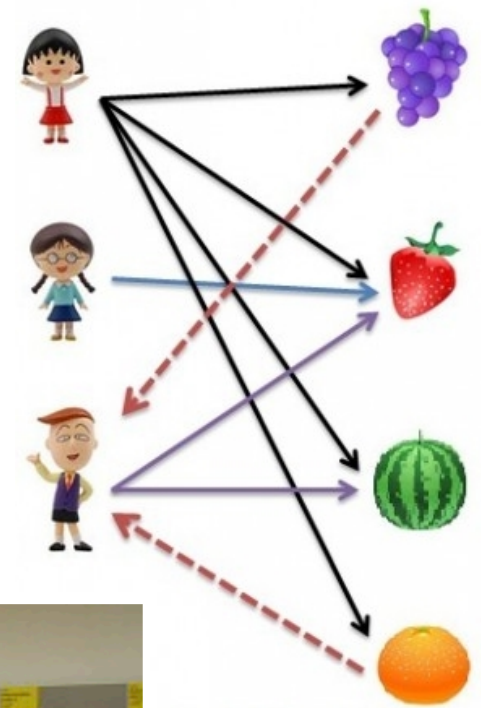


# **Lecture 5:**

# **Recommendation**

# Example

- Diapers and beer



# Examples

- Netflix Prize 2009



- Amazon Recommendation Engine

## Customers Who Bought This Item Also Bought



[30 Rock: Seasons 1-3](#)  
DVD ~ Tracy Morgan  
★★★★★ (7)  
\$60.49



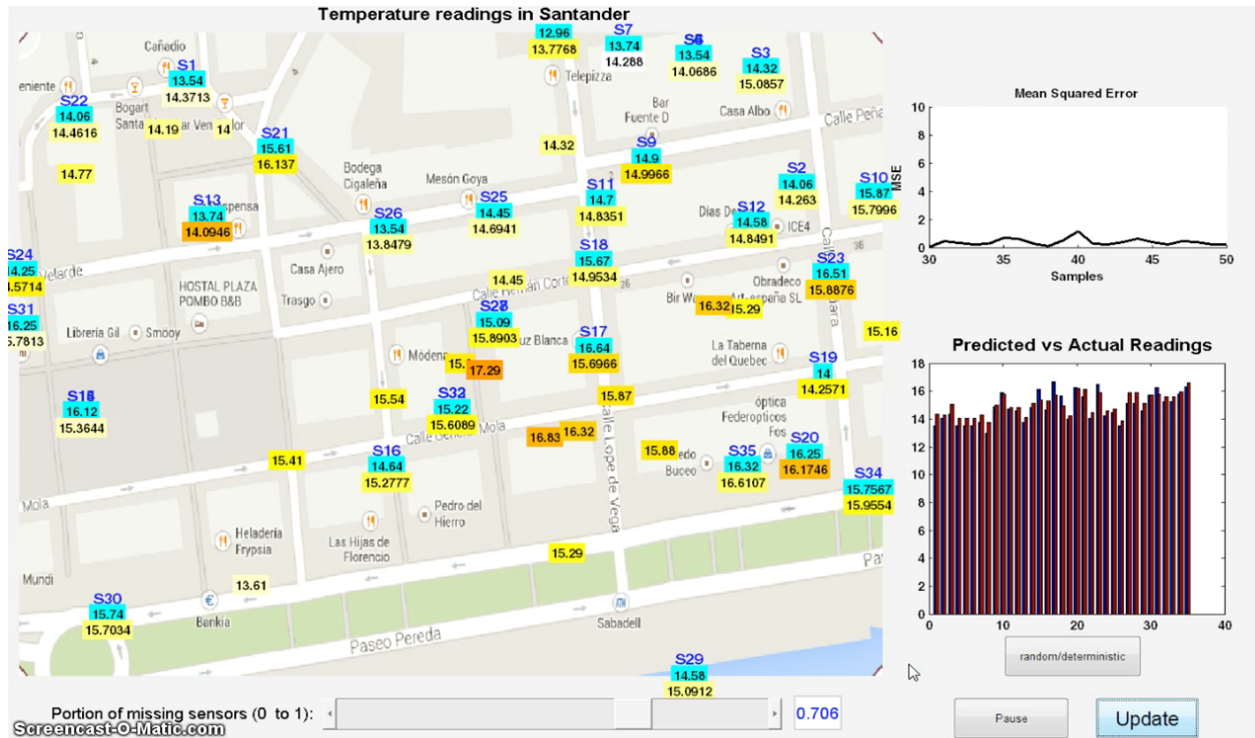
[Desperate Housewives: The Complete Seasons 1-5](#)  
DVD ~ Teri Hatcher  
★★★★☆ (2)  
\$179.99



