

Recommendation Systems

CS771: Introduction to Machine Learning

Purushottam Kar



Outline of discussion

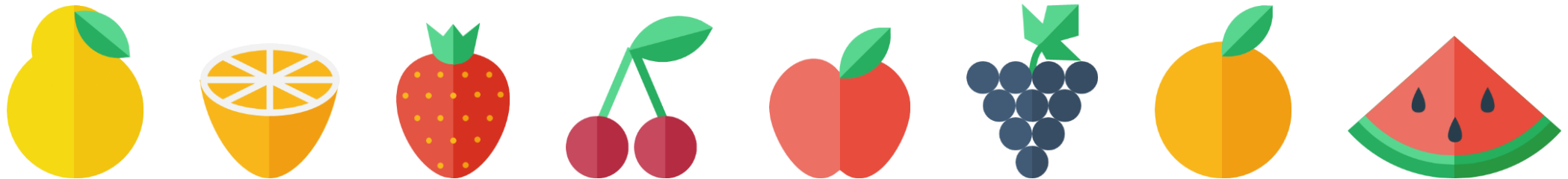
- Three major techniques
 - Collaborative Filtering
 - Extreme Classification
 - Multi-armed Bandits
- Barely scratching the surface
- Entire conference RecSys dedicated to recommendation systems
- Look for papers in other conferences too!
NIPS, ICML, KDD, WWW, WSDM

Problem Statement

- There are users ...



- ... and there are items ...



- ... and we need to recommend, to each user, items s/he will “like”
- Users usually come to us one at a time in no particular order

A Powerful Abstraction

User

- Amazon.com user
- Facebook user
- Student
- Patient
- Search query
- ...

Items

- Amazon.com products
- Facebook posts
- Study material
- Medicines
- Internet webpages
- ...

- Users are in millions and thousands come each second
- Items are in millions and each user likes only 5-10 items
- Some items are popular but most items liked by only 5-10 users

Collaborative Filtering

The Matrix Completion Problem

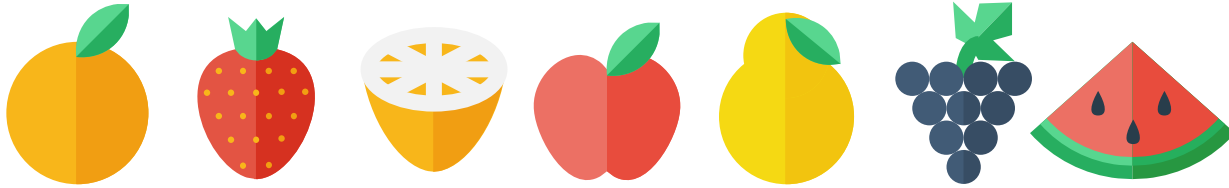
- m users n items
- Rating matrix $X \in \mathbb{R}^{m \times n}$
- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?

The Matrix Completion Problem



- m users n items
- Rating matrix $X \in \mathbb{R}^{m \times n}$
- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?

