

Collaborative filtering

1 / 14

Recommender systems

Netflix problem: 480000 users, 18000 movies.

- ▶ Each user rates some subset of the movies with a score in $\{1, 2, \dots, 5\}$.
On average, each user rates around 200 movies, though the variance is high (e.g., some user has rated >17000 movies).
- ▶ Goal is to predict how users would rate movies they haven't seen.
- ▶ Common to reduce $\{1, 2, \dots, 5\}$ to $\{-1, +1\}$ (e.g., $\{4, 5\} \mapsto +1$).

Common supervised learning approach:

- ▶ Get *features* for each user i and movie j (e.g., $\mathbf{x}_i \in \mathbb{R}^d$ and $\tilde{\mathbf{x}}_j \in \mathbb{R}^d$); goal is to predict rating as function of $(\mathbf{x}_i, \tilde{\mathbf{x}}_j)$.
- ▶ **Linear function:** $\mathbf{x}_i \tilde{\mathbf{x}}_j^\top \mapsto \mathbf{x}_i^\top \mathbf{W} \tilde{\mathbf{x}}_j$ for $\mathbf{W} \in \mathbb{R}^{d \times d}$.
Can use SVM, logistic regression, etc.

What if you don't have any features?

2 / 14

Collaborative filtering

Collaborative filtering (CF):

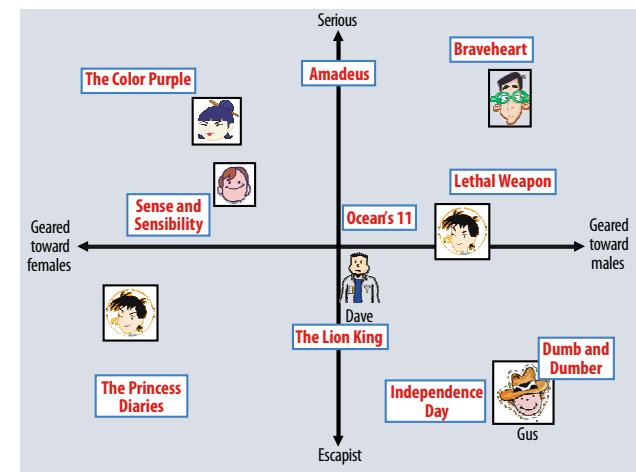
- ▶ Users who rate the same movie similarly are likely to be similar.
- ▶ Movies that are rated similarly by the same user are likely to be similar.

Can we formulate a model that expresses this intuition?

- ▶ **Side goal for model:** learn *feature representations* of users and movies that are semantically meaningful.

3 / 14

User/movie parameters



(Graphic is from Koren, Bell, and Volinsky.)

- ▶ Each user i represented by two-dimensional parameter $\mathbf{u}_i = (u_{i,1}, u_{i,2})$.
- ▶ Each movie j represented by two-dimensional parameter $\mathbf{v}_j = (v_{j,1}, v_{j,2})$.

4 / 14

Simple statistical model for CF

- **Model parameters:** $\theta := (\mathbf{u}_1, \dots, \mathbf{u}_m, \mathbf{v}_1, \dots, \mathbf{v}_n)$.
 - Parameter vector for user i : $\mathbf{u}_i \in \mathbb{R}^k$.
 - Parameter vector for movie j : $\mathbf{v}_j \in \mathbb{R}^k$.

(Assume $k \leq \min\{m, n\}$.)

- **Distribution of ratings:** the $A_{i,j}$ are independent, and

$$A_{i,j} \sim \mathcal{N}(\langle \mathbf{u}_i, \mathbf{v}_j \rangle, 1).$$

$$\mathbb{E} \left\{ \begin{bmatrix} A_{1,1} & A_{1,2} & \cdots & A_{1,n} \\ A_{2,1} & A_{2,2} & \cdots & A_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1} & A_{m,2} & \cdots & A_{m,n} \end{bmatrix} \right\} = \underbrace{\begin{bmatrix} \leftarrow & \mathbf{u}_1^\top & \rightarrow \\ \leftarrow & \mathbf{u}_2^\top & \rightarrow \\ \vdots & \vdots & \vdots \\ \leftarrow & \mathbf{u}_m^\top & \rightarrow \end{bmatrix}}_{m \times k \text{ matrix}} \underbrace{\begin{bmatrix} \uparrow & \mathbf{v}_1 & \downarrow & \vdots & \uparrow \\ \uparrow & \mathbf{v}_2 & \downarrow & \vdots & \uparrow \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \uparrow & \mathbf{v}_n & \downarrow & \vdots & \uparrow \end{bmatrix}}_{k \times n \text{ matrix}}.$$