[Scrubs: The Complete Seasons 1-8](#)  
DVD ~ Zach Braff  
★★★★★ (2)  
\$148.49

# Examples

- Missing sensor data



# Collaborative Filtering

- 400,000 users  
17,000 movies  
But only a few ratings (1%)
- User  $a$   
Movie  $i$   
Rating  $Y_{ai} \in \mathbb{R}$  (e.g. 1-5)
- Training data is the incomplete matrix  $Y$
- Goal is to predict unobserved ratings

$m$  movies

$n$  users

5	5						5		
		3	5	1	3	4	4		4
	4	2			2				
		5							5
4	5							4	
4							4		
5		4	5	1		4			
	4								
5				4					
5						4			
		5				5		3	

$Y_{ai}$

# Collaborative Filtering

- **Matrix or tensor completion problems**
- Collaborative: cross-users Filtering: prediction
- Dimensionality reduction

A tensor is a multi-dimensional array.  
e.g.  $n \times m$  2-tensor  
e.g.  $p \times q \times r$  3-tensor  
(2-tensor = matrix)

$m$  movies

5	5						5		
		3	5	1	3	4	4		4
	4	2			2				
		5							5
4	5							4	
4							4		
5		4	5	1		4			
	4								
5				4					
5						4			
		5				5		3	

$n$  users

$Y_{ai}$

Users who share similar interests on an item in the past are more likely to hold similar opinions on other items compared with a randomly chosen user.

# Types

**Model based** (Treat  $Y_{ai}$  as responses) – Matrix factorization

- Given training data of the form  $((a, i), Y_{ai})$ ,  
find a function  $f: \{1, \dots, n\} \times \{1, \dots, m\} \rightarrow \mathbb{R}$ .

**Memory based** (Treat  $Y_{ai}$  as features) – K nearest neighbors

- Given (incomplete) user ratings  $Y_a \in \mathbb{R}^m$ ,  
find structure in the data to predict missing values.

Recommendation is a *task* for which you can use either model based or memory based learning algorithms.

# K-Nearest Neighbors

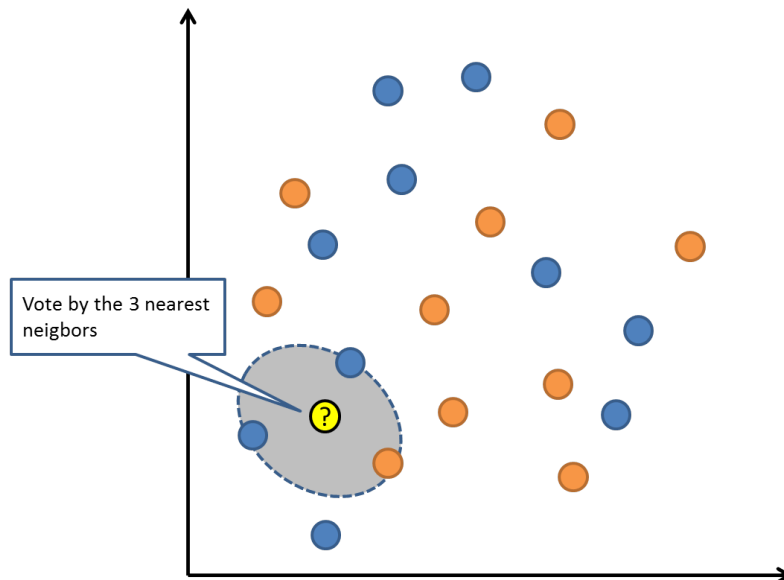
Memory based



# k-Nearest Neighbors

## Main Idea.

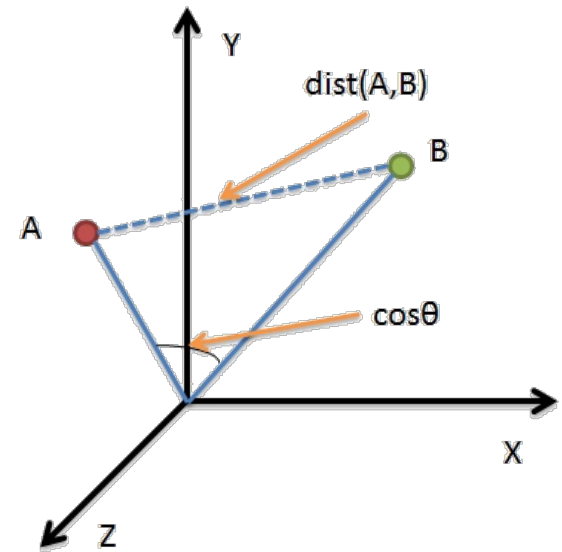
- Find **a few users**  $b_1, \dots, b_k$  (neighbors) that are similar to user  $a$
- Use information from users  $b_1, \dots, b_k$  to predict ratings of user  $a$



# Correlation Coefficient

## Cosine Similarity

$$\cos(x, y) = \frac{x^T y}{\|x\| \|y\|}$$



## Correlation Coefficient

$$\text{corr}(x, y) = \cos(x - \bar{x}, y - \bar{y}) = \frac{(x - \bar{x})^T (y - \bar{y})}{\|x - \bar{x}\| \|y - \bar{y}\|}$$

where  $\bar{x}, \bar{y}$  are the averages of the entries of  $x, y$ .