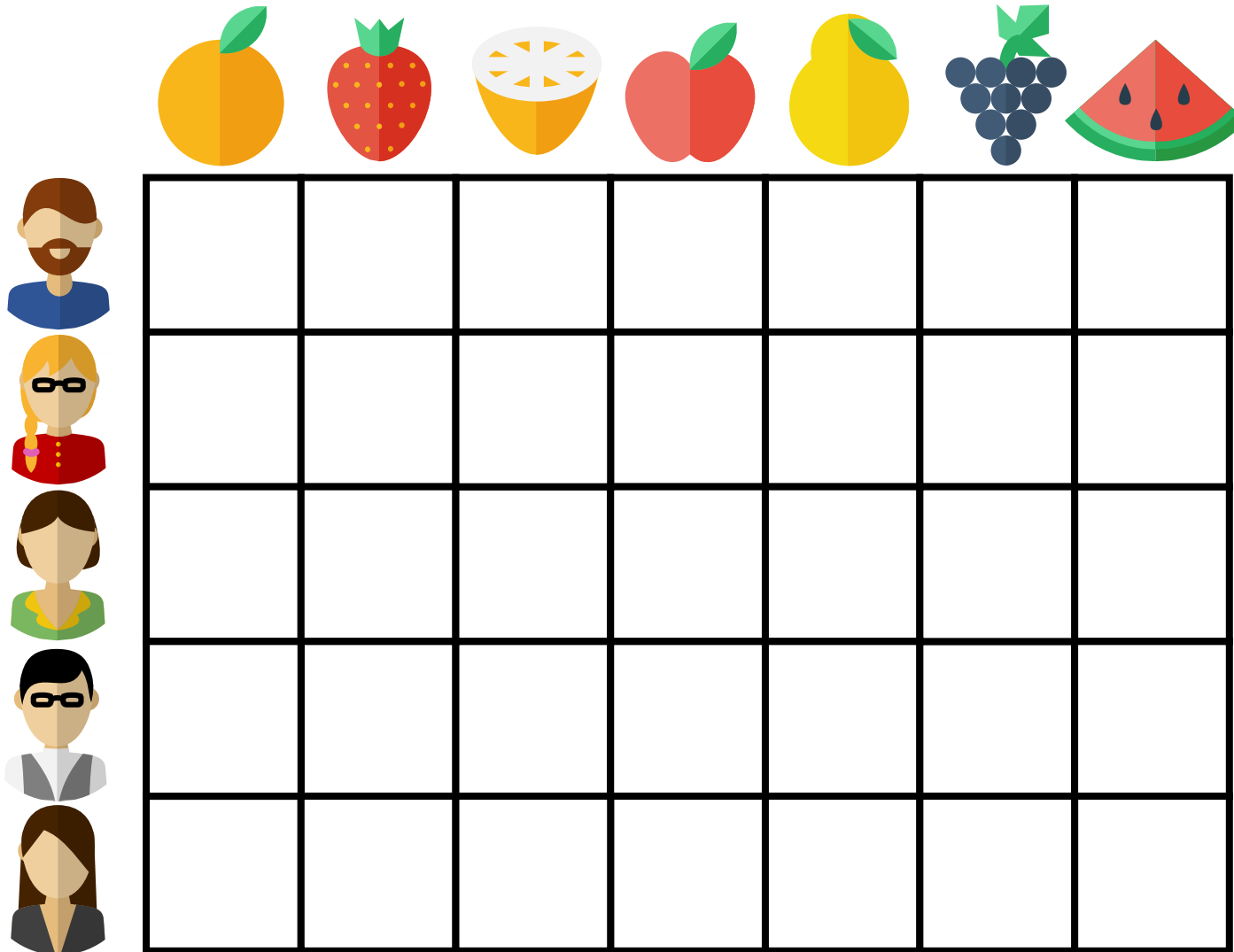
The Matrix Completion Problem


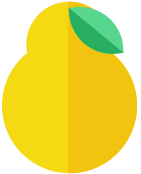
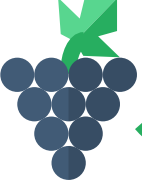




- m users n items
- Rating matrix $X \in \mathbb{R}^{m \times n}$
- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?



