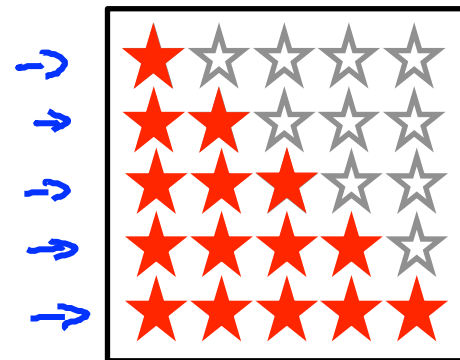# Recommender Systems

## Problem formulation

Machine Learning

# Example: Predicting movie ratings

→ User rates movies using ~~one~~ to five stars

zero

| Movie | Alice (1) | Bob (2) | Carol (3) | Dave (4) |
|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 |
| Romance forever | 5 | ? 4.5 | ? 0 | 0 |
| Cute puppies of love | ? 5 | 4 | 0 | ? 0 |
| Nonstop car chases | 0 | 0 | 5 | 4 |
| Swords vs. karate | 0 | 0 | 5 | ? 4 |

$n_u = 4$          $n_m = 5$

$n_u$ = no. users

$n_m$ = no. movies

$r(i, j)$ = 1 if user $j$ has rated movie $i$

$y^{(i,j)}$ = rating given by user $j$ to movie $i$ (defined only if $r(i, j) = 1$)

$0, ..., 5$

Andrew Ng

# Recommender Systems

## Content-based recommendations

Machine Learning

# Content-based recommender systems

$n_u = 4, \quad n_m = 5$

$x_0 = 1$

$$x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$$

| Movie | Alice (1) $\theta^{(1)}$ | Bob (2) $\theta^{(2)}$ | Carol (3) $\theta^{(3)}$ | Dave (4) $\theta^{(4)}$ |
|---|---|---|---|---|
| Love at last  1 | 5 | 5 | 0 | 0 |
| Romance forever  2 | 5 | ? | ? | 0 |
| Cute puppies of love  3 | ? 4.95 | 4 | 0 | ? |
| Nonstop car chases  4 | 0 | 0 | 5 | 4 |
| Swords vs. karate  5 | 0 | 0 | 5 | ? |

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$

$n = 2$

For each user $j$, learn a parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user $j$ as rating movie $i$ with $(\theta^{(j)})^T x^{(i)}$ stars.

$\theta^{(j)} \in \mathbb{R}^{n+1}$

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \longleftrightarrow \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \qquad (\theta^{(1)})^T x^{(3)} = 5 \times 0.99$$
$$= 4.95$$

**Problem formulation**

→ $r(i, j) = 1$ if user $j$ has rated movie $i$ (0 otherwise)

→ $y^{(i,j)} =$ rating by user $j$ on movie $i$ (if defined)

→ $\theta^{(j)}$ = parameter vector for user $j$

→ $x^{(i)}$ = feature vector for movie $i$

→ For user $j$, movie $i$, predicted rating: $\underline{\underline{(\theta^{(j)})^T(x^{(i)})}}$     $\Theta^{(j)} \in \mathbb{R}^{n+1}$

→ $\underline{m^{(j)}}$ = no. of movies rated by user $j$

To learn $\underline{\theta^{(j)}}$:

$$\min_{\Theta^{(j)}} \frac{1}{2\cancel{m^{(j)}}} \sum_{i: r(i,j)=1} \underline{\left((\Theta^{(j)})^T(x^{(i)}) - y^{(i,j)}\right)^2} + \frac{\lambda}{2\cancel{m^{(j)}}} \sum_{k=1}^{n} (\Theta_k^{(j)})^2$$

**Optimization objective:**

To learn $\theta^{(j)}$ (parameter for user $j$ ):

$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i:r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

To learn $\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(n_u)}$:

$$\min_{\theta^{(1)}, \ldots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