# User Similarity



CR = “Common Ratings”

To compute the similarity between users  $a$  and  $b$ ,

1. Find  $CR(a, b)$ , the set of movies rated by both  $a$  and  $b$ .
2. Let  $Z_a, Z_b$  be the vector of ratings in  $CR(a, b)$  for each user.
3. Compute the correlation coefficient between  $Z_a$  and  $Z_b$ .

$$\text{sim}(a, b) = \text{corr}(Z_a, Z_b) \in [-1, 1]$$

You could also subtract the average of  $Y_a$

# User Similarity

Movie	1	2	3	4	5	6
User a	4	5			1	
User b	2	4	1			2

$CR(a, b) = \text{Movie 1 and Movie 2}$

User a is represented by  $x = (4, 5)$ , User b is represented by  $y = (2, 4)$

$$x - \bar{x} = (4, 5) - 4.5 = (-0.5, 0.5)$$

$$y - \bar{y} = (2, 4) - 3 = (-1, 1)$$

$$\text{sim}(a, b) = \frac{\cos[(-0.5, 0.5), (-1, 1)]}{\sqrt{0.5^2 + 0.5^2} \sqrt{(-1)^2 + 1^2}} = \frac{0.5 + 0.5}{\sqrt{0.5} \sqrt{2}} = 1$$

# Weighted Prediction

Let  $\bar{Y}_b$  denote the average of movie ratings by user  $b$ .

To predict the rating  $Y_{ai}$  of user  $a$  for movie  $i$ ,

1. Rank users  $b$  who rated movie  $i$  according to value of  $\text{sim}(a, b)$ .
2. Let  $\text{kNN}(a, i)$  be the set of highest  $k$  users.
3. Let  $(Y_{ai} - \bar{Y}_a)$  be weighted average of  $(Y_{bi} - \bar{Y}_b)$ ,  $b \in \text{kNN}(a, i)$ .
4. Let weight for  $(Y_{bi} - \bar{Y}_b)$  be proportional to  $\text{sim}(a, b)$ .

# Anti-Correlated Users

If  $(Y_{bi} - \bar{Y}_b)$  is strongly anti-correlated with  $(Y_{ai} - \bar{Y}_a)$ , i.e.

$$\text{sim}(a, b) \ll 0,$$

then  $-(Y_{bi} - \bar{Y}_b)$  is strongly correlated with  $(Y_{ai} - \bar{Y}_a)$ .

We should exploit this information because ratings are rare.

# Weighted Prediction

Let  $\bar{Y}_b$  denote the average of movie ratings by user  $b$ .

To predict the rating  $Y_{ai}$  of a user  $a$  for a movie  $i$ ,

1. Rank users  $b$  who rated movie  $i$  according to value of  $|\text{sim}(a, b)|$ .
2. Let  $\text{kNN}(a, i)$  be the set of highest  $k$  users.
3. Let  $(Y_{ai} - \bar{Y}_a)$  be weighted average of  $\pm(Y_{bi} - \bar{Y}_b)$ ,  $b \in \text{kNN}(a, i)$ .
4. Let weight for  $\pm(Y_{bi} - \bar{Y}_b)$  be proportional to  $|\text{sim}(a, b)|$ .

$$\hat{Y}_{ai} - \bar{Y}_a = \frac{\sum_{b \in \text{kNN}(a, i)} \text{sim}(a, b)(Y_{bi} - \bar{Y}_b)}{\sum_{b \in \text{kNN}(a, i)} |\text{sim}(a, b)|}$$

# Discussion

$$\hat{Y}_{ai} - \bar{Y}_a = \frac{\sum_{b \in \text{kNN}(a,i)} \text{sim}(a, b)(Y_{bi} - \bar{Y}_b)}{\sum_{b \in \text{kNN}(a,i)} |\text{sim}(a, b)|}$$

- Formula is not sensitive to the bias (mean) of each user, but it is sensitive to the spread (variance) of each user.
- There is no training loss, no training algorithm for kNN.