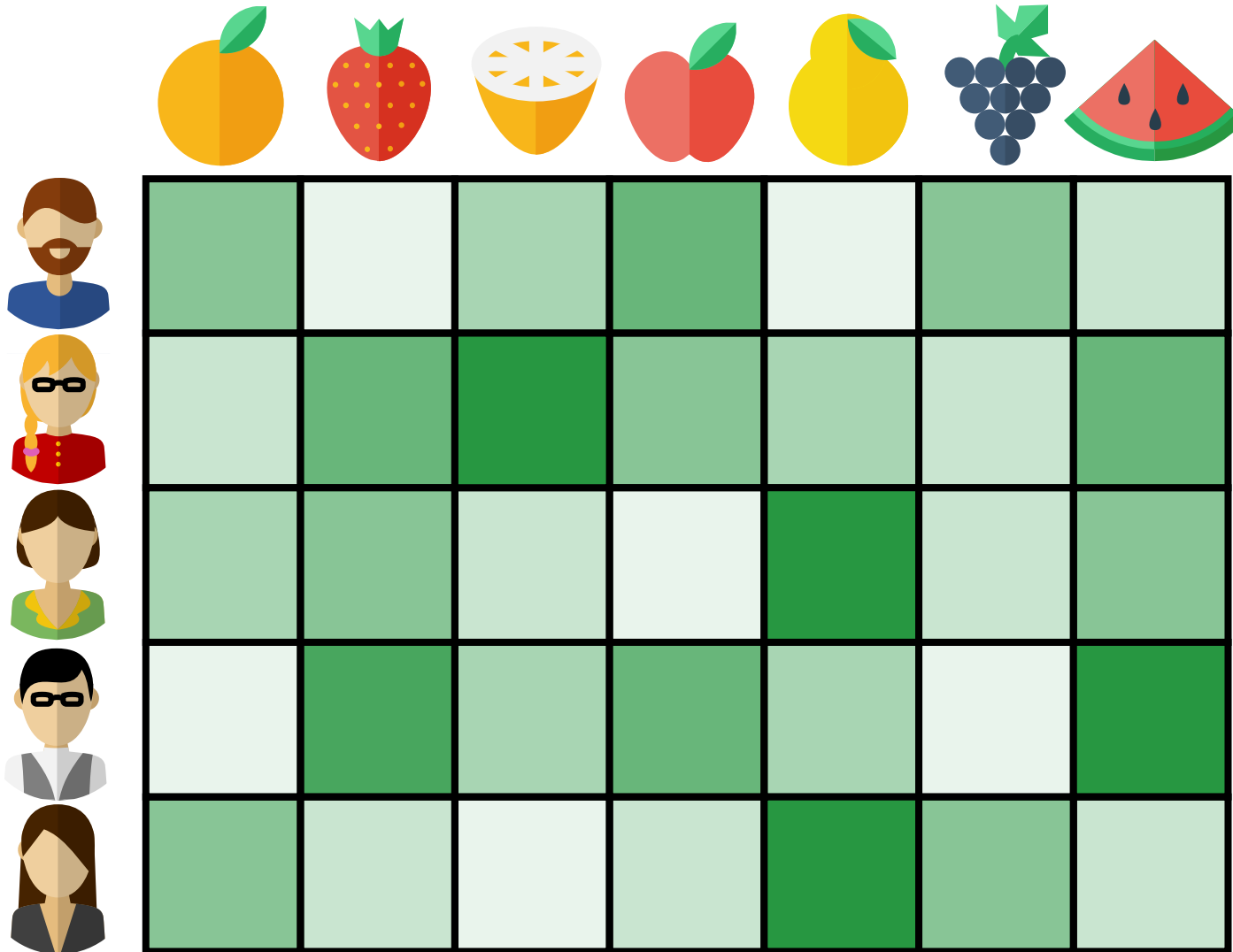
The Matrix Completion Problem



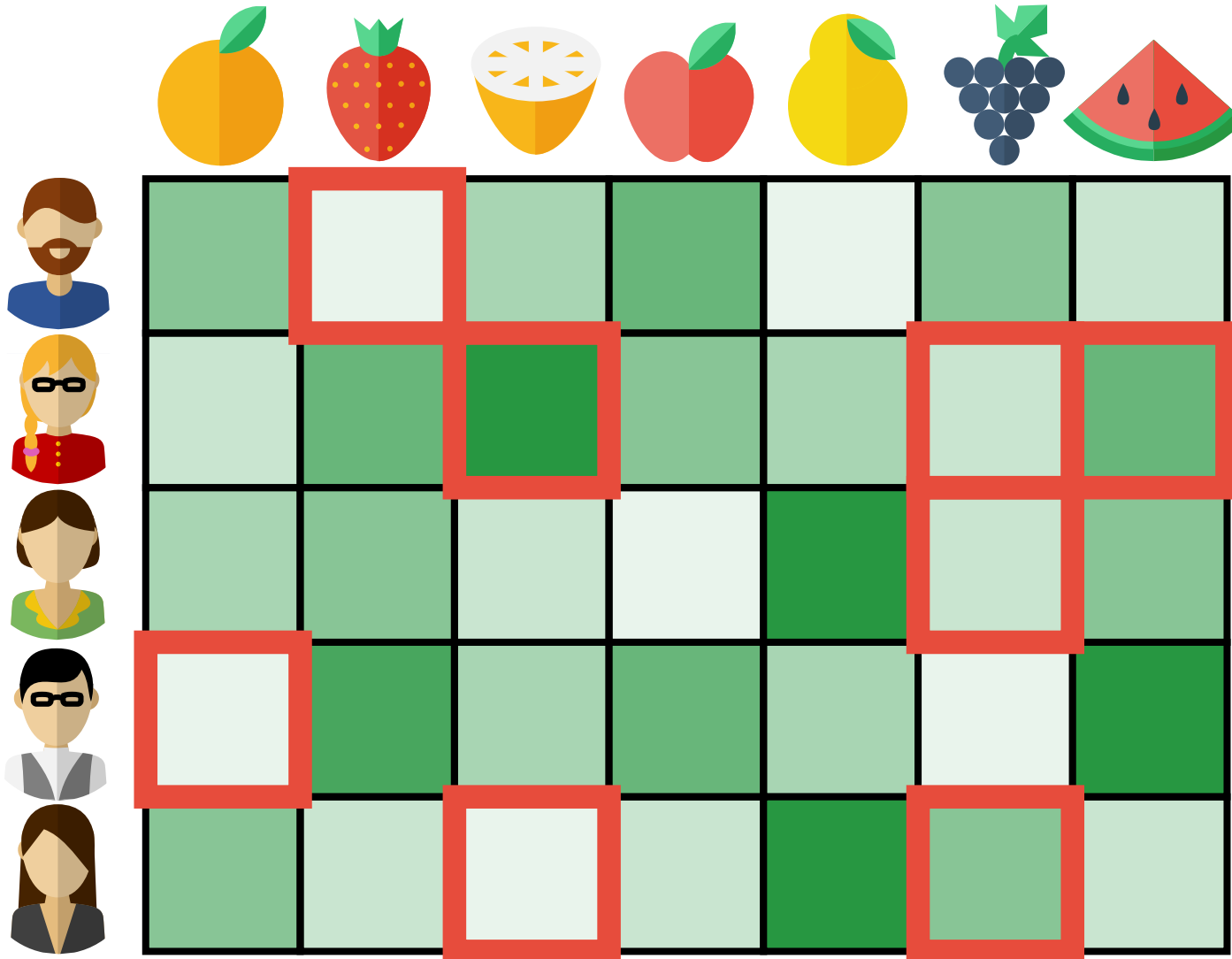
- m users n items
- Rating matrix $X \in \mathbb{R}^{m \times n}$
- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?

The Matrix Completion Problem



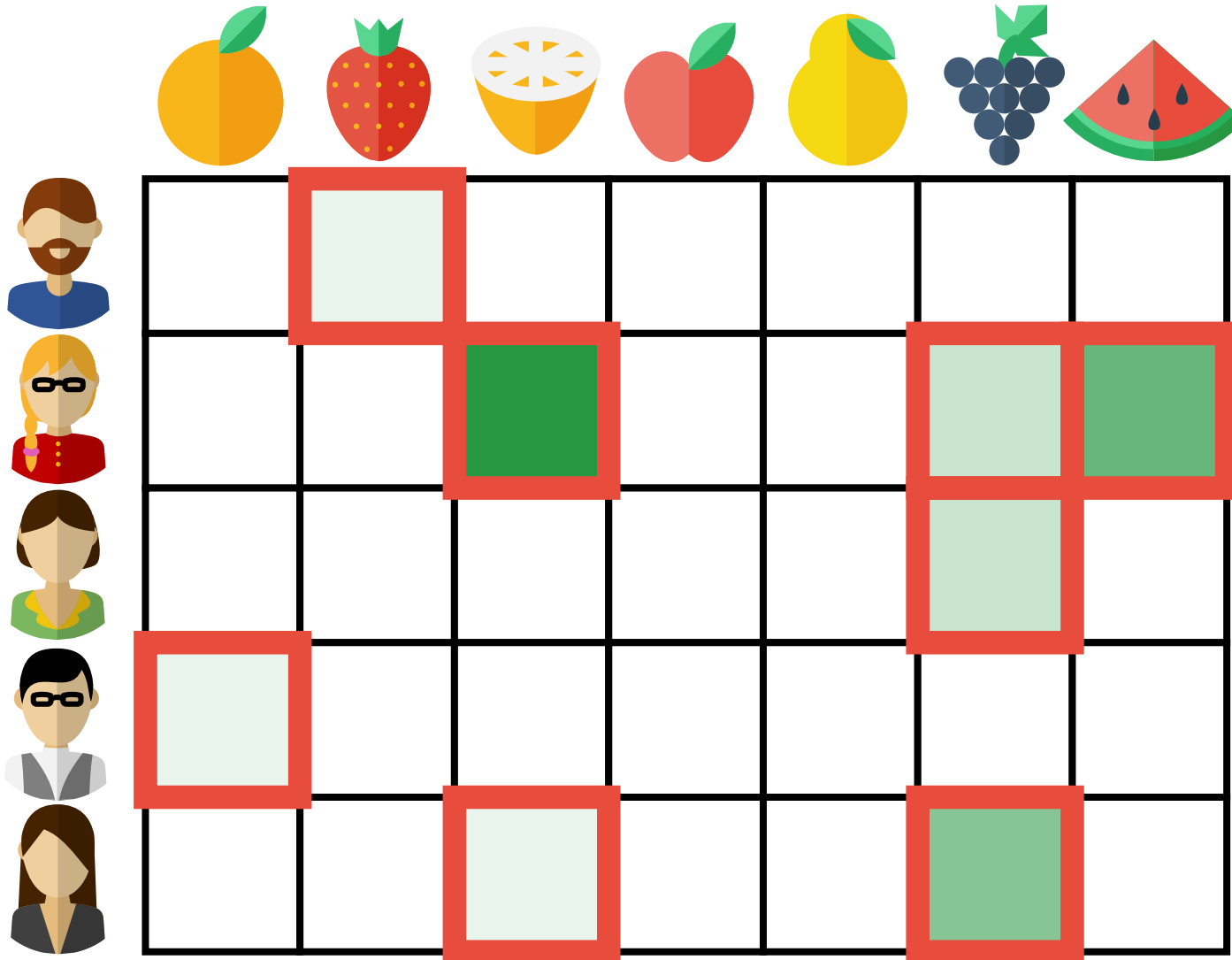
- m users n items
- Rating matrix $X \in \mathbb{R}^{m \times n}$
- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?

