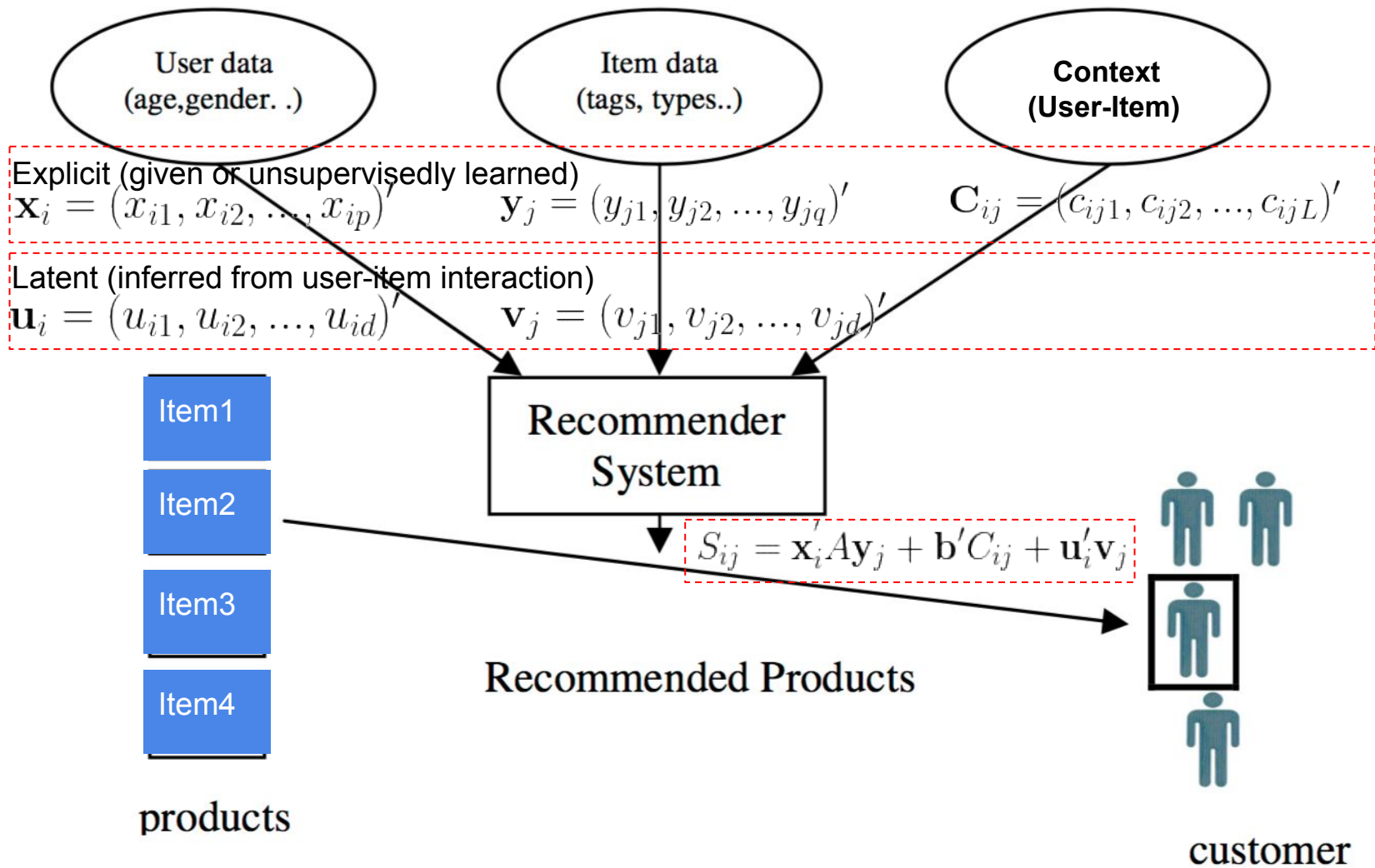


# A Statistical Modeling Framework for Developing Recommenders

David Yang  
xdyang70@gmail.com



# Recommender via Statistical Models

Using Rating Data with Gaussian Distribution for Illustration

Likelihood:

$$R_{ij} \sim \text{Normal}(S_{ij}, \sigma^2)$$
$$S_{ij} = \mathbf{x}_i' \mathbf{A} \mathbf{y}_j + \mathbf{b}' \mathbf{C}_{ij} + \mathbf{u}_i' \mathbf{v}_j$$