# Matrix Factorization

Model based

# Subspace Learning

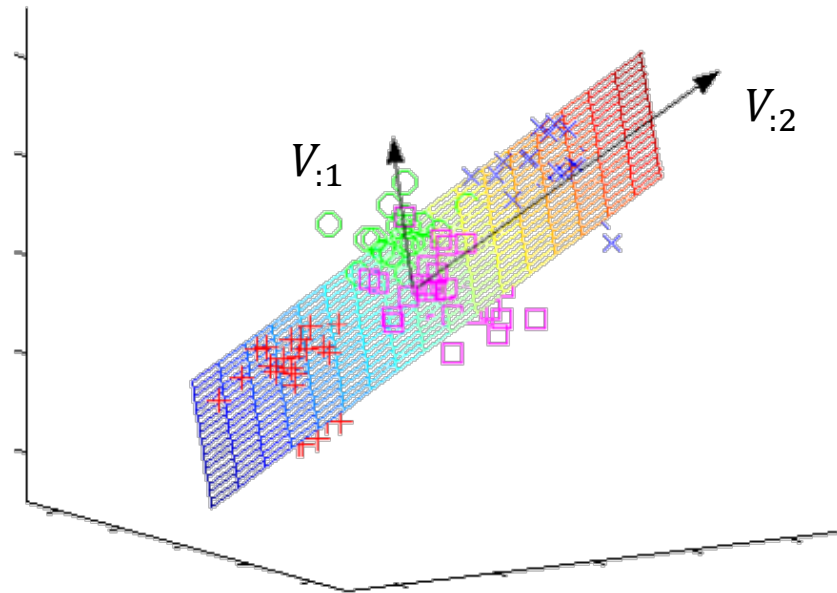
$k \ll m$  means  $k$  is  
much smaller than  $m$

## Main idea.

- *Completed* rating vectors  $\hat{Y}_1, \dots, \hat{Y}_n \in \mathbb{R}^m$  lie in some  $k$ -dimensional subspace,  $k \ll m$ .
- i.e. there exists vectors  $V_{:1}, \dots, V_{:k} \in \mathbb{R}^m$  such that for all users  $a$ ,

$$\hat{Y}_a = U_{a1} V_{:1} + U_{a2} V_{:2} + \dots + U_{ak} V_{:k}$$

for some coefficients  $U_{a1}, \dots, U_{ak} \in \mathbb{R}$ .



# Matrix Factorization

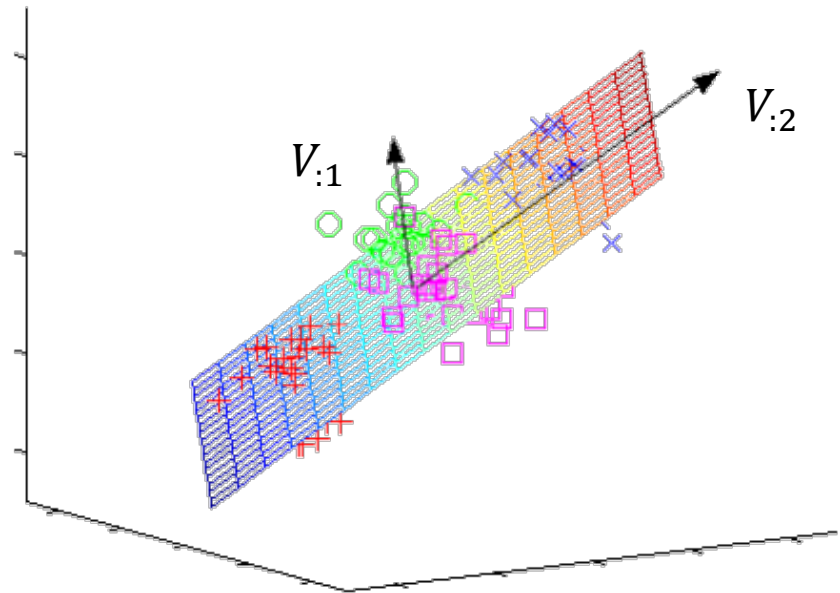
$$\begin{aligned}\hat{Y}_1^\top &= U_{11} V_{:1}^\top + \cdots + U_{1k} V_{:k}^\top \\ \hat{Y}_2^\top &= U_{21} V_{:1}^\top + \cdots + U_{2k} V_{:k}^\top \\ &\vdots \\ \hat{Y}_n^\top &= U_{n1} V_{:1}^\top + \cdots + U_{nk} V_{:k}^\top\end{aligned} \quad \Rightarrow \quad \begin{aligned}\hat{Y} &= U_{:1} V_{:1}^\top + \cdots + U_{:k} V_{:k}^\top \\ &= (U_{:1}, \dots, U_{:k}) \begin{pmatrix} V_{:1}^\top \\ \vdots \\ V_{:k}^\top \end{pmatrix} \\ &= (U_{:1}, \dots, U_{:k})(V_{:1}, \dots, V_{:k})^\top \\ &= UV^\top\end{aligned}$$

where  $U_{:1}, \dots, U_{:k}$  are the columns of  $U \in \mathbb{R}^{n \times k}$ ,  
and  $V_{:1}, \dots, V_{:k}$  are the columns of  $V \in \mathbb{R}^{m \times k}$ .