$\theta^{(1)}, \ldots, \theta^{(n_u)}$

# Optimization algorithm:

$$\min_{\theta^{(1)},\ldots,\theta^{(n_u)}} \frac{1}{2}\sum_{j=1}^{n_u}\sum_{i:r(i,j)=1}\left((\theta^{(j)})^T x^{(i)} - y^{(i,j)}\right)^2 + \frac{\lambda}{2}\sum_{j=1}^{n_u}\sum_{k=1}^{n}(\theta_k^{(j)})^2$$

$$J(\theta^{(1)},\ldots,\theta^{(n_u)})$$

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1}((\theta^{(j)})^T x^{(i)} - y^{(i,j)})x_k^{(i)} \quad (\text{for } k = 0)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha\left(\sum_{i:r(i,j)=1}((\theta^{(j)})^T x^{(i)} - y^{(i,j)})x_k^{(i)} + \lambda\theta_k^{(j)}\right) \quad (\text{for } k \neq 0)$$

$$\frac{1}{m^{(j)}}$$

$$\frac{\partial}{\partial\theta_k^{(j)}} J(\theta^{(1)},\ldots,\theta^{(n_u)})$$

Machine Learning

# Recommender Systems

## Collaborative filtering

# Problem motivation

| Movie | Alice (1) | Bob (2) | Carol (3) | Dave (4) | $x_1$ (romance) | $x_2$ (action) |
|---|---|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 | 0.9 | 0 |
| Romance forever | 5 | ? | ? | 0 | 1.0 | 0.01 |
| Cute puppies of love | ? | 4 | 0 | ? | 0.99 | 0 |
| Nonstop car chases | 0 | 0 | 5 | 4 | 0.1 | 1.0 |
| Swords vs. karate | 0 | 0 | 5 | ? | 0 | 0.9 |

# Problem motivation



| Movie | Alice (1) $\theta^{(1)}$ | Bob (2) $\theta^{(2)}$ | Carol (3) $\theta^{(3)}$ | Dave (4) $\theta^{(4)}$ | $x_1$ (romance) | $x_2$ (action) |
|---|---|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 | 1.0 | 0.0 |
| Romance forever | 5 | ? | ? | 0 | ? | ? |
| Cute puppies of love | ? | 4 | 0 | ? | ? | ? |
| Nonstop car chases | 0 | 0 | 5 | 4 | ? | ? |
| Swords vs. karate | 0 | 0 | 5 | ? | ? | ? |

$x_0 = 1$

$$x^{(1)} = \begin{bmatrix} 1 \\ 1.0 \\ 0.0 \end{bmatrix}$$

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \quad \theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \quad \theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} \quad \theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

$x^{(1)}$

$\theta^{(j)}$

$(\theta^{(1)})^T x^{(1)} \approx 5$

$(\theta^{(2)})^T x^{(1)} \approx 5$

$(\theta^{(3)})^T x^{(1)} \approx 0$

$(\theta^{(4)})^T x^{(1)} \approx 0$

Andrew Ng

# Optimization algorithm

Given $\theta^{(1)}, \ldots, \theta^{(n_u)}$, to learn $x^{(i)}$:

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^{n} (x_k^{(i)})^2$$

Given $\theta^{(1)}, \ldots, \theta^{(n_u)}$, to learn $x^{(1)}, \ldots, x^{(n_m)}$:

$$\min_{x^{(1)}, \ldots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2$$

# Collaborative filtering

Given $x^{(1)}, \ldots, x^{(n_m)}$ (and movie ratings),
can estimate $\theta^{(1)}, \ldots, \theta^{(n_u)}$

Given $\theta^{(1)}, \ldots, \theta^{(n_u)}$,
can estimate $x^{(1)}, \ldots, x^{(n_m)}$

$\sigma^{(i,j)}$
$y^{(i,j)}$

Guess $\quad \theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \ldots$

# Recommender Systems

## Collaborative filtering algorithm

Machine Learning

# Collaborative filtering optimization objective

$(i,j) : r(i,j) = 1$

$x \in \mathbb{R}^n$

$x \in \mathbb{R}^{n+1}$

$\theta \in \mathbb{R}^n$

$x_1 = 1$

Given $x^{(1)}, \ldots, x^{(n_m)}$, estimate $\theta^{(1)}, \ldots, \theta^{(n_u)}$:

$$\min_{\theta^{(1)}, \ldots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

Given $\theta^{(1)}, \ldots, \theta^{(n_u)}$, estimate $x^{(1)}, \ldots, x^{(n_m)}$:

$$\min_{x^{(1)}, \ldots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2$$

Minimizing $x^{(1)}, \ldots, x^{(n_m)}$ and $\theta^{(1)}, \ldots, \theta^{(n_u)}$ simultaneously:

$$J(x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

$$\min_{\substack{x^{(1)}, \ldots, x^{(n_m)} \\ \theta^{(1)}, \ldots, \theta^{(n_u)}}} J(x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)})$$