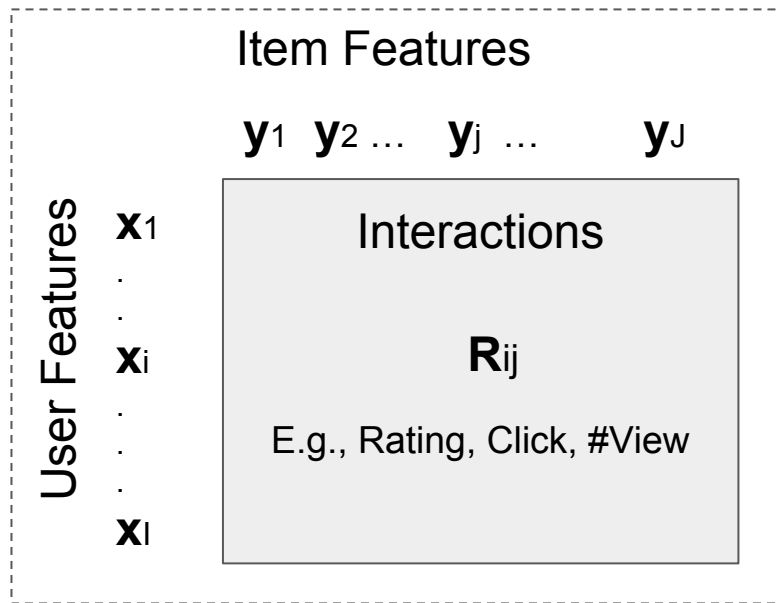
Prior Distribution:

$$Pr(\boldsymbol{\Theta}) = Pr(\{\mathbf{A}, \mathbf{b}, \mathbf{U}, \mathbf{V}, \sigma^2\})$$

# A Framework for Recommenders

Recommendations jointly made based on (Hybrid Solution)

- What we know about the items (Item-based)
- What we know about the user (User-based)
- How users interact with items (Collaborative)



Features

- User Attributes  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})'$
- Item Attributes  $\mathbf{y}_j = (y_{j1}, y_{j2}, \dots, y_{jq})'$

Interactions (Explicit: e.g., rating, Implicit: e.g., click):

- $\mathbf{R}_{ij} = (R_{ij1}, R_{ij2}, \dots, R_{ijK})'$   
[default  $K=1$ ]

Contexts of interactions (e.g., time, position):

- $\mathbf{C}_{ij} = (c_{ij1}, c_{ij2}, \dots, c_{ijL})'$

Scores - Modeling Result:

- $S_{ij} = f(\mathbf{x}_i, \mathbf{y}_j, \mathbf{c}_{ij}, \Theta)$

Modeling:

- Minimize  $\sum_{i,j} Dist(R_{ij}, S_{ij}) + \lambda r(\Theta)$

# Matrix-Vector Format (using Gaussian for illustration)

	$\mu_{11}, \dots, \mu_{1p}$	$\mu_{I1}, \dots, \mu_{Ip}$
$x_{11}, \dots, x_{1p}$	$R_{11}, \dots, R_{1J}$	
$x_{i1}, \dots, x_{ip}$	$R_{i1}, \dots, R_{iJ}$	
$x_{I1}, \dots, x_{Ip}$	$R_{I1}, \dots, R_{IJ}$	

$$R_{ij} = \underbrace{\mathbf{x}_i \mathbf{A} \mathbf{y}_j + \mathbf{b} \mathbf{c}_{ij}}_{\text{Fixed Effects}} + \underbrace{\mathbf{u}'_i \mathbf{v}_j + \epsilon_{ij}}_{\text{Random Effects}}$$

$$\mathbf{u}_i \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_{d \times d})$$

$$\mathbf{v}_j \sim N(\mathbf{0}, \sigma_v^2 \mathbf{I}_{d \times d})$$

$$\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$$

# Alternative Formulation

Predictors (I.V.)	DV
$x_{11}, \dots, x_{1p}; y_{11}, \dots, y_{1q}; x_{11}y_{11}, \dots, x_{1p}y_{1q}; c_{111}, \dots, c_{11l}; \dots$	$R_{11}$
.....	.
$x_{i1}, \dots, x_{ip}; y_{j1}, \dots, y_{jq}; x_{i1}y_{j1}, \dots, x_{ip}y_{jq}; c_{ij1}, \dots, c_{ijl}; \dots$	$R_{ij}$
.....	.
$x_{I1}, \dots, x_{Ip}; y_{J1}, \dots, y_{Jq}; x_{I1}y_{J1}, \dots, x_{Ip}y_{Jq}; c_{IJ1}, \dots, c_{IJl}; \dots$	$R_{IJ}$

$$R_{ij} = \underbrace{\beta_1 x_{11} + \dots + \beta_m c_{ijL}}_{\text{Fixed Effects}} + \underbrace{\mathbf{u}'_i \mathbf{v}_j + \epsilon_{ij}}_{\text{Random Effects}}$$