The Matrix Completion Problem



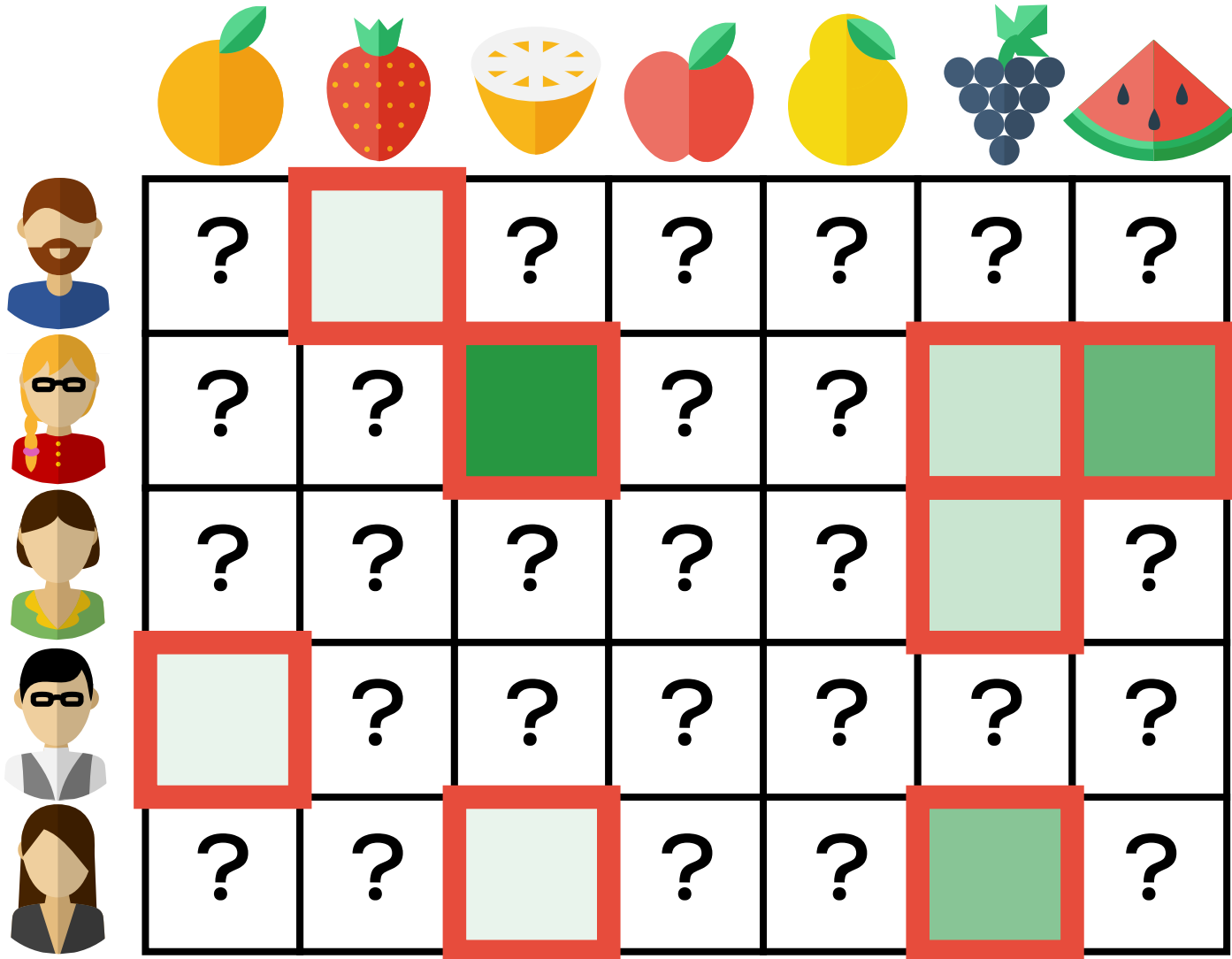
- m users n items
- Rating matrix $X \in \mathbb{R}^{m \times n}$
- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?


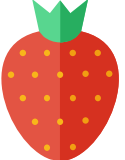


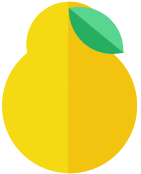
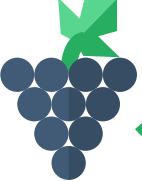






The Matrix Completion Problem



- m users n items
- Rating matrix $X \in \mathbb{R}^{m \times n}$
- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?