$\theta \to x \to \theta \to x \to \ldots$

Andrew Ng

**Collaborative filtering algorithm**

1. Initialize $x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)}$ to small random values.

2. Minimize $J(x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)})$ using gradient descent (or an advanced optimization algorithm). E.g. for every $j = 1, \ldots, n_u, i = 1, \ldots, n_m$ :

$$x_k^{(i)} := x_k^{(i)} - \alpha \left( \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})\theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$

3. For a user with parameters $\theta$ and a movie with (learned) features $x$, predict a star rating of $\theta^T x$ .

$x_0 = 1$    $x \in \mathbb{R}^n, \theta \in \mathbb{R}^n$

$\theta_1$
$\vdots$
$\theta_n$

$\dfrac{\partial}{\partial x_k^{(i)}} J(\ldots)$

$(\theta^{(j)})^T (x^{(i)})$

Andrew Ng

Machine Learning

# Recommender Systems

## Vectorization:
## Low rank matrix factorization

# Collaborative filtering

| Movie | Alice (1) | Bob (2) | Carol (3) | Dave (4) |
|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 |
| Romance forever | 5 | ? | ? | 0 |
| Cute puppies of love | ? | 4 | 0 | ? |
| Nonstop car chases | 0 | 0 | 5 | 4 |
| Swords vs. karate | 0 | 0 | 5 | ? |

$$n_m = 5$$
$$n_u = 4$$

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

$y^{(i,j)}$

# Collaborative filtering

$$X (\Theta)^T \leftarrow$$

$$(\Theta^{(j)})^T (x^{(i)})$$

$$(i,j) \nearrow$$

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

## Predicted ratings:

$$\begin{bmatrix} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \dots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \dots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \dots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{bmatrix}$$

$$X = \begin{bmatrix} -(x^{(1)})^T- \\ -(x^{(2)})^T- \\ \vdots \\ -(x^{(n_m)})^T- \end{bmatrix} \qquad \Theta = \begin{bmatrix} -(\Theta^{(1)})^T- \\ -(\Theta^{(2)})^T- \\ \vdots \\ -(\Theta^{(n_u)})^T- \end{bmatrix}$$

$$\rightarrow \text{Low rank matrix factorization}$$

Andrew Ng

## Finding related movies

For each product $i$, we learn a feature vector $x^{(i)} \in \mathbb{R}^n$.

$\rightarrow$ $x_1 = $ romance, $x_2 = $ action, $x_3 = $ comedy, $x_4 = \dots$

How to find movies $j$ related to movie $i$?

Small $\|x^{(i)} - x^{(j)}\|$ $\rightarrow$ Movies $j$ and $i$ are "similar"

5 most similar movies to movie $i$:
Find the 5 movies $j$ with the smallest $\|x^{(i)} - x^{(j)}\|$.

Machine Learning

# Recommender Systems

Implementational detail: Mean normalization

# Users who have not rated any movies

| Movie | Alice (1) | Bob (2) | Carol (3) | Dave (4) | Eve (5) |
|---|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 | ? |
| Romance forever | 5 | ? | ? | 0 | ? |
| Cute puppies of love | ? | 4 | 0 | ? | ? |
| Nonstop car chases | 0 | 0 | 5 | 4 | ? |
| Swords vs. karate | 0 | 0 | 5 | ? | ? |

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

$$\min_{\substack{x^{(1)},\ldots,x^{(n_m)} \\ \theta^{(1)},\ldots,\theta^{(n_u)}}} \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

$$n=2 \qquad \theta^{(5)} \in \mathbb{R}^2 \qquad \theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \qquad \frac{\lambda}{2}\left[ (\theta_1^{(5)})^2 + (\theta_2^{(5)})^2 \right]$$

$$(\theta^{(5)})^T x^{(i)} = 0$$

# Mean Normalization:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

2.5
2.5
2
⋮

$$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

learn $\theta^{(j)}, x^{(i)}$

For user $j$, on movie $i$ predict:

$$\rightarrow (\theta^{(j)})^T (x^{(i)}) + \mu_i$$

User 5 (Eve):

$$\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\underbrace{(\theta^{(5)})^T (x^{(i)})}_{=0} + \boxed{\mu_i}$$