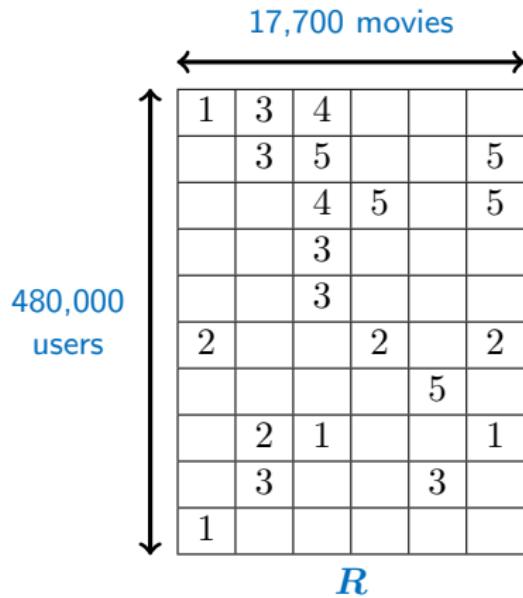


Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. Principal Component Analysis (PCA)
3. Singular Value Decomposition (SVD)
4. **Matrix Factorization**
 - Motivation & Approach
 - Regularization & Sparsity
 - Further Factorization Models
5. Neighbor Graph Methods
6. Autoencoders (Nonlinear Dimensionality Reduction)

Motivation: The Netflix Prize

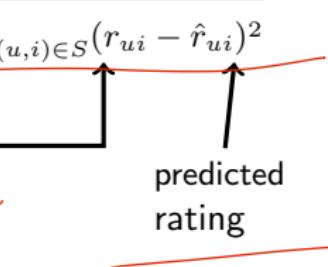
- Training data
 - 100 million ratings, 480,000 users, 17,770 movies
 - 6 years of data: 2000-2005
- Test data
 - Last few ratings of each user (2.8 million)
 - Root Mean Square Error (RMSE)
 - Netflix's system RMSE: 0.9514
 $\text{trival} = \overbrace{1.054}$
- Competition
 - 2,700+ teams
 - \$1 million prize for 10% improvement



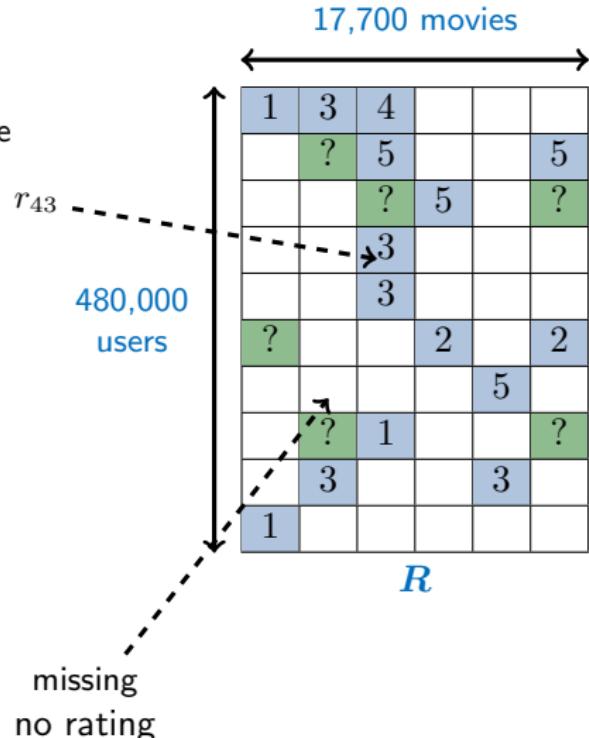
Evaluating Recommender Systems

- $S = \text{set of tuples } (u, i) \text{ of users } u \text{ that have rated item } i \text{ with a rating of } r_{ui}$