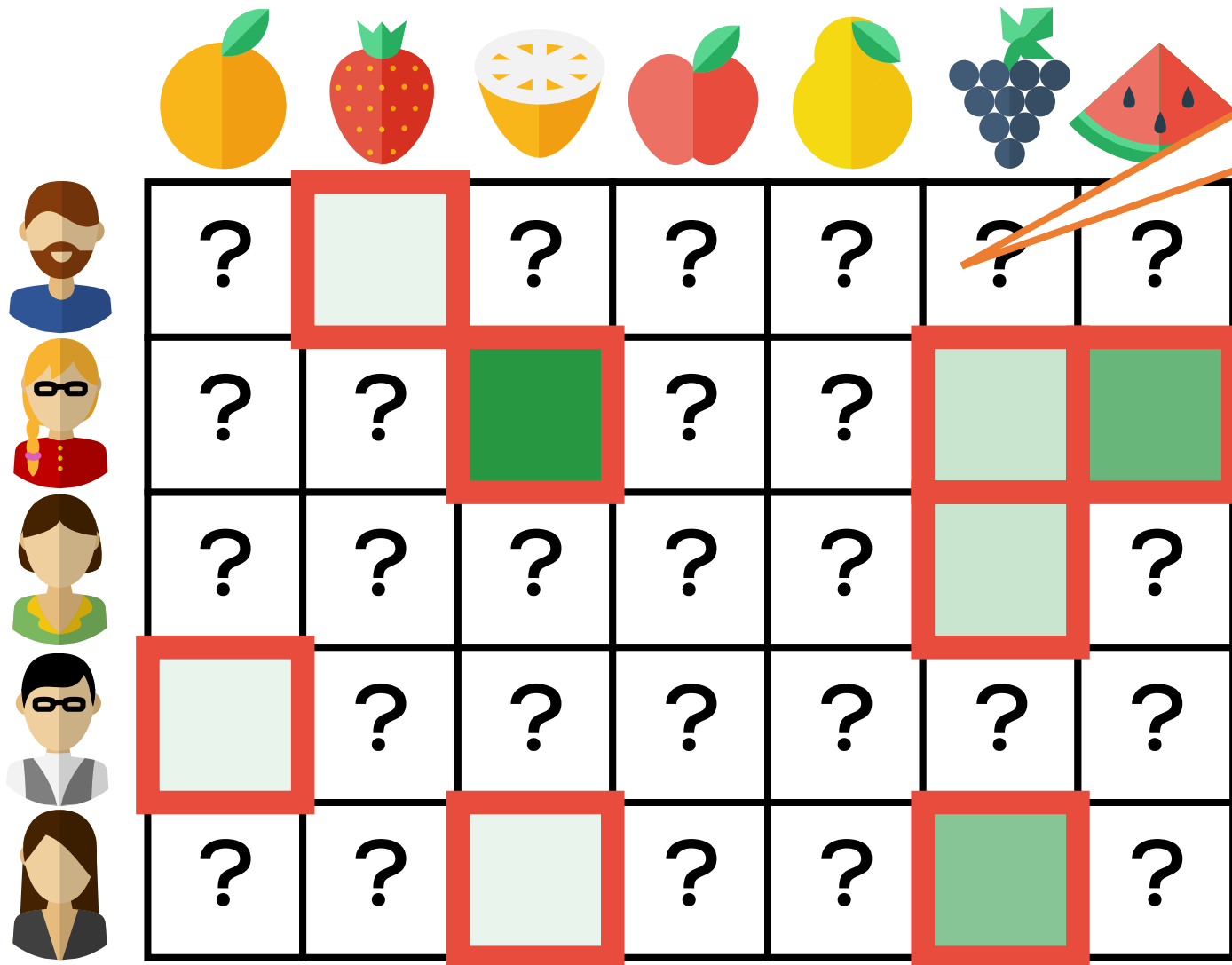
The Matrix Completion Problem


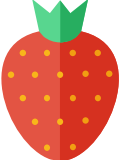


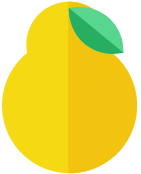
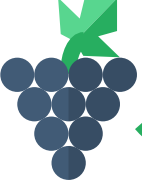








							
	?		?	?	?	?	?
	?	?		?	?		
	?	?	?	?	?		?
		?	?	?	?	?	?
	?	?		?	?		?

- m users n items
- Rating matrix $X \in \mathbb{R}^{m \times n}$
- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?

The Matrix Completion Problem



							
	?		?	?	?	?	?
	?	?		?	?		
	?	?	?	?	?		?
		?	?	?	?	?	?
	?	?		?	?		?

Knowing the entire matrix would reveal the most liked items for every user and allow me to recommend those!

- X_{ij} : rating indicating how much user i likes item j
- Deeper shade \Rightarrow more like
- Sometimes users rate some items (rarely though)
- So we get to see those entries of X , but only those
- Can we recover the unseen entries i.e. complete X ?

The Low-rank Matrix Completion Problem

- Unfortunately the matrix completion problem is ill-posed
- No reason why observed entries of the matrix should tell us anything about the unobserved entries
- Absolutely necessary to impose some structure on the matrix
- Popular assumption: the rating matrix X is low rank (say rank k)
- So need to find a low-rank matrix \hat{X} such that X and \hat{X} agree on the observed entries. Let $\Omega \subset [m] \times [n]$ be the locations observed

$$\arg \min_{\text{rank}(\hat{X}) \leq k} \sum_{(i,j) \in \Omega} (X_{ij} - \hat{X}_{ij})^2$$

- However, every rank k matrix \hat{X} can be written as $\hat{X} = UV^T$ where $U \in \mathbb{R}^{m \times k}$ and $V \in \mathbb{R}^{n \times k}$ (show this using singular value decomp!)