I.V. = Independent Variable;    DV = Dependent Variable

# Item/User Feature Vectors

Example of Item Features	$Y_j$		Example of User Features	$X_i$
Category: Business	0.0	Common Features	Interest: Business	1.0
Category: Entertainment	0.4		Interest: Entertainment	1.0
...	...		...	...
Category: Science	0.1		Interest: Science	0.0
Words: best	0.0		Words: best	0.3
Words: worst	0.2		Words: worst	0.1
...	...		...	...
Words: Surprise	0.3		Words: Surprise	0.1
Doc2Vec			Profile2Vec	
Vec_1	0.9		Vec_1	0.7
...	...		...	...
Vec_100	-0.5		Vec_100	-0.3
Other: Length	100		Other: #views	80
...	...		...	...
Other: Aging	20		Other: Demographics	30

# Approaches Creating Item Feature Vector: $Y_j$

- Categorization & Description  
Industry Coding; Man-made Tags; Sources; Location; Image/Audio Features
- Bag-of-Words:  
TF, TF-IDF; Phrase & Entities; Synonym Expansion; Dimension Reduction (Feature Selection with L1/L2 Regularization; Singular Value Decomposition; Random Projection)
- Topic Modeling  
Latent Dirichlet Allocation (LDA); Hierarchical Dirichlet Process (LDA); lda2Vec
- Semantics Modeling:  
Word2Vec; lda2Vec



# Approaches Creating User Feature Vector: $\mathbf{X}_i$

- Declared Profile  
Demographics; Declared/implied Interest (e.g., Netflix and StitchFix)
- Past Interaction with Content  
 $\mathbf{X}_i = F(\{ \mathbf{Y}_j: j \in \overline{\mathcal{I}}_i \})$  where  $\overline{\mathcal{I}}_i$  contains all items interacted by user  $i$ .  $F$  can be simple or weighted averaging (give higher weights to items recently interacted).
- Other User-related information  
Current Location; Usage-based features (web visits; devices used to access web); Search History; Item Set (similar to Bag of Words).

# User/Item-based Methods

## Unsupervised

- User-Item Similarity  
 $S(\mathbf{X}_i, \mathbf{Y}_j)$ :  $i = 1, \dots, I$ ;  $j = 1, \dots, J$ .
- Item-Item Similarities  
 $S(\mathbf{Y}_{j1}, \mathbf{Y}_{j2})$ :  $j1, j2 = 1, \dots, J$ .
- User-User Similarity  
 $S(\mathbf{X}_{i1}, \mathbf{X}_{i2})$ :  $i1, i2 = 1, \dots, I$ .

### Similarity:

- Cosine, Pearson/Spearman Correlation Coefficient, Okapi BM25, Jaccard.
- Weighted:  $S(\mathbf{X}_i, \mathbf{Y}_j) = \mathbf{X}_i' \mathbf{W} \mathbf{Y}_j$   
Problem: how to determine  $\mathbf{W}$ ?

## Supervised

$$S_{ij} = S(\mathbf{x}_i, \mathbf{y}_j) = \mathbf{x}_i' \mathbf{A} \mathbf{y}_j + \mathbf{b}' \mathbf{C}_{ij}$$

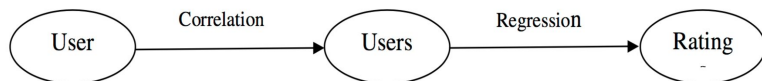
- Binary Ratings (logistic model)  
 $R_{ij} \sim \text{Bernoulli}((1 + \exp\{-S_{ij}\})^{-1})$
- Numerical Ratings (Gaussian Model)  
 $R_{ij} \sim \text{Normal}(S_{ij}, \sigma^2)$
- Ordinal Ratings (Cumulative Logit)  
 $R_{ij} \sim \text{Multinomial}(\pi_{ij1}, \dots, \pi_{ijK})$
- Pairwise Preference Scores  
 $R_{ijk} \sim \text{Bernoulli}((1 + \exp\{-(S_{ij} - S_{ik})\})^{-1})$
- Regularized M.L.E.  
$$\arg \max_{\mathbf{A}, \mathbf{\Theta}} (\log \Pr(R | \mathbf{A}, \mathbf{\Theta}) - \lambda r(\mathbf{A}))$$