- $RMSE = \sqrt{\frac{1}{|S|} \sum_{(u,i) \in S} (r_{ui} - \hat{r}_{ui})^2}$

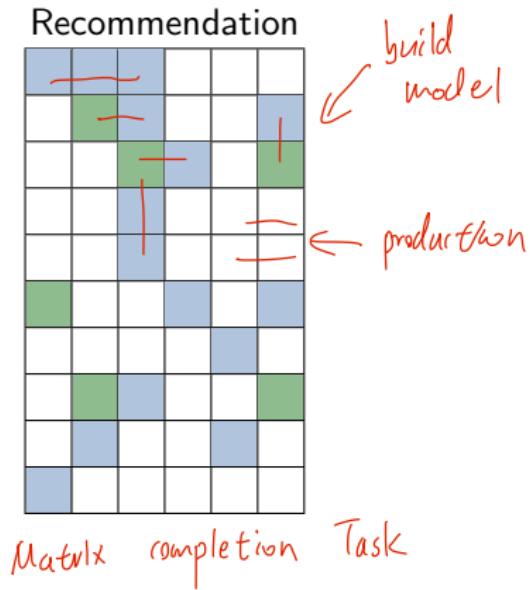
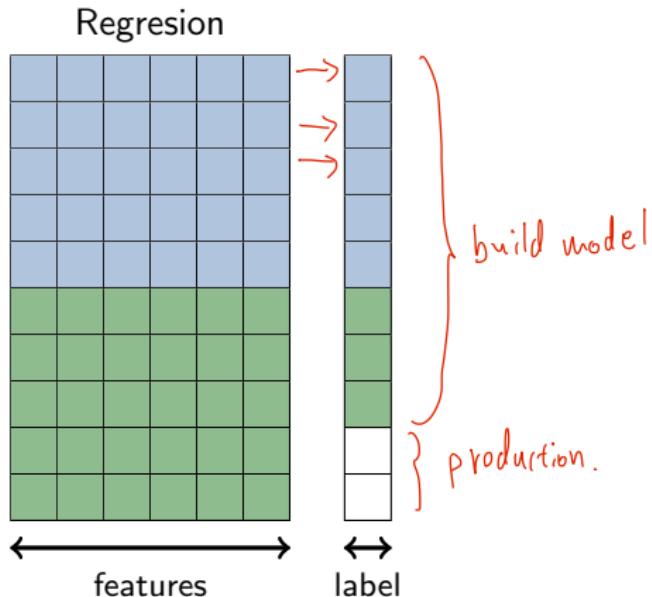


- Legend:
 - Training and validation data
 - Test data
 - Missing data



Regression vs. Recommendation

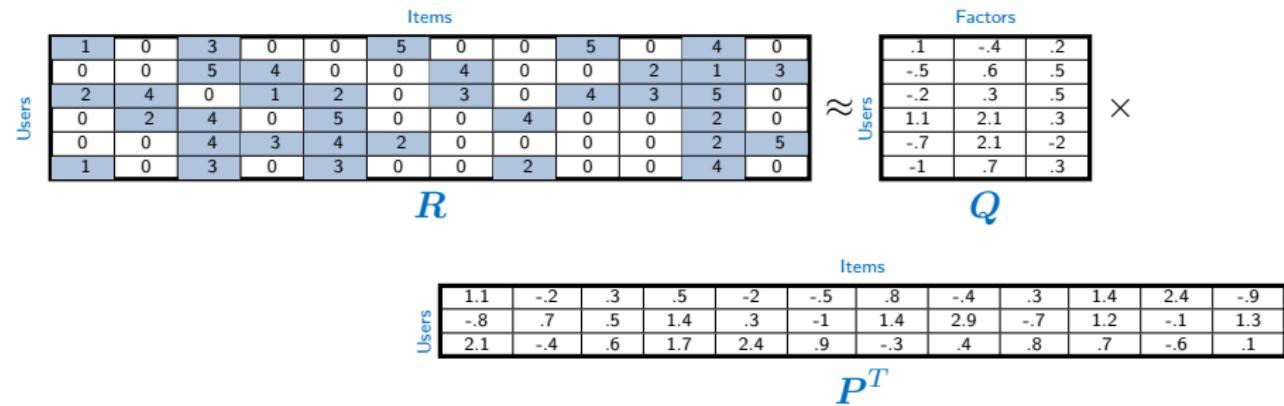
- Legend:
 - Training and validation data
 - Test data
 - Missing data