# Low-Rank Approximation

## Main Idea.

- Find low-rank *complete* matrix  $\hat{Y}$  that is closest to *incomplete* matrix  $Y$ .
- $\hat{Y}$  is called a *low-rank approximation* of  $Y$ .



# Overall idea – movie recommendation

- $\hat{Y}_{ai} = U_a^\top V_i$
- When watching movies, user  $a$  has preferences on
  - Comedy, sci-fi, action, romance ( $k$  factors)
  - $U_a \in \mathbb{R}^k$  characterizes user  $a$ 's preference on these  $k$  factors
- Movie  $i$  can also be described by
  - Comedy, sci-fi, action, romance ( $k$  factors)
  - $V_i \in \mathbb{R}^k$  describes how the  $k$  factors distribute in movie  $i$
- When we do the multiplication  $U_a^\top V_i$ , we obtain how much user  $a$  like movie  $i$

# Prediction

For an unknown rating  $Y_{ai}$  of user  $a$  for movie  $i$ , we predict

$$\hat{Y}_{ai} = (UV^\top)_{ai} = U_a V_i^\top$$

where  $U_a$  is the  $a$ -th row of  $U$  and  $V_i$  is the  $i$ -th row of  $V$ .

$U_a$ : factor vector of user  $a$ .

$V_i$ : factor vector of movie  $i$ .

# Training Loss

We regularize  
because there  
are infinitely  
many solutions  
 $U' = cU, V' = c^{-1}V$

Apply squared loss to each observed rating  $Y_{ai}$  .

$$\begin{aligned}\mathcal{L}_{n,k,\lambda}(U, V; Y) &= \sum_{(a,i) \in D} \frac{1}{2} (Y_{ai} - (UV^\top)_{ai})^2 + \frac{\lambda}{2} \|U\|^2 + \frac{\lambda}{2} \|V\|^2 \\ &= \sum_{(a,i) \in D} \frac{1}{2} (Y_{ai} - U_a^\top V_i)^2 + \frac{\lambda}{2} \sum_a \|U_a\|^2 + \frac{\lambda}{2} \sum_i \|V_i\|^2\end{aligned}$$

where  $D$  is the set of  $(a, i)$  where  $Y_{ai}$  is observed.

**Definition.** The *Frobenius* norm  $\|U\|$  of a matrix  $U \in \mathbb{R}^{n \times k}$  is

$$\|U\| = \sqrt{\sum_{a=1}^n \sum_{b=1}^k U_{ab}^2} .$$

# Alternating Least Squares

**Coordinate Descent** (optimization).

Repeat until convergence:

1. Fix  $V$  and minimize  $\mathcal{L}_{n,k,\lambda}(U, V; Y)$  over  $U$ .
2. Fix  $U$  and minimize  $\mathcal{L}_{n,k,\lambda}(U, V; Y)$  over  $V$ .



# Alternating Least Squares

1. Initialize  $V_1, V_2, \dots, V_m \in \mathbb{R}^k$  randomly.
2. Repeat until convergence:

- a. For each user  $a$ , find  $U_a$  that minimizes

$$\sum_{i: (a,i) \in D} \frac{1}{2} (Y_{ai} - U_a^\top V_i)^2 + \frac{\lambda}{2} \|U_a\|^2$$

- b. For each movie  $i$ , find  $V_i$  that minimizes

$$\sum_{a: (a,i) \in D} \frac{1}{2} (Y_{ai} - U_a^\top V_i)^2 + \frac{\lambda}{2} \|V_i\|^2$$

These are  
standard  
linear  
regression  
problems.

# Discussion

## User Bias

- Do we subtract the average ratings of each user?

## Optimization.

- Like k-means, the algorithm converges to a local minimum.
- Perform multiple initializations, and pick best result.