5 / 14

Prediction and parameter estimation

Data: $\mathcal{A} := \{a_{i,j} \in \mathbb{R} : (i,j) \in \Omega\}$, for subset of user/movie pairs
 $\Omega \subseteq \{1, 2, \dots, m\} \times \{1, 2, \dots, n\}$.

Prediction

Given parameters $\hat{\theta} = (\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_m, \hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_n)$, can *predict* $a_{i,j}$ for $(i,j) \notin \Omega$ (i.e., user/movie pairs you don't have ratings for):

$$\hat{a}_{i,j} := \langle \hat{\mathbf{u}}_i, \hat{\mathbf{v}}_j \rangle.$$

Maximum likelihood estimation

Log-likelihood of $\theta = (\mathbf{u}_1, \dots, \mathbf{u}_m, \mathbf{v}_1, \dots, \mathbf{v}_n)$ given $\mathcal{A} := \{a_{i,j} : (i,j) \in \Omega\}$:

$$\mathcal{L}(\theta; \mathcal{A}) = -\frac{1}{2} \sum_{(i,j) \in \Omega} (a_{i,j} - \langle \mathbf{u}_i, \mathbf{v}_j \rangle)^2 + (\text{terms not involving } \theta).$$

Unfortunately, this is generally hard to maximize.


6 / 14

Special case: complete set of ratings


Special case: have ratings for all user/movie pairs
(i.e., $\Omega = \{1, \dots, m\} \times \{1, \dots, n\}$).


$$\mathcal{L}(\theta; \mathcal{A}) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n (a_{i,j} - \langle \mathbf{u}_i, \mathbf{v}_j \rangle)^2 + (\text{terms not involving } \theta).$$


MLE is given by rank- k singular value decomposition (SVD):


$$\mathbf{A}$$
$$(m \times n)$$

\approx


$$\mathbf{L}$$
$$(m \times k)$$


$$\mathbf{S}$$
$$(k \times k)$$


$$\mathbf{R}^\top$$
$$(k \times n)$$

Use rows of $\mathbf{L}\mathbf{S}^{1/2}$ as the \mathbf{u}_i , and rows of $\mathbf{R}\mathbf{S}^{1/2}$ as the \mathbf{v}_j .

(Solution is not always unique!)

7 / 14

Matrix completion

General case: MLE is equivalent to (low rank) matrix completion problem:

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \quad & \sum_{(i,j) \in \Omega} (a_{i,j} - X_{i,j})^2 \\ \text{s.t.} \quad & \text{rank}(\mathbf{X}) \leq k. \end{aligned}$$

Objective is not convex due to rank constraint.

8 / 14

Known user parameters

$$\mathbb{E} \left\{ \begin{bmatrix} A_{1,1} & A_{1,2} & \cdots & A_{1,n} \\ A_{2,1} & A_{2,2} & \cdots & A_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1} & A_{m,2} & \cdots & A_{m,n} \end{bmatrix} \right\} = \begin{bmatrix} \leftarrow & \mathbf{u}_1^\top & \rightarrow \\ \leftarrow & \mathbf{u}_2^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{u}_m^\top & \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_n \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$$

Hypothetical situation: suppose all user parameters are already known.

MLE for movie parameters \mathbf{v}_j given by ordinary least squares:

$$\mathbf{v}_j := \arg \min_{\mathbf{v} \in \mathbb{R}^k} \sum_{(i,j) \in \Omega} (\langle \mathbf{u}_i, \mathbf{v} \rangle - a_{i,j})^2.$$

Analogous if, instead, we suppose all movie parameters are already known:
get MLE for user parameters \mathbf{u}_j via ordinary least squares.

Idea: alternate between the two ...