SVD on Rating Data

- Goal: Make good recommendations
 - Good performance on observed (user, item) ratings, i.e. low RMSE
 - Generalize to the unseen test data
- Can we use SVD to obtain the solution?
 - SVD on the rating matrix $R \in \mathbb{R}^{n \times d}$ where we replace **missing** entries with **zeros**
 - $R \approx Q \cdot P^T$ // SVD: $R = U\Sigma V^T \rightarrow Q = U\Sigma, P = V$

$$\boxed{\square} = \boxed{\text{?}} \quad \text{DRAFT!}$$



SVD on Rating Data

- SVD gives **minimum reconstruction error** (Sum of Squared Errors):

$$\min_{U, \Sigma, V} \sum_{\substack{u=1 \dots n \\ i=1 \dots d}} \underbrace{\left(R_{ui} - [U \Sigma V^T]_{ui} \right)^2}_{B} \quad \begin{array}{l} \min_B \|R - B\|_F^2 \\ \text{rank}(B) = k \end{array}$$

- SSE and RMSE are monotonically related:
 - $RMSE = \frac{1}{const.} \sqrt{SSE}$
 - Great news: SVD is minimizing RMSE
- **Complication:**
 - The sum in the SVD error term is over all entries (no rating = zero rating)
 - But our R has **missing** entries!
 - (Classical) SVD isn't defined when entries are missing!
- Also: We actually don't care about orthogonality and normalization

Discussion: SVD/PCA

- Optimal low rank approximation
 - In terms of Frobenius norm (and also in terms of the spectral norm)
- Missing elements (we have no information/not observed) are treated as 0 (a very low rating)
 - Critical for many applications (e.g. recommender systems)
 - General problem: two kinds of "zeros" (not observed/missing $\neq 0$)
- Lack of Sparsity
 - Singular vectors are dense
 - Potential interpretability problem
- Orthogonality really required/useful?

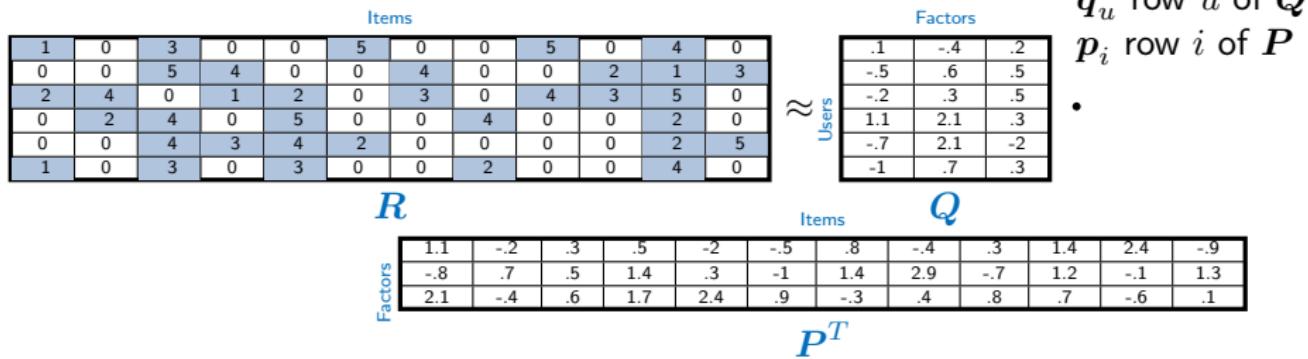
$$\begin{matrix} \bullet & \bullet \\ \bullet & \bullet \\ \bullet & \bullet \\ \bullet & \bullet \end{matrix} = \begin{matrix} U \end{matrix} \Sigma \begin{matrix} V^T \end{matrix}$$

Latent Factor Models

- Our goal: Find $\mathbf{Q} \in \mathbb{R}^{n \times k}$ and $\mathbf{P} \in \mathbb{R}^{d \times k}$ such that:

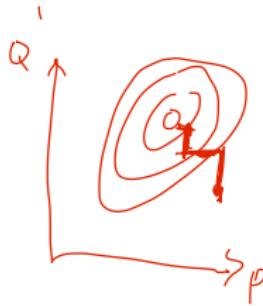
$$\min_{\mathbf{P}, \mathbf{Q}} \sum_{(u,i) \in S} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2$$

- We only sum over existing entries, i.e. the set $S = \{(u, i) | r_{ui} \neq \text{missing}\}$
- We don't require columns of \mathbf{P}, \mathbf{Q} to be orthogonal/ unit length
- Use standard optimization techniques to solve this problem



Finding the Latent Factors (I)

- Goal: $\min_{\mathbf{P}, \mathbf{Q}} \sum_{(u,i) \in S} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2 =: \min_{\mathbf{P}, \mathbf{Q}} f(\mathbf{P}, \mathbf{Q})$
- One approach: **Alternating Optimization** //a.k.a. block coordinate minimization
 - Pick initial values for \mathbf{P} and \mathbf{Q}
 - Alternately keep one variable fix and optimize for the other
 - Repeat until convergence
- Pseudo Code:
 1. Initialize $\mathbf{P}^{(0)}, \mathbf{Q}^{(0)}, t = 0$
 2. $\mathbf{P}^{(t+1)} = \arg \min_{\mathbf{P}} f(\mathbf{P}, \mathbf{Q}^{(t)})$
 3. $\mathbf{Q}^{(t+1)} = \arg \min_{\mathbf{Q}} f(\mathbf{P}^{t+1}, \mathbf{Q})$
 4. $t = t + 1$
 5. goto 2 until convergence



Finding the Latent Factors (II)

- Initialization of P and Q :
 - Use, e.g., SVD where missing entries are replaced by 0 (or overall mean)

- How to solve

$$\mathbf{P}^{(t+1)} = \arg \min_{\mathbf{P}} f(\mathbf{P}, \mathbf{A}^{(t)}) = \arg \min_{\mathbf{P}} \sum_{(u,i) \in S} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2$$

- Observation 1: Since Q is fixed, the problem can be solved for **each** vector \mathbf{p}_i **independently**

$$\min_{\mathbf{P}} \sum_{(u,i) \in S} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2 = \sum_{i=1 \dots d} \min_{\mathbf{p}_i} \sum_{u \in S_{*,i}} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2$$

where $S_{*,i} = \{u | (u, i) \in S\}$ // = only users u who have rated item i

- Equivalently for vector \mathbf{q}_u based on the set $S_{*,i} = \{u | (u, i) \in S\}$

$$\Rightarrow = \sum_{(u,1) \in S} (r_{u1} - q_u \cdot p_1^T)^2 + \sum_{(u,2) \in S} (r_{u2} - q_u \cdot p_2^T)^2 + \dots$$

Finding the Latent Factors (III)

- *Observation 2:* $\min_{\mathbf{p}_i} \sum_{u \in S_{*,i}} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2$ is an **ordinary least squares regression** problem
 - $\min_{\mathbf{w}} \sum_{j=1}^n (y_j - \mathbf{w}^T \mathbf{x}_j)$ // y_j is a scalar, \mathbf{x}_j and \mathbf{w} (column) vectors
 - Optimal solution in closed form:
$$\mathbf{w} = \left(\frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^T \right)^{-1} \cdot \left(\frac{1}{n} \sum_{j=1}^n \mathbf{x}_j y_j^T \right)^{-1}$$
$$W = \left(\frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^T \right)^{-1} \cdot \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j y_j$$
$$= (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$
 - In our case: $\mathbf{p}_i^T = \left(\frac{1}{|S_{*,i}|} \sum_{u \in S_{*,i}} \mathbf{q}_u^T \mathbf{q}_u \right)^{-1} \cdot \frac{1}{|S_{*,i}|} \sum_{u \in S_{*,i}} \mathbf{q}_u^T r_{ui}$
 - Equivalently for \mathbf{q}_u
- Computation of $\arg \min_{\mathbf{P}, \mathbf{Q}} \left(\mathbf{P}, \mathbf{Q}^{(t)} \right)$ reduces to a standard problem