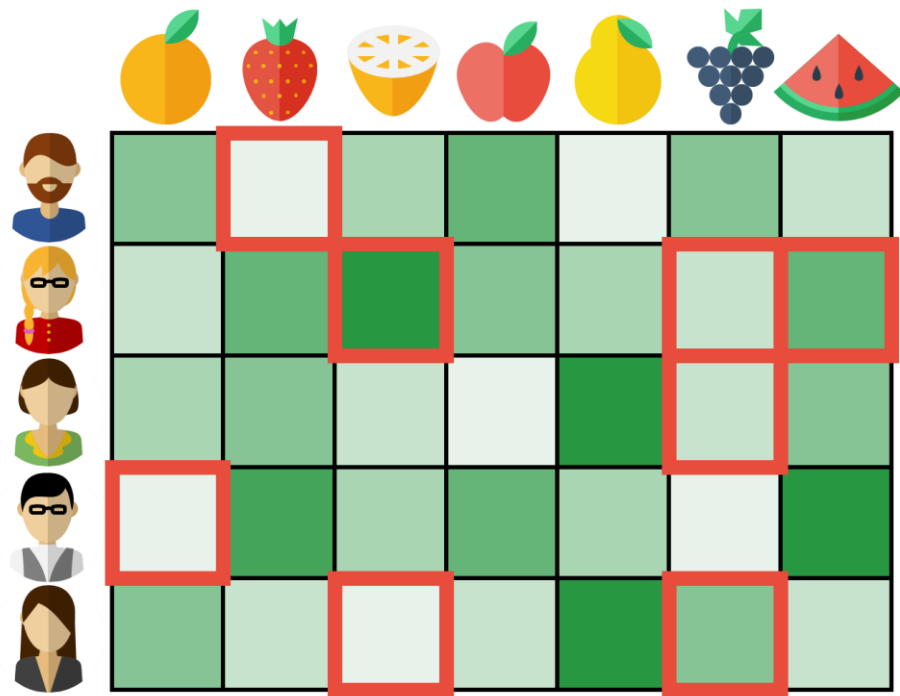
The Low-rank Matrix Completion Problem

- Unfamiliar intuition: Can encode intuitions like “similar items rated similarly by the same user/similar users”
- Noisy entries: k tuned based on how many entries observed – more entries, larger k
- Absolutely necessary to impose some structure on the matrix
- Popular assumption: the rating matrix X is low rank (say rank k)
- So need to find a low-rank matrix \hat{X} such that X and \hat{X} agree on the observed entries. Let $\Omega \subset [m] \times [n]$ be the locations observed

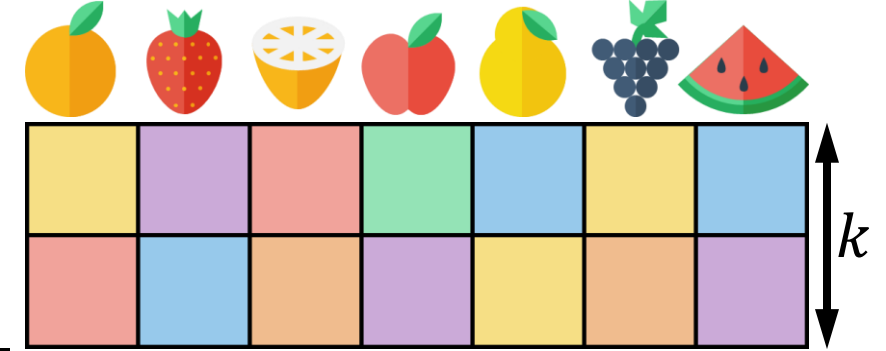
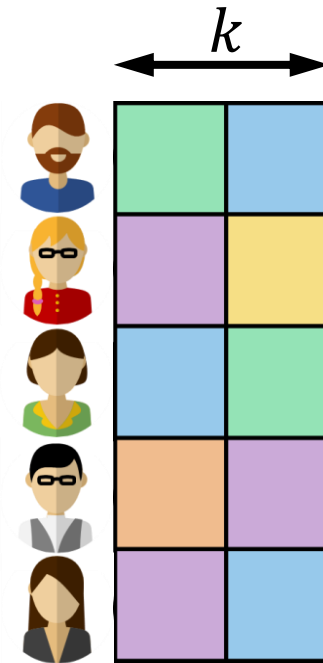
$$\arg \min_{\text{rank}(\hat{X}) \leq k} \sum_{(i,j) \in \Omega} (X_{ij} - \hat{X}_{ij})^2$$

- However, every rank k matrix \hat{X} can be written as $\hat{X} = UV^T$ where $U \in \mathbb{R}^{m \times k}$ and $V \in \mathbb{R}^{n \times k}$ (show this using singular value decomp!)

Low-rank Matrix Completion



\approx

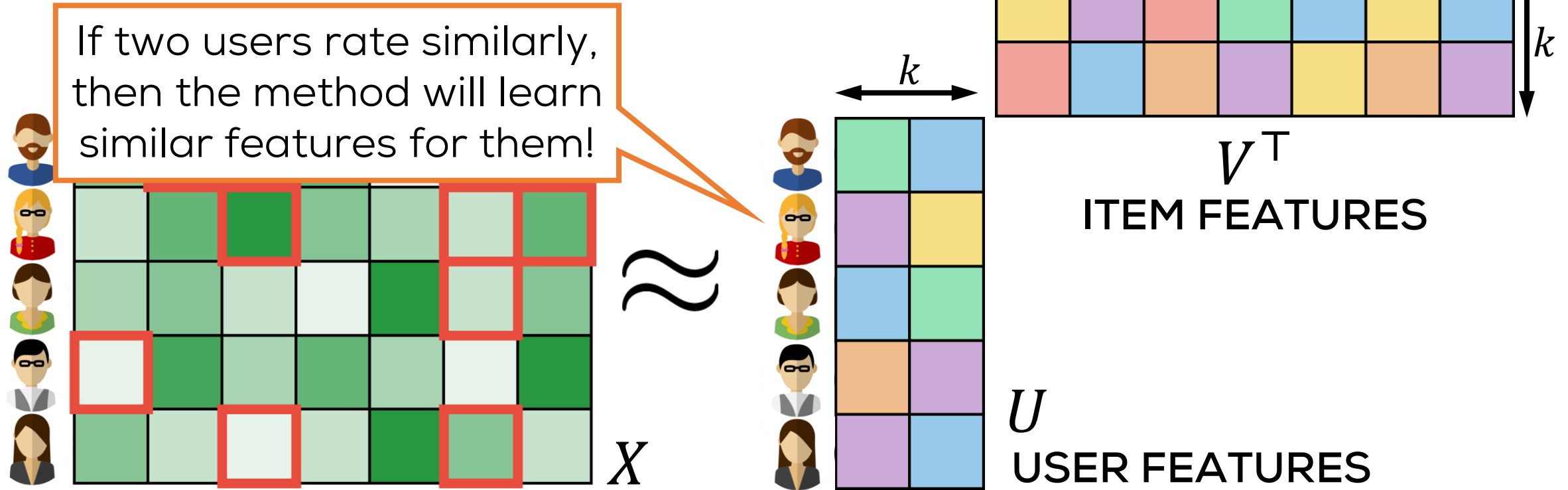


ITEM FEATURES

USER FEATURES

- So we have $X_{ij} \approx \langle \mathbf{u}^i, \mathbf{v}^j \rangle$ where $U = [\mathbf{u}^1, \dots, \mathbf{u}^m]^\top$ and $V = [\mathbf{v}^1, \dots, \mathbf{v}^n]^\top$
- Nice! As a side effect, vector representations of all users and items
- Such features can be very useful in recommendation (will soon see)
- Can also violate privacy if details users didn't tell are revealed
- E.g. if one of the k features turns out to be correlated with age