Alternating (regularized) least squares

- Initialize $\hat{\mathbf{u}}_i \in \mathbb{R}^k$ for each user i and $\hat{\mathbf{v}}_j \in \mathbb{R}^k$ for each movie j .
- For $t = 1, 2, \dots$:

- For each user $i = 1, 2, \dots, m$,

$$\hat{\mathbf{u}}_i := \arg \min_{\mathbf{u} \in \mathbb{R}^k} \sum_{(i,j) \in \Omega} (\langle \mathbf{u}, \hat{\mathbf{v}}_j \rangle - a_{i,j})^2 + \lambda \|\mathbf{u}\|_2^2.$$

- For each movie $j = 1, 2, \dots, n$,

$$\hat{\mathbf{v}}_j := \arg \min_{\mathbf{v} \in \mathbb{R}^k} \sum_{(i,j) \in \Omega} (\langle \hat{\mathbf{u}}_i, \mathbf{v} \rangle - a_{i,j})^2 + \lambda \|\mathbf{v}\|_2^2.$$

(Could also switch or randomize order of updates.)

Alternating stochastic gradient method

When $|\Omega|$ is very large, each iteration can be expensive.

Alternating stochastic gradient method (for $\lambda = 0$)

- Initialize $\hat{\mathbf{u}}_i \in \mathbb{R}^k$ for each user i and $\hat{\mathbf{v}}_j \in \mathbb{R}^k$ for each movie j .
- For $t = 1, 2, \dots$:
 - Pick $(i_t, j_t) \in \Omega$ uniformly at random.
 - Update:

$$\begin{aligned} \hat{\mathbf{u}}_{i_t} &:= \hat{\mathbf{u}}_{i_t} - 2\eta (\langle \hat{\mathbf{u}}_{i_t}, \hat{\mathbf{v}}_{j_t} \rangle - a_{i_t, j_t}) \hat{\mathbf{v}}_{j_t}, \\ \hat{\mathbf{v}}_{j_t} &:= \hat{\mathbf{v}}_{j_t} - 2\eta (\langle \hat{\mathbf{u}}_{i_t}, \hat{\mathbf{v}}_{j_t} \rangle - a_{i_t, j_t}) \hat{\mathbf{u}}_{i_t}. \end{aligned}$$

(Could also switch or randomize order of updates.)

Example

Ran alternating (regularized) least squares on 800000 ratings of movies from $m = 6040$ users and $n = 3952$ movies.

Movie parameters $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n \in \mathbb{R}^k$ give *feature representations* of movies.
Are they semantically meaningful?

Some nearest-neighbor pairs $(\mathbf{v}_i, \text{NN}(\mathbf{v}_i))$:

- Toy Story (1995), Toy Story 2 (1999)
- Sense and Sensibility (1995), Emma (1996)
- Heat (1995), Carlito's Way (1993)
- The Crow (1994), Blade (1998)
- Forrest Gump (1994), Dances with Wolves (1990)
- Mrs. Doubtfire (1993), The Bodyguard (1992) ???
- ...

- **Initialization:** can't initialize with all $\mathbf{u}_i = \mathbf{v}_j = \mathbf{0}$.
(Try random instead.)
- How to pick k or λ ? Cross validation.
- **Model with global/per-user/per-movie biases:**
Parameters: $\theta = (\mathbf{u}_1, \dots, \mathbf{u}_m, \mathbf{v}_1, \dots, \mathbf{v}_n, b_1, \dots, b_m, c_1, \dots, c_n, \mu)$

$$A_{i,j} \sim \text{N}(\mu + b_i + c_j + \langle \mathbf{u}_i, \mathbf{v}_j \rangle, 1).$$

- **Combination with user/movie features $\mathbf{x}_i, \tilde{\mathbf{x}}_j \in \mathbb{R}^d$:**
Parameters: $\theta = (\mathbf{u}_1, \dots, \mathbf{u}_m, \mathbf{v}_1, \dots, \mathbf{v}_n, b_1, \dots, b_m, c_1, \dots, c_n, \mu, \mathbf{W})$
- $$A_{i,j} \mid \mathbf{X} = \mathbf{x}_i \wedge \tilde{\mathbf{X}} = \tilde{\mathbf{x}}_j \sim \text{N}(\mu + b_i + c_j + \langle \mathbf{u}_i, \mathbf{v}_j \rangle + \mathbf{x}_i^\top \mathbf{W} \tilde{\mathbf{x}}_j, 1).$$
- Many other variations!

1. Simple statistical model for CF; two methods to (attempt to) compute the MLE.
2. Some possible generalizations of the CF model.