Alternating Optimization: Discussion

- May provide solution to difficult optimization problems
- Here: sequence of simple OLS problems
 - Overall algorithm can be implemented in a few lines in, e.g., Python
 - Quite efficient: since data is sparse, the sets $S_{*,i}$ (and $S_{u,*}$) are relatively small
- **Drawback** of Alternating Optimization:
 - Solution is only an approximation
 - No guarantee that close to the optimal solution
 - Highly depends on initial solution

SGD for Matrix Factorization

- Stochastic Gradient Descent as an alternative
- SGD on the objective $\mathcal{L} := \sum_{(u,i) \in S} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2$
 - Pick a random user u and a random item i with rating r_{ui} (batch size 1)
 - Compute the gradients w.r.t. the parameters $\frac{\partial \mathcal{L}}{\partial \mathbf{q}_u}$ and $\frac{\partial \mathcal{L}}{\partial \mathbf{p}_i}$
 - Update the parameters $\mathbf{q}_u \leftarrow \mathbf{q}_u - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{q}_u}$, $\mathbf{p}_i \leftarrow \mathbf{p}_i - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{p}_i}$
 - η - learning rate
- In this case the gradient update step is rather simple
 - $e_{ui} \leftarrow r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T$ \\ helper variable, the current error
 - $\mathbf{q}_u \leftarrow \mathbf{q}_u + 2\eta(e_{ui}\mathbf{p}_i)$
 - $\mathbf{p}_i \leftarrow \mathbf{p}_i + 2\eta(e_{ui}\mathbf{q}_u)$

Rating Prediction

- How to estimate the missing rating of user u for item i ?

		Items								
		1	2	3	4	5	6	7	8	9
Users	1				5			5		4
	2			5		4		2	1	3
	2	4		1		3		4	3	5
		2	4		5		4		2	
			4	3	4	2			2	5
	1		3		3		2		4	

$$\hat{r}_{ui} = \mathbf{q}_u \cdot \mathbf{p}_i^T = \sum_k \mathbf{q}_{uk} \cdot \mathbf{p}_{ik}$$

\mathbf{q}_u row u of \mathbf{Q}
 \mathbf{p}_i row i of \mathbf{P}

\approx

R

Factors

.1	-.4	.2
-.5	.6	.5
-.2	.3	.5
1.1	2.1	.3
-.7	2.1	-2
-1	.7	.3

Q

• Factors

Items												
	1.1	-.2	.3	.5	-.2	-.5	.8	-.4	.3	1.4	2.4	-.9
1	-.8	.7	.5	1.4	.3	-1	1.4	2.9	-.7	1.2	-.1	1.3
2	2.1	-.4	.6	1.7	2.4	.9	-.3	.4	.8	.7	-.6	.1

P^T

Rating Prediction

- How to estimate the missing rating of user u for item i ?

Users	Items								
	1		3	5	4	5	2	1	3
2	4	5	1	?	3	4	3	5	
3	2	4		5		4		2	
4		4	3	4	2			2	5
5	1		3		2			4	

$$\hat{r}_{ui} = \mathbf{q}_u \cdot \mathbf{p}_i^T = \sum_k \mathbf{q}_{uk} \cdot \mathbf{p}_{ik}$$

\mathbf{q}_u row u of \mathbf{Q}
 \mathbf{p}_i row i of \mathbf{P}

R

Factors User 2

Users	Factors	
	1	2
1	-.4	.2
2	.6	.5
3	.3	.5
4	2.1	-3
5	2.1	-2
6	.7	.3

• Factors Item 2

Users	Items								
	1.1	-.2	.3	.5	-2	-.5	.8	-.4	.3
1	-.8	.7	.5	1.4	.3	-1	1.4	2.9	-.7
2	2.1	-.4	.6	1.7	2.4	.9	-.3	.4	.8
3									
4									
5									
6									

P^T

Rating Prediction

- How to estimate the missing rating of user u for item i ?