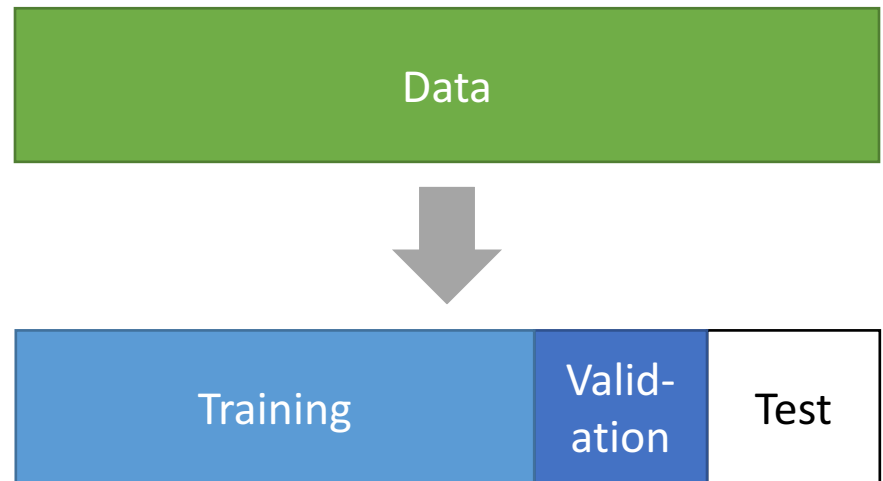
## Generalization.

- Use **validation** to pick right **hyperparameters**  $k$  and  $\lambda$ .

# Validation Set

Split the data into

- **Test set**  $\mathcal{S}_*$   
For evaluating, reporting performance at the end
- **Training set**  $\mathcal{S}_n$   
For training optimal parameters in a model
- **Validation set**  $\mathcal{S}_{\text{val}}$   
For model selection, e.g. picking  $k$  in k-means, picking  $\lambda$  in ridge regression. Acts as a proxy for test set.



# Validation Loss

The *validation error* is the test loss applied to the validation set.

The *training error* is the test loss applied to the training set, and it may be different from the training loss used for optimization.

**Example.** Ridge Regression

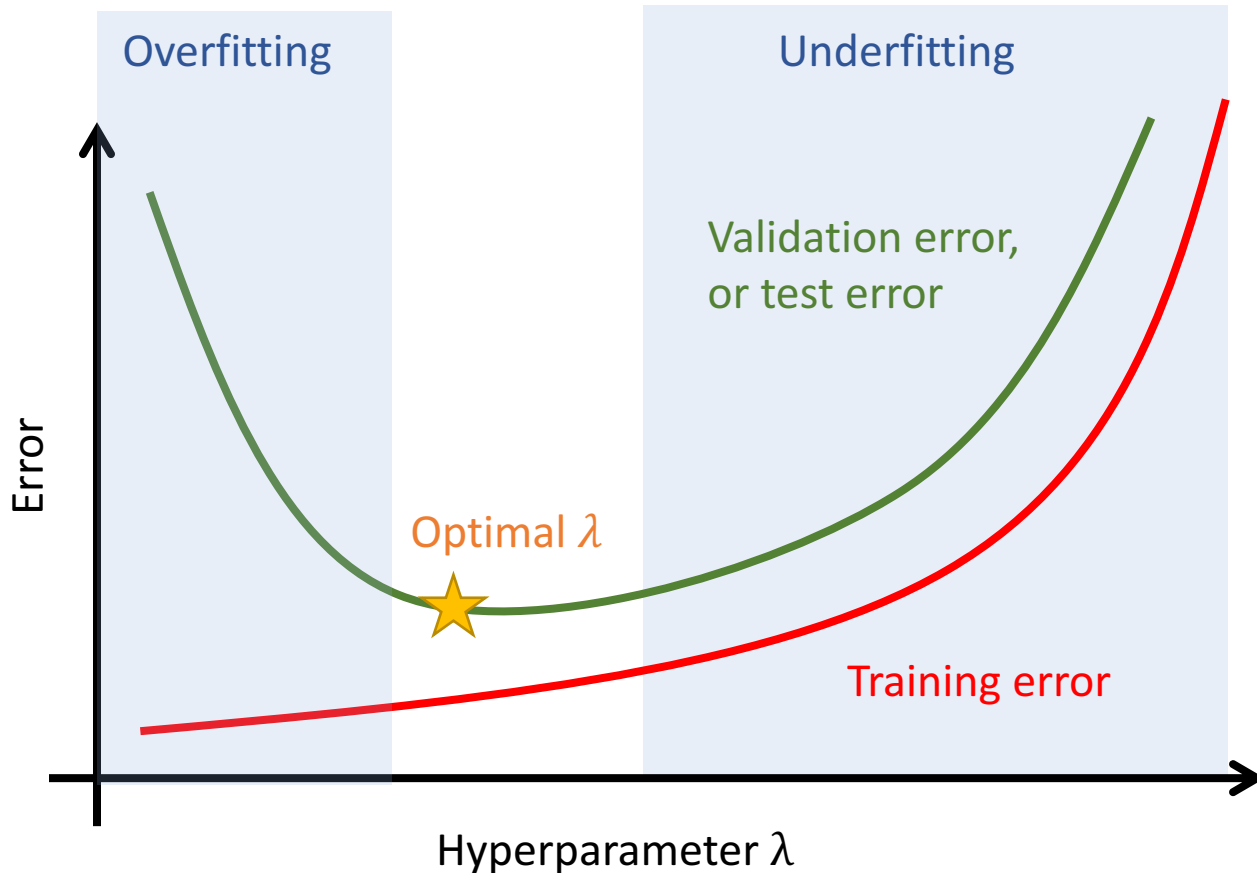
Test loss/error  $\mathcal{R}(\hat{\theta}; \mathcal{S}_*) = \frac{1}{n} \sum_{(x,y) \in \mathcal{S}_*} \frac{1}{2} (y - \hat{\theta}^\top x)^2$

