k} \sum_{i, j \in N(i)} \kappa_{ij} \left(L_{k,i} - L_{k,j} \right)^{2} + \sum_{(x,y) \in S} \ell \left(y, B_{b(x)}^{T} L_{l(x)} - F_{\perp} \right)$$

$$\operatorname{Log-Loss of Latent Factor Score}$$



$$S = \{(\boldsymbol{x}, \boldsymbol{y})\}$$
$$= \exp\{-d(i, j)^2 / \sigma\}$$

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

Optimization via Gradient Descent

$$\underset{B \ge 0, L \ge 0, F_{\perp}}{\operatorname{argmin}} \lambda_{1} \left(\left\| B \right\|^{2} + \left\| L \right\|^{2} \right) + \lambda_{2} \sum_{k} \sum_{i, j \in N(i)} \kappa_{ij} \left(L_{k,i} - L_{k,j} \right)^{2} + \sum_{(x,y) \in S} \ell \left(y, B_{b(x)}^{T} L_{l(x)} - F_{\perp} \right)$$

$$\partial_{L_i} = 2\left(\lambda_1 L_i + \lambda_2 \sum_{k} \sum_{j \in N(i)} \left(L_{k,j} - L_{k,i}\right)\right) - \sum_{(\boldsymbol{x}, y) \in S} \frac{\partial \log P(y \mid \boldsymbol{x})}{\partial L_i}$$

п

v

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

Lecture 15: Recent Applications of Latent Factor Models

Player Factors

Initialization

- Need to initialize B & L before doing gradient descent
 - Need to preserve spatial structure
 - Random initialization often doesn't work
- Train F_s first w/o factorization

 Then factorize F into B*L
 - Non-negative Matrix Factorization
 - Use as initialization



D

Μ

General Prediction Problem (Multiclass Logistic Regression)

$$P(y|\mathbf{x}) = \frac{1}{Z(\mathbf{x}|F)} \exp\{F(y|\mathbf{x})\}$$

Scoring Function

 $Z(\mathbf{x}|F) = \sum_{y' \in Y(\mathbf{x})} \exp\{F(y'|\mathbf{x})\}$

 $F(y|\mathbf{x}) = \begin{cases} F_s(\mathbf{x}) & \text{if } y = s \\ F_p(i, \mathbf{x}) & \text{if } y = p_i \\ F_{\perp}(\mathbf{x}) & \text{if } y = \perp \end{cases}$

All possible y's

Multiclass Classification

Partition Function

Shot Pass to teammate Hold on to ball

Latent Factor Model (Passing)



• Passing Score:
$$F_p(i, \mathbf{x}) = P_i^T L_{l(i, \mathbf{x})}$$

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

Where are Players Likely to Receive Passes?



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



Visualizing Location Factors M



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

Latent Factor Model (Passing) #2



Additive Decomposition:

• **Passing Score:**
$$F_p(i, x) = P_i^T L_{l(i,x)} + Q_{1,l(x)}^T Q_{2,l(i,x)}$$

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

How do passes tend to flow?





 Q_1



 Q_2

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

How do passes tend to flow?



Passing To "X"

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

Visualizing Defender Factors (Subtractive Decomposition)

.













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

Visualizing Time-of-Possession Factors



How long the ball handler has held on to the ball.

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

Recap: Latent Factor Spatial Models

- Multiclass prediction
- Spatial Model
 - Captures spatial structure
- Multitask prediction
- Interpretable Model

Gradient Descent

E.g., shot, pass, etc.

Discretize into fine-grained cells (regularize neighbors to be similar) (compose different spatial factors)

One task per player

Latent factor model Shared representation across players Non-negative coefficients

Requires good initialization (NNMF of full model)

Latent Collaborative Clustering

Motivation: Richer User Interfaces



Information Bottleneck

We can only learn as much as interface permits (We can only help user as much as interface permits)



"ReGroup: Interactive Machine Learning for On-Demand Group Creation in Social Networks" [Ameershi et al., CHI 2012] Pinterest



Ŧ±

Sweet Somethings

Unfollow Board

Send Board

73 Pins 101 F





Cinnamon Cake with Cinnamon-Cream Cheese Frosting - Recipes, Dinner Ideas, Healthy Recipes & Food Guide