# Collaborative Filtering

## User-User Similarity

$$S_{ij} = \bar{R}_i + \frac{\sum_{l \in I_j(i)} w(i, l)(R_{lj} - \bar{R}_l)}{\sum_{l \in I_j(i)} |w(i, l)|}$$

- Similarity Function (e.g., Pearson Correlation)
- Neighborhood Selection
- Weighting



## Item-Item Similarity



Matrix Factorization:  $S_{ij} = \mathbf{u}_i' \mathbf{v}_j$

$$\begin{matrix} \boxed{\mathbf{S}} \\ I \times J \end{matrix} = \begin{matrix} \boxed{\mathbf{U}} \\ I \times L \end{matrix} \times \begin{matrix} \boxed{\mathbf{V}} \\ L \times J \end{matrix}$$

( $L \ll I, L \ll J$ )

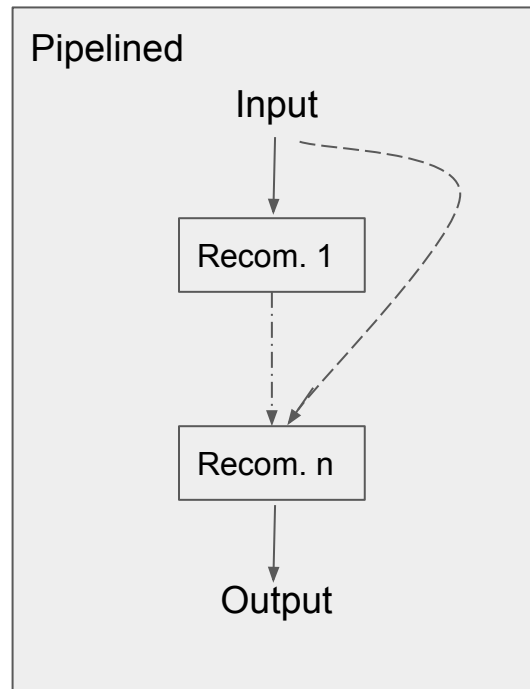
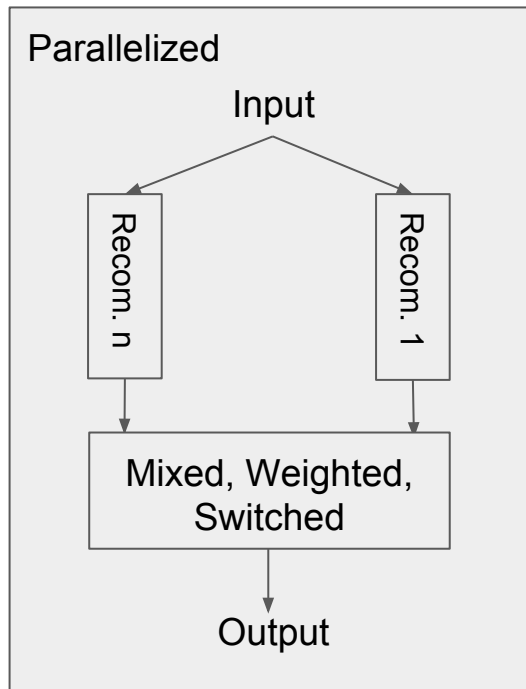
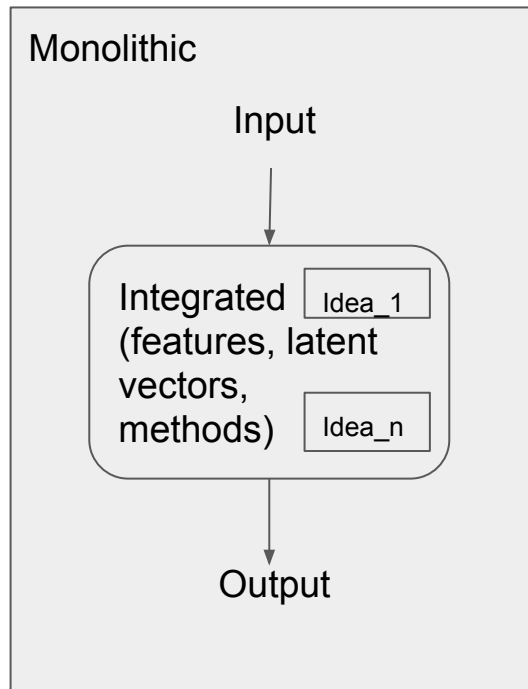
E.g., for numerical ratings