Validation loss/error  $\mathcal{R}(\hat{\theta}; \mathcal{S}_{\text{val}}) = \frac{1}{n} \sum_{(x,y) \in \mathcal{S}_{\text{val}}} \frac{1}{2} (y - \hat{\theta}^\top x)^2$

Training error  $\mathcal{R}(\hat{\theta}; \mathcal{S}_n) = \frac{1}{n} \sum_{(x,y) \in \mathcal{S}_n} \frac{1}{2} (y - \hat{\theta}^\top x)^2$

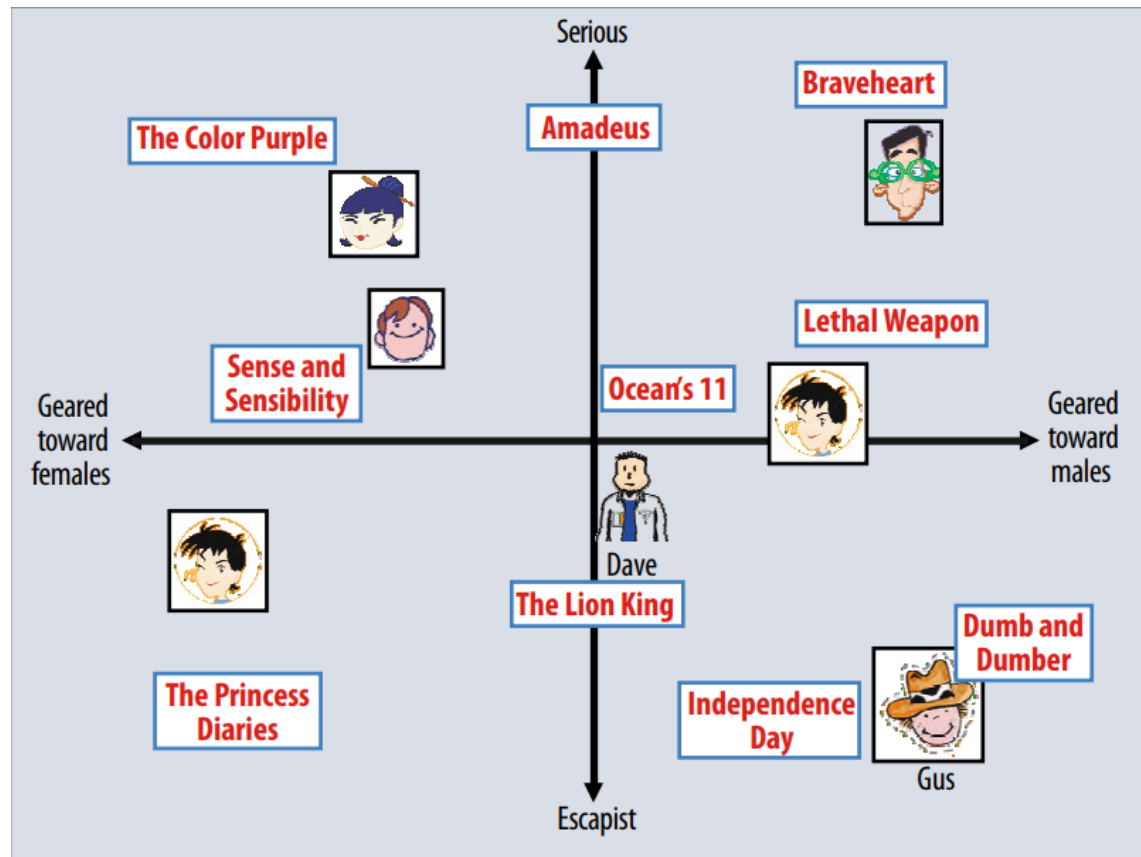
Training loss  $\mathcal{L}_{n,\lambda}(\theta; \mathcal{S}_n) = \frac{1}{n} \sum_{\text{data } (x,y)} \frac{1}{2} (y - \theta^\top x)^2 + \frac{\lambda}{2} \|\theta\|^2$