Pinned from piarecipes.com





Tips For Perfect Chocolate Chip Cookies Say good-bye to flat and hard cookies! SweetLittleDluebird.com





from The Girl Who Ate Everything
 Cookies and Cream
 Cheesecakes

平士 7

Pinned from the-girl-who-ate-everything.com





Great visual aid on piping tips and their result for cake decorating. $\mp \pm 2$

Pinned from media-cache-ak0.pinimg.com

http://www.pinterest.com



"Apolo: Making Sense of Large Network Data by Combining Rich User Interaction and Machine Learning" [Chau et al., CHI 2011]

Demo (Goal Oriented Data Browsing)



Planning a vacation

Planning meals for the week

Literature review for grant proposal

Recap: Collaborative Filtering

- Training Data: N items, M users
- $Y = \{Y_1 \dots Y_M\}$



Recap: Collaborative Filtering

• Learning Goal: $F(m,i) \approx y_{mi}$

(+1 or -1)



•
$$F(m, \overline{\mathbf{F}}_{\text{Pails Garler}}) \approx +1$$

- F(m, ,) ≈ -1
- Prediction goal: F(m,i) for new items

Collaborative Clustering

- Training Data: N items, M users
- $Y = \{Y_1 \dots Y_M\}$



http://www.yisongyue.com/publications/www2014_collab_cluster.pdf

Collaborative Clustering

• Learning goal: $F(m,i,j) \approx y_{mij}$ (+1 or -1)



• Prediction goal: F(m,i,j) for new items

http://www.yisongyue.com/publications/www2014_collab_cluster.pdf

Recap: Latent Factor Models

- Collaborative Filtering: N items, M users
- Users rate items (y_{mi})



Prediction: F(m,i) = <u_m, x_i>

Learning for Collaborative Filtering



$$L(\mathbf{U},\mathbf{X}) = \sum_{m} \sum_{i \in Y_m} (F(m,i) - y_{mi})^2$$

$$F(m,i) = x_i^T u_m$$

Recap: Collaborative Filtering

	Les Invalides	Palais Garnier - Opera National de Paris	Promenade plantee,	River Seine	Musee Rodin	Napoleon's Tomb,	Luxembourg Gardens	Pont-Neuf
User 1		A	Ð					
User 2		Ð	Solution		P			
User 3	Ð						E)	
User 4				Ð				Ð
User 5			Ð				Ð	Ð

Can now predict missing values

$$F(m,i) = x_i^T u_m$$

Also known as matrix completion





Latent Collaborative Clustering

•
$$F(m,i,j) = x_i \circ u_m \circ x_j$$

• $F(m,i,j) = \bigcup_{Diag(u_m)} U_m \circ x_j$

$$F(m, i, j) = x_i^T \operatorname{Diag}(u_m) x_j$$
$$= \left\langle x_i, x_j \right\rangle_{\operatorname{Diag}(u_m)}$$

Diag(u_m) = user-specific transform
 Metric learning (i.e., Mahanalobis)

http://www.yisongyue.com/publications/www2014_collab_cluster.pdf

Actual model slightly more complicated



(optimization details in paper)

http://www.yisongyue.com/publications/www2014_collab_cluster.pdf



- Alternating least squares optimization
 (convex in **U** and each x, but not jointly)
- Iterative closed form solutions
 - (use ADMM to solve $U \ge 0$)

(optimization details in paper)

http://www.yisongyue.com/publications/www2014_collab_cluster.pdf

Relationship to Metric Learning

- Metric Learning Learns a Transformation U:
 - Such that transforming features x by U minimizes loss:

Normal Metric Learning: $\operatorname{argmin}_{\mathbf{U} \ge \mathbf{0}} \lambda_1 \| \mathbf{U} \|^2 + L(\mathbf{U}, \mathbf{X})$ $\begin{aligned} \text{Transformed Inner Product:} \\ \left\langle x_i, x_j \right\rangle_U = x_i^T U x_j \end{aligned}$

- Latent Collaborate Clustering Learns both U and Features x
 - Useful if there many tasks, and few examples per task
 - Or if don't trust raw features x



Data Collection

New User Study: Hypothetical trip to Paris

Mechanical Turk

Cluster interesting attractions

250 attractions (from TripAdvisor)

218 users18.7 items per user4.5 clusters per user

INTERACTIVE SURVEY: Categorizing Attractions in Paris

• Pretend you're learning and organizing information about Paris for a potential trip there.