Low-rank Matrix Completion



- So we have $X_{ij} \approx \langle \mathbf{u}^i, \mathbf{v}^j \rangle$ where $U = [\mathbf{u}^1, \dots, \mathbf{u}^m]^T$ and $V = [\mathbf{v}^1, \dots, \mathbf{v}^n]^T$
- Nice! As a side effect, vector representations of all users and items
- Such features can be very useful in recommendation (will soon see)
- Can also violate privacy if details users didn't tell are revealed
- E.g. if one of the k features turns out to be correlated with age

Solving the LRMC problem

Solving the LRMC problem

$$\arg \min_{\text{rank}(\hat{X}) \leq k} \sum_{(i,j) \in \Omega} (X_{ij} - \hat{X}_{ij})^2$$

Solving the LRMC problem

$$\arg \min_{\substack{U \in \mathbb{R}^{m \times k} \\ V \in \mathbb{R}^{n \times k}}} \sum_{(i,j) \in \Omega} (X_{ij} - \langle \mathbf{u}^i, \mathbf{v}^j \rangle)^2$$

Solving the LRMC problem

$$\arg \min_{\substack{U \in \mathbb{R}^{m \times k} \\ V \in \mathbb{R}^{n \times k}}} \sum_{(i,j) \in \Omega} (X_{ij} - \langle \mathbf{u}^i, \mathbf{v}^j \rangle)^2 + \lambda \cdot (\|U\|_F^2 + \|V\|_F^2)$$

Solving the LRMC problem

Nice to
regularize

$$\arg \min_{\substack{U \in \mathbb{R}^{m \times k} \\ V \in \mathbb{R}^{n \times k}}} \sum_{(i,j) \in \Omega} (X_{ij} - \langle \mathbf{u}^i, \mathbf{v}^j \rangle)^2 + \lambda \cdot (\|U\|_F^2 + \|V\|_F^2)$$

Solving the LRMC problem

$$\arg \min_{\substack{U \in \mathbb{R}^{m \times k} \\ V \in \mathbb{R}^{n \times k}}} \sum_{(i,j) \in \Omega} (X_{ij} - \langle \mathbf{u}^i, \mathbf{v}^j \rangle)^2 + \lambda \cdot \left(\sum_{i=1}^m \|\mathbf{u}^i\|_2^2 + \sum_{j=1}^n \|\mathbf{v}^j\|_2^2 \right)$$

Nice to regularize

Solving the LRMC problem

Nice to
regularize

$$\arg \min_{\substack{U \in \mathbb{R}^{m \times k} \\ V \in \mathbb{R}^{n \times k}}} \sum_{(i,j) \in \Omega} (X_{ij} - \langle \mathbf{u}^i, \mathbf{v}^j \rangle)^2 + \lambda \cdot \left(\sum_{i=1}^m \|\mathbf{u}^i\|_2^2 + \sum_{j=1}^n \|\mathbf{v}^j\|_2^2 \right)$$

- Note that loss calculated only over observed entries i.e. $(i,j) \in \Omega$
- How to solve this ? Non convex problem! ☹
- Can use projected gradient descent on $\arg \min_{\text{rank}(\hat{X}) \leq k} \sum_{(i,j) \in \Omega} (X_{ij} - \hat{X}_{ij})^2$
- A faster way exists using alternating minimization
- AltMin is very versatile, gives us EM, Lloyd's algo, and now this
- Current state-of-the-art for matrix completion is AltMin

AltMin for LRMC

$$\arg \min_{\substack{U \in \mathbb{R}^{m \times k} \\ V \in \mathbb{R}^{n \times k}}} \sum_{(i,j) \in \Omega} (X_{ij} - \langle \mathbf{u}^i, \mathbf{v}^j \rangle)^2 + \lambda \cdot \left(\sum_{i=1}^m \|\mathbf{u}^i\|_2^2 + \sum_{j=1}^n \|\mathbf{v}^j\|_2^2 \right)$$

- For now, fix all \mathbf{v}^j and all \mathbf{u}^i except \mathbf{u}^1 . The problem now becomes

$$\mathbf{u}^1 = \arg \min_{\mathbf{u} \in \mathbb{R}^k} \sum_{(1,j) \in \Omega} (X_{1j} - \langle \mathbf{u}, \mathbf{v}^j \rangle)^2 + \lambda \cdot \|\mathbf{u}\|_2^2$$

- Wait ... this is just a ridge regression problem!
- What happens if I fix all \mathbf{u}^i and all \mathbf{v}^j except \mathbf{v}^1

$$\mathbf{v}^1 = \arg \min_{\mathbf{v} \in \mathbb{R}^k} \sum_{(i,1) \in \Omega} (X_{i1} - \langle \mathbf{u}^i, \mathbf{v} \rangle)^2 + \lambda \cdot \|\mathbf{v}\|_2^2$$

- Nice ... again just a ridge regression problem!

AltMin for LRMC

Notice that summation is only over items that user 1 has rated – usually very small

$V \in \mathbb{R}^{n \times n}$

$\langle \mathbf{u}^i, \mathbf{v}^j \rangle)^2 + \lambda \cdot \left(\mathbf{v}^j \text{ act as the "feature" vectors and } X_{1j} \text{ as "responses" in this ridge regression problem} \right)$

- For now, fix all \mathbf{v}^j and all \mathbf{u}^i except \mathbf{u}^1 . The problem now becomes

$$\mathbf{u}^1 = \arg \min_{\mathbf{u} \in \mathbb{R}^k} \sum_{(1,j) \in \Omega} (X_{1j} - \langle \mathbf{u}, \mathbf{v}^j \rangle)^2 + \lambda \cdot \|\mathbf{u}\|_2^2$$

Notice that summation is only over users that rated item 1 – usually small for most items

- Wait ... this is just a ridge regression problem
- What happens if I fix all \mathbf{u}^i and all \mathbf{v}^j except \mathbf{v}^1

$$\mathbf{v}^1 = \arg \min_{\mathbf{v} \in \mathbb{R}^k} \sum_{(i,1) \in \Omega} (X_{i1} - \langle \mathbf{u}^i, \mathbf{v} \rangle)^2 + \lambda \cdot \|\mathbf{v}\|_2^2$$

This means that these problems can be solved cheaply

- Nice ... again just a ridge regression problem!