# Model Selection



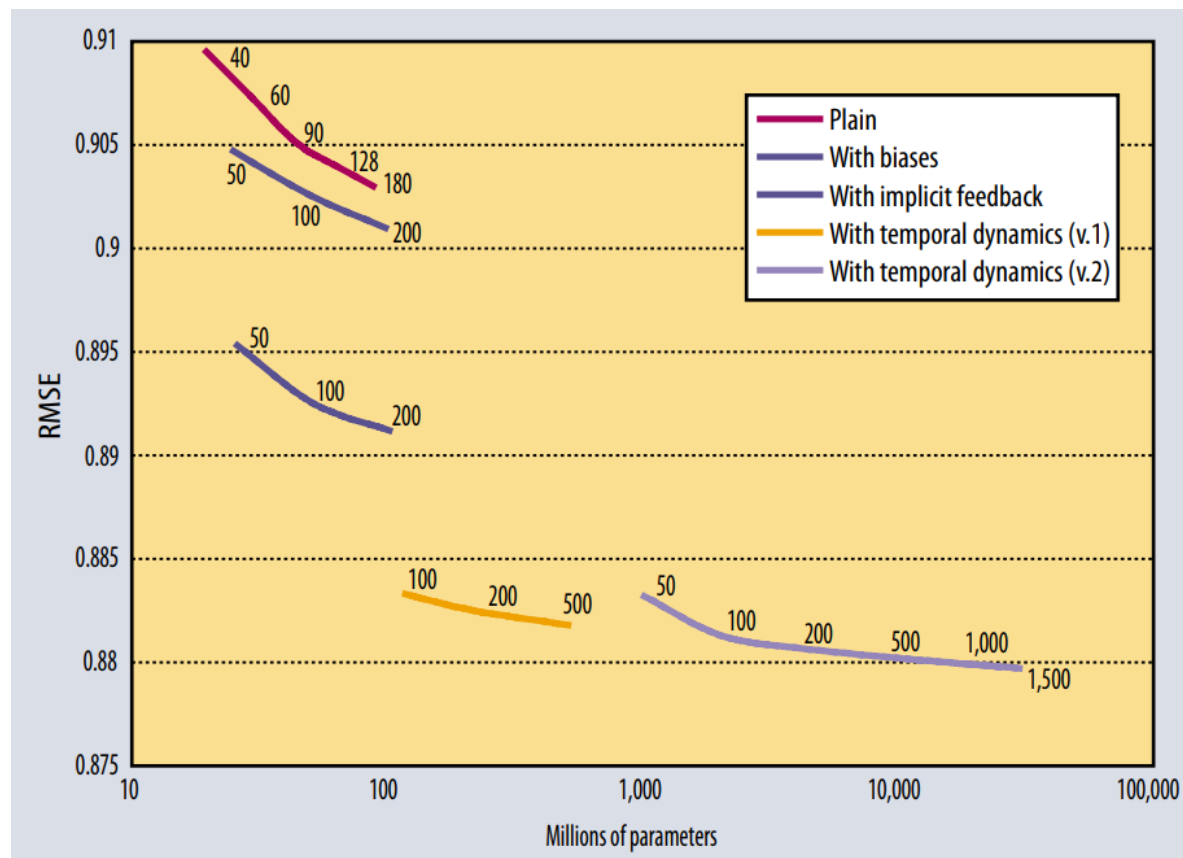
# Netflix Results

Plot each movie  $i$  according to first two coordinates of  $V_i \in \mathbb{R}^k$ , i.e.  $V_{i1}$ ,  $V_{i2}$ .



# Netflix Results

Winning team  
*Bellkor's  
Pragmatic Chaos*  
essentially used  
a matrix  
factorization  
model with  
additional  
parameters  
for biases,  
user attributes  
and time  
dependencies.



# Summary

- Matrix Completion
  - Recommender Systems
  - Collaborative Filtering
- $k$ -Nearest Neighbors
  - Correlation Coefficient
  - Weighted Prediction
- Matrix Factorization
  - Subspace Learning
  - Low-rank Approximation
  - Regularized Training Loss
  - Alternating Least Squares



# Intended Learning Outcomes

## Matrix Completion

- Recognize that the key problem in collaborative filtering is matrix completion. Give examples of collaborative filtering.
- Explain how dimensionality reduction helps matrix completion.

## k-Nearest Neighbors

- Compute the user similarity using the correlation coefficient.
- Compute the k-nearest neighbors ranked by user similarity.
- Predict an unknown rating using weighted prediction.

# Intended Learning Outcomes

## Matrix Factorization

- Describe the relationship between subspace learning, low-rank approximation and matrix factorization.
- Write down the regularized training loss. Explain why regularization is necessary for matrix factorization.
- Describe the alternating least squares algorithm, and explain why it minimizes the training loss. Apply the algorithm in a recommendation problem.
- Describe how validation can be used to select suitable hyperparameters  $k$  and  $\lambda$  for good generalization.