Round 1/4, Part A: Tell Us What You Find Interesting in Paris!

- · Please inspect the attractions below.
- Select the attractions you find interesting by clicking on them.
- Selected attractions will turn red-bordered.
- You can right-click on attractions to view more information about them.
- When you are done, click on "Next" at the bottom.
- *NOTE* -- only meaningful results will be approved.



NEXT >>

Data Collection

New User Study: Hypothetical trip to Paris

Mechanical Turk

Cluster interesting attractions

250 attractions (from TripAdvisor)

218 users 18.7 items per user 4.5 clusters per user

INTERACTIVE SURVEY: Categorizing Attractions in Paris

• Pretend you're learning and organizing information about Paris for a potential trip there.

Round 1/4, Part B: Tell Us How You'd Categorize Them!

- Please categorize the attractions into groups that are meaningful to you.
- Drag each attraction into one of the blue shaded group regions.
- You may re-arrange existing groups by dragging attractions to new groups.
- · You can create as many groups as you like (by clicking "Add New Group").
- Please label each of the groups you use.
- When you are done, click on "Next" at the bottom.
- *NOTE* -- only meaningful results will be approved.

Add New Group Remove Empty Groups



Data Collection

"Very interesting - I feel like I learned a bit about Paris!"

"I have never been to Paris but now I want to go even more."

Cluster interesting

"We would just about sell our souls to get there! Thank you for letting me travel there vicariously through this HIT!!!" (HIT refers to "Human Intelligence Task")

"Great format. Have been to Paris several times but still learned about many great new places to visit"

1.5 clusters per user

Features for Attractions

2 sets of features

Wikipedia tf-idf (~3.5K features)

Mechanical Turk Tagging (right) (~40 features)

Used for baselines

Keyword Tagging Attractions in Paris!

- Please inspect the attraction below.
- SELECT ALL keywords that are appropriate for this attraction.
- Selected keywords will turn RED.
- The right pane below displays additional information (e.g., wikipedia page) for your convenience.



Place de la Madeleine

Incient Ruin	Palace / Mansion	= Sear
rchitecture	Performance	-
irt	Plaza / Open Area	LaM
Bridge	 Recreational 	Law
Cabaret	Relaxing / Leisure	
Cemetary	Religious	
Comedy	Scenic Nature	
Culture	Scenic Urban	
Dining	Scenic Water	
ountain	Shopping	
arden / Park	 Sightseeing 	
listorical	 Spa / Massage 	
arge Building	Sports	
femorial	Street	
fonument / Statue	Theater / Opera	
fuseum Art	Tour	L'église
luseum Other	Transportation	Madelein
lightlife	 Walking / Strolling 	formally,
Dutdoors	Zoo / Aquarium	comman
		The Mad



(Submit)

E
 E
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C
 C

Evaluation

• For each test user:

- Hold out some attractions
- Predict cluster membership (or new cluster)

• Feature-based baselines:

- Diagonal metric (per user)

[Wagstaff & Cardie, 2000] [Xing et al., 2002] [Schultz & Joachims, 2003] [Davis et al., 2007] [Parameswaran & Weinberger, 2010]

- Latent feature transform (significantly slower) [Blitzer & Weston, 2012]

Prediction Tasks



Predict cluster membership

- Existing cluster:



Cluster with highest average F(m,i,j). (conventional classification)





All clusters have negative F(m,i,j). (predicting novelty)



Hold 25% from each cluster (no empty clusters) Predict cluster membership (conventional classification)





Hold 50% of attractions at random (some clusters empty) Predict cluster membership (or new cluster)

Recap: Latent Collaborative Clustering

- Multitask Metric Learning & Feature Learning
 Trained on partial clusterings created by individual users
- Individual factors x hard to visualize
 Maybe easier if enforced non-negativity
- Maybe better served as an embedding model:

$$F(m,i,j) = \|x_i - x_j\|_{\text{Diag}(u_m)}^2$$

Recap: Latent Factor Models

- Great way to compactly represent data
 - Share across many tasks
- Can be very interpretable
 - At least the simpler versions
 - Tradeoff between interpretability and accuracy
- Particularly useful if you don't trust (or don't have) raw features
- Next Week: Deep Learning