Alternating Optimization Revisited

ALTERNATING OPTIMIZATION

1. Observed locations Ω and entries $\{X_{ij} : (i, j) \in \Omega\}$
2. Initialize $\mathbf{v}^{1,0}, \mathbf{v}^{2,0}, \dots, \mathbf{v}^{n,0}$
3. For $t = 1, 2, \dots$

1. Update the user vectors for $i = 1, \dots, m$

$$\mathbf{u}^{i,t} \leftarrow \left(\sum_{(i,j) \in \Omega} \mathbf{v}^{j,t-1} (\mathbf{v}^{j,t-1})^\top + \lambda \cdot I \right)^{-1} \left(\sum_{(i,j) \in \Omega} X_{ij} \cdot \mathbf{v}^{j,t-1} \right)$$

2. Update the item vectors for $j = 1, \dots, n$

$$\mathbf{v}^{j,t} \leftarrow \left(\sum_{(i,j) \in \Omega} \mathbf{u}^{i,t} (\mathbf{u}^{i,t})^\top + \lambda \cdot I \right)^{-1} \left(\sum_{(i,j) \in \Omega} X_{ij} \cdot \mathbf{u}^{i,t} \right)$$

4. Repeat until convergence

Alternating Optimization Revisited

ALTERNATING OPTIMIZATION

1. Observed locations Ω and entries $\{X_{ij} : (i, j) \in \Omega\}$
2. Initialize $\mathbf{v}^{1,0}, \mathbf{v}^{2,0}, \dots, \mathbf{v}^{n,0}$
3. For $t = 1, 2, \dots$

Inexpensive for most items and users

1. Update the user vectors for $i = 1, \dots, m$

Don't solve these RR problems optimally, just take an SGD step!

$$\mathbf{v}^{j,t} \leftarrow \left(\mathbf{v}^{j,t-1} (\mathbf{v}^{j,t-1})^\top + \lambda \cdot I \right)^{-1} \left(\sum_{(i,j) \in \Omega} X_{ij} \cdot \mathbf{v}^{j,t-1} \right)$$

$\mathcal{O}(k^3 + k^2(m+n))$ time per iteration

Can speed up even more

2. Update the item vectors for $j = 1, \dots, n$

$$\mathbf{v}^{j,t} \leftarrow \left(\sum_{(i,j) \in \Omega} \mathbf{u}^{j,t} (\mathbf{u}^{j,t})^\top + \lambda \cdot I \right)^{-1} \left(\sum_{(i,j) \in \Omega} X_{ij} \cdot \mathbf{u}^{j,t} \right)$$

4. Repeat until convergence

Alternating Optimization Revisited

SGD ALTERNATING OPTIMIZATION

1. Observed locations Ω and entries $\{X_{ij} : (i, j) \in \Omega\}$
2. Initialize $\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^m$ and $\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^n$
3. For $t = 1, 2, \dots$

1. Choose a random observed entry $(i^t, j^t) \in \Omega$
2. Update the user i^t , and item j^t

$$\mathbf{u}^{i^t} \leftarrow \mathbf{u}^{i^t} - 2\eta_t \cdot \left(\lambda \cdot \mathbf{u}^{i^t} - \left(X_{i^t j^t} - \langle \mathbf{u}^{i^t}, \mathbf{v}^{j^t} \rangle \right) \cdot \mathbf{v}^{j^t} \right)$$
$$\mathbf{v}^{j^t} \leftarrow \mathbf{v}^{j^t} - 2\eta_t \cdot \left(\lambda \cdot \mathbf{v}^{j^t} - \left(X_{i^t j^t} - \langle \mathbf{u}^{i^t}, \mathbf{v}^{j^t} \rangle \right) \cdot \mathbf{u}^{i^t} \right)$$

3. All other user and item vectors stay the same
4. Repeat until convergence

Alternating Optimization Revisited

SGD ALTERNATING OPTIMIZATION

Using simpler notation to avoid clutter

1. Observed locations Ω and entries $\{x_{ij}: (i, j) \in \Omega\}$
2. Initialize $\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^m$ and $\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^n$
3. For $t = 1, 2, \dots$

1. Choose a random observed entry $\mathcal{O}(k)$ time per iteration
2. Update the user i^t , and item j^t

$$\mathbf{u}^{i^t} \leftarrow \mathbf{u}^{i^t} - 2\eta_t \cdot \left(\lambda \cdot \mathbf{u}^{i^t} - \left(X_{i^t j^t} - \langle \mathbf{u}^{i^t}, \mathbf{v}^{j^t} \rangle \right) \cdot \mathbf{v}^{j^t} \right)$$
$$\mathbf{v}^{j^t} \leftarrow \mathbf{v}^{j^t} - 2\eta_t \cdot \left(\lambda \cdot \mathbf{v}^{j^t} - \left(X_{i^t j^t} - \langle \mathbf{u}^{i^t}, \mathbf{v}^{j^t} \rangle \right) \cdot \mathbf{u}^{i^t} \right)$$

3. All other user and item vectors stay the same
4. Repeat until convergence

Pros and Cons

Pros

- Relies purely on behavioural data
- Does not require users to submit personal information
- Does not require sellers to submit item information
- Collaborative method: lazy users benefit from active users

Cons

- Relies purely on behavioural data
- Cannot utilize user and item information even if present
- Expensive to recompute matrix factorization again and again
- Adding new users and items not simple

Please give your Feedback

<http://tinyurl.com/ml17-18afb>