$$\arg \min_{\mathbf{u}_i, \mathbf{v}_j} \sum_{i,j} (R_{ij} - S_{ij})^2$$

## Optimization Methods

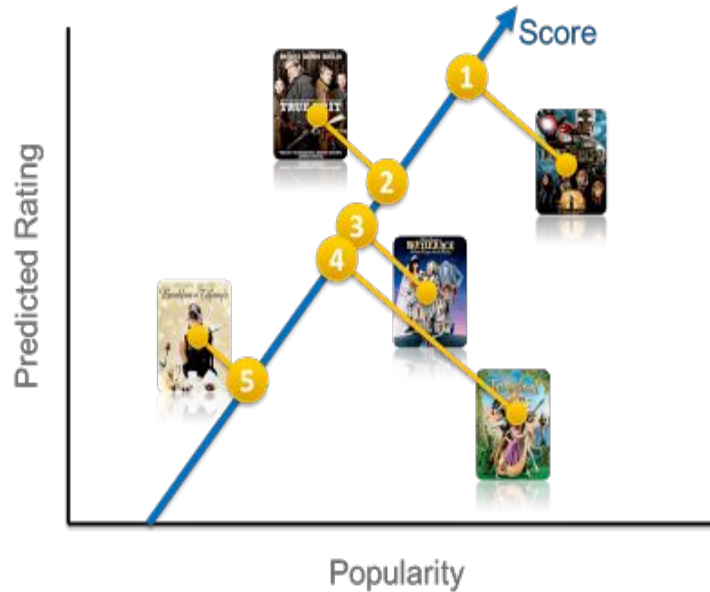
- Alternating Least Square
- Stochastic Gradient Descent

# Putting Them Together - Hybridization (Ensemble)

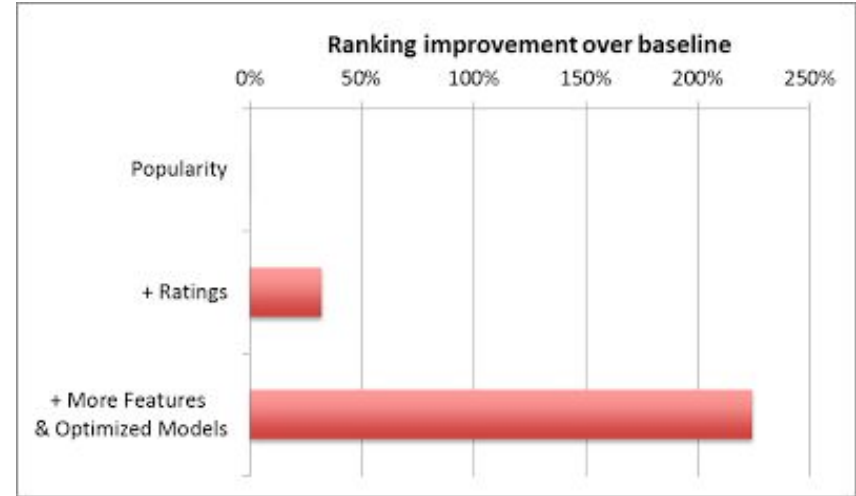


# More on Hybridizations

Weighted	The ratings of several recommendation techniques are combined together to produce a single recommendation
Switching	The system switches between recommendation techniques depending on the current situation
Mixed	Recommendations from several different recommenders are presented at the same time
Feature Combination	Features from different recommendation data sources are thrown together into a single recommendation algorithm
Cascade	One recommender refines the recommendations given by another
Feature Augmentation	Output from one technique is used as an input feature to another
Metal-Level	The model learned by one recommender is used as input to another



$f_{\text{rank}}(u,v) = w_1 p(v) + w_2 r(u,v) + b$ , where  $u$ =user,  $v$ =video item,  $p$ =popularity and  $r$ =predicted rating



A new feature does not show value because the model cannot learn it? Or, a more powerful model is not useful simply because you don't have the feature space that exploits its benefits?