

Plan for today:

- submittable stuff ✓
- Derive matrix factorization gradient descent
- Talk regularization
- Code mode
- Collaborative filtering
- Netflix challenge
- Adapting MF to CF
- More code mode
- Considerations

Matrix factorization gradient

Gradient descent

$$A = \underbrace{UV^T}_{P} = P$$

$$a_{n+1} = a_n - \alpha \nabla f(a_n)$$

α = learning rate

a_n = ~ features at current step

a_{n+1} = features for next step

$\nabla f(a_n)$ = how our error changes based on current

Really: we consider how our error changes with changing input

→ then we move in the direction of greatest decrease

or greatest decrease

$$\min_{U, V} (A - UVU^T)^2$$

$$\rightarrow \min_{u, v} \sum_{i,j} (a_{ij} - u_i v u_j^T)^2$$

↑
our error e_{ij}

$$e_{ij}^2 = (a_{ij} - \underbrace{u_i v u_j^T})^2$$

$$\frac{de_{ij}^2}{du_i} = (a_{ij} - u_i v u_j^T)(-v u_j^T)$$
$$= -\underbrace{e_{ij} v u_j^T}$$

$$\frac{de_{ij}^2}{dv} = -e_{ij} u_i v$$

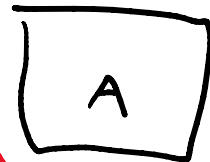
$$\frac{de_{ij}^2}{dV} = -e_{ij} u_i u_j^T$$

→ How our error is changing w.r.t. u, v

Our update equations:

$$u_i = u_i + \alpha e_{ij} v u_j^T$$

$$u_j = u_j + \alpha e_{ij} u_i v$$



$$U_j = U_j + \alpha e_{ij} U_i V$$

$$V = V + \alpha e_{ij} U_i U_j^T$$

$O(k^2)$
 $O(k^2)$

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

generally: $k \ll n$

$$U \in \mathbb{R}^{n \times k}$$

k : # of latent features

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

Note: when naively training on A

→ if we train on all $a_{ij} \in A$ including zeros

⇒ our predictor will just output 0 for everything

$O(n^2)$

→ if we train only on nonzeros

⇒ we'll just output 1 for all

$O(n)$

Workaround 1: we weight zeros and nonzeros separately

$$\min_{U, V} \sum_{\text{nonzeros in } A} (a_{ij} - U_i V U_j^T)^2 + w \sum_{\text{zeros in } A} (a_{ij} - U_i V U_j^T)^2$$

weighting to minimize training impact of zeros



Workaround 2: only train on a sample of zeros

Issue: we aren't constraining values of U, V
 \Rightarrow they blow up

Regularization

\rightarrow we'll also include values of U, V within our optimization function

$$\min_{U, V} \sum_{\text{non-zero}} (e_{ij} - U_i V_j^T)^2 + w \sum_{\text{zeros}} (a_{ij} - U_i V_j^T)^2$$

$$+ \beta_1 \|U\|_{\text{fro}} + \beta_2 \|V\|_{\text{fro}}$$

\nwarrow sum of squared values
 parameter to control relative impact of this on overall optimization
 \rightarrow additional "loss terms"

- βU_i
- βU_j
- βV

$$U - \lambda$$
$$- B V$$

Our new update equation:

$$U_i = U_i + \alpha (e_{ij} U U_j^T - B_1 U_i)$$

$$U_j = U_j + \alpha (e_{ij} U_i V - B_1 U_j)$$

$$V = V + \alpha (e_{ij} U_i U_j^T - B_2 V)$$

Collaborative filtering

One common approach for
recommender systems

Definition: Making a selection from available options (filtering) for a specific user given user preferences in general (collaborative)

Similar to predictive stuff in HW 1
→ filtering what products for a user
→ Based off purchases of 'similar' users defined via Jaccard and other

and other

General approaches for C.F.:

→ like HW: define explicit similarities

→ here today: CF via MF

→ machine learning:

- define explicit features
- can better capture non-linear behavior that MF can't
- Issue for us: lack of quality data

Reason: privacy concern,
easy to de-anonymize

Netflix Challenge

~2007 or so

Netflix: we have a predictor
for user-movie preferences

→ predicts (1-5) rating for given user

Challenge:

→ improve predictor by 10%

↑ → a hundred 100s of millions of ratings

- improve p
- Released 100s of millions of ratings
- Had an internal test set for rating
- If improve by 10% → you get \$1,000,000
- Took a couple years for teams to do so
 - winners used aggregation of ML approaches
 - MF itself got up to 7%

Interesting notes:

Temporal effects of ratings

- some movies if rated immediately vs. later got higher (Patch Adams)
- some had opposite effect (Memento)

Naive approach:

$$a_{ij} = 0.5(\underset{\substack{\uparrow \\ \text{avg for user } i}}{j}}{\text{avg}(a_{ij})} + \underset{\substack{\uparrow \\ \text{avg for movie } j}}{i}}{\text{avg}(a_{ij})})$$

Adapting MF for CF (MF4CF)

consider user-movie matrix

↳ our update equations:

$$U_i = U_i + \alpha (e_{ij} V_j - \beta_1 U_i)$$

$$V_j = V_j + \alpha (e_{ij} U_i - \beta_2 V_j)$$

Considerations and Challenges

Sparse data

→ much fewer ratings than non-ratings

Scale of data

→ how to adapt

Generalization:

→ how do we work on novel data

→ be careful not to overfitting

"Cold-start":

→ How to predict for a new user?

Black sheep:

→ some users you can't predict

→ all they are is noise