Users	Items								
	1		3		5		5		4
		5	4		4		2	1	3
2	4		1		3		4	3	5
	2	4		5		4		2	
	4	3	4	2				2	5
1		3		3		2		4	

$$\hat{r}_{ui} = \mathbf{q}_u \cdot \mathbf{p}_i^T = \sum_k \mathbf{q}_{uk} \cdot \mathbf{p}_{ik}$$

\mathbf{q}_u row u of \mathbf{Q}
 \mathbf{p}_i row i of \mathbf{P}

\approx

R

$$\begin{aligned} & -0.5 \times -2 + 0.6 \times 0.3 + 0.5 \times 2.4 \\ & = 1 + 0.18 + 1.2 = 2.38 \end{aligned}$$

Factors

.1	-.4	.2
-.5	.6	.5
-.2	.3	.5
1.1	2.1	.3
-.7	2.1	-2
-1	.7	.3

Q

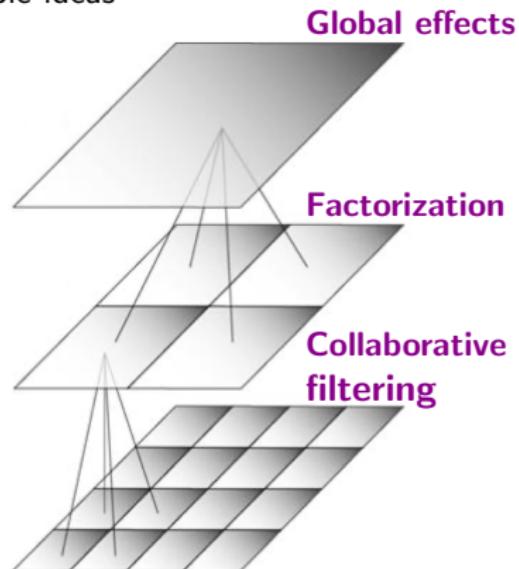
• Factors

Users	Items											
	1.1	-.2	.3	.5	-2	-.5	.8	-.4	.3	1.4	2.4	-.9
	-.8	.7	.5	1.4	.3	-1	1.4	2.9	-.7	1.2	-.1	1.3
	2.1	-.4	.6	1.7	2.4	.9	-.3	.4	.8	.7	-.6	.1

P^T

Side note: BellKor Recommender System

- The winner of the Netflix Challenge uses matrix factorization as **one** building block
 - The overall model is a combination of multiple ideas
- Multi scale modeling of the data:
Combine top level, 'regional' modeling of the data, with a refined, local view:
 - **Global:**
 - Overall deviations of users/movies
 - **Factorization:**
 - Addressing "regional" effects
 - **Collaborative filtering:**
 - Extract local patterns



User and Item Biases

- Certain users might give overly optimistic or pessimistic rating
- Certain movies tend to always have low ratings
- We introduce additional bias terms to capture these effects
 - Each user u can have a bias term b_u
 - Each user i can have a bias term b_i
 - Additional global bias term b shared by everyone
- The resulting optimization problem can be easily solved with SGD

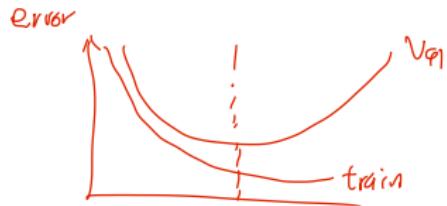
$$\min_{\substack{\mathbf{P}, \mathbf{Q} \\ b_u, b_i, b}} \sum_{(u,i) \in S} (r_{ui} - (\mathbf{q}_u \cdot \mathbf{p}_i^T + b_u + b_i + b))^2$$

- The **cold start** problem is another important issue

Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. Principal Component Analysis (PCA)
3. Singular Value Decomposition (SVD)
4. **Matrix Factorization**
 - Motivation & Approach
 - **Regularization & Sparsity**
 - Further Factorization Models
5. Neighbor Graph Methods
6. Autoencoders (Nonlinear Dimensionality Reduction)

Challenge: Overfitting



- Final Goal: Minimize SSE for **unseen test data**
- Proxy: Minimize SSE on **training data**
 - Want large k (number of factors) to capture all the signals
 - But, SSE on **test data** begins to rise for larger k
 - Why?
- Classical example of overfitting
 - With too many degrees freedom (too many free parameters) the model starts fitting noise
 - The model fits too well the training data but does not generalize well to unseen test data

Challenge: Overfitting

- Problem can easily be seen in our scenario:
 - In each step of the alternating optimization we solve an OLS regression
 - Number of regression parameters = k = number of latent factors
 - Number of data points used for regression = cardinality of $S_{*,i}/S_{u,*}$
 - Lots of parameters but not enough data points → regression can overfit
 - If k is large this might even lead to an **underdetermined** system of equations
 - Problem is ill posed
- Solution: Regularization
 - Interpretation for underdetermined systems: "In the absence of any other information, the parameter vector should be small (i.e. only small effect of features)"
 - Equivalently: We put a prior on the weights

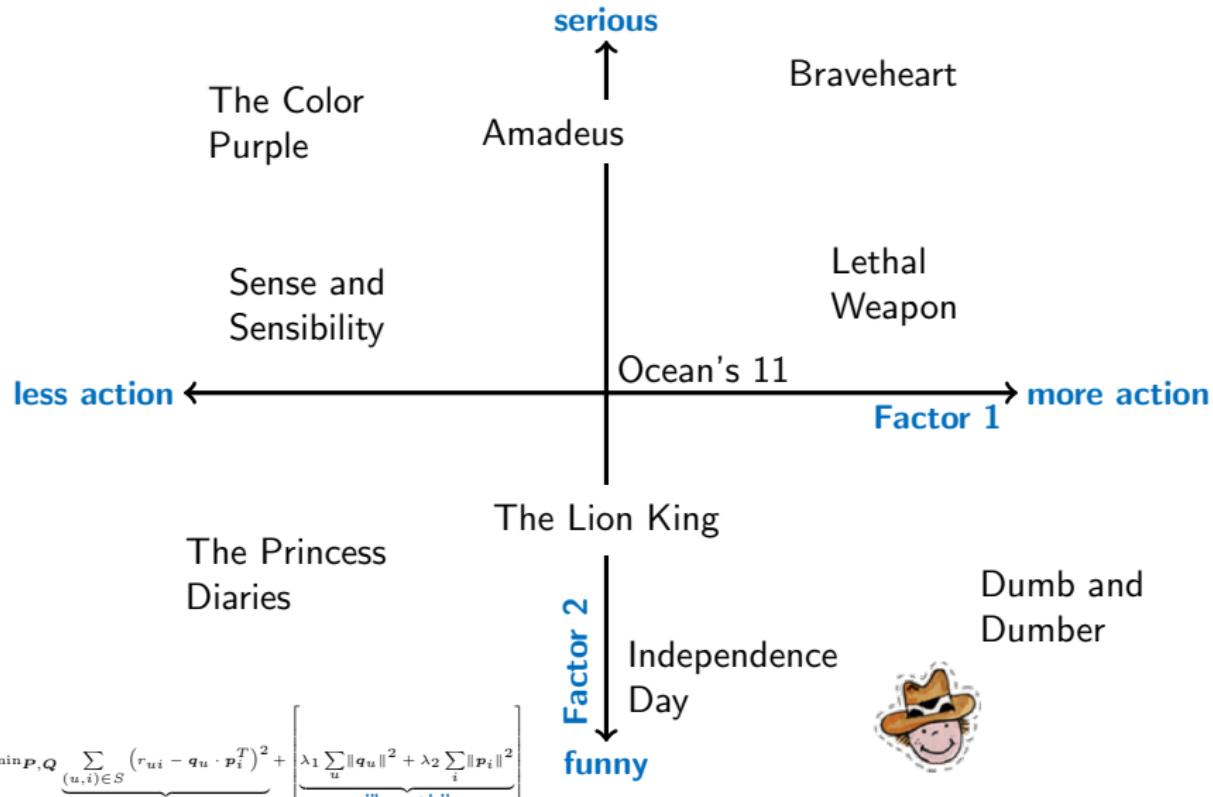
Regularization

- To solve overfitting we introduce regularization:
 - Allow rich model where we have sufficient data
 - Shrink aggressively where data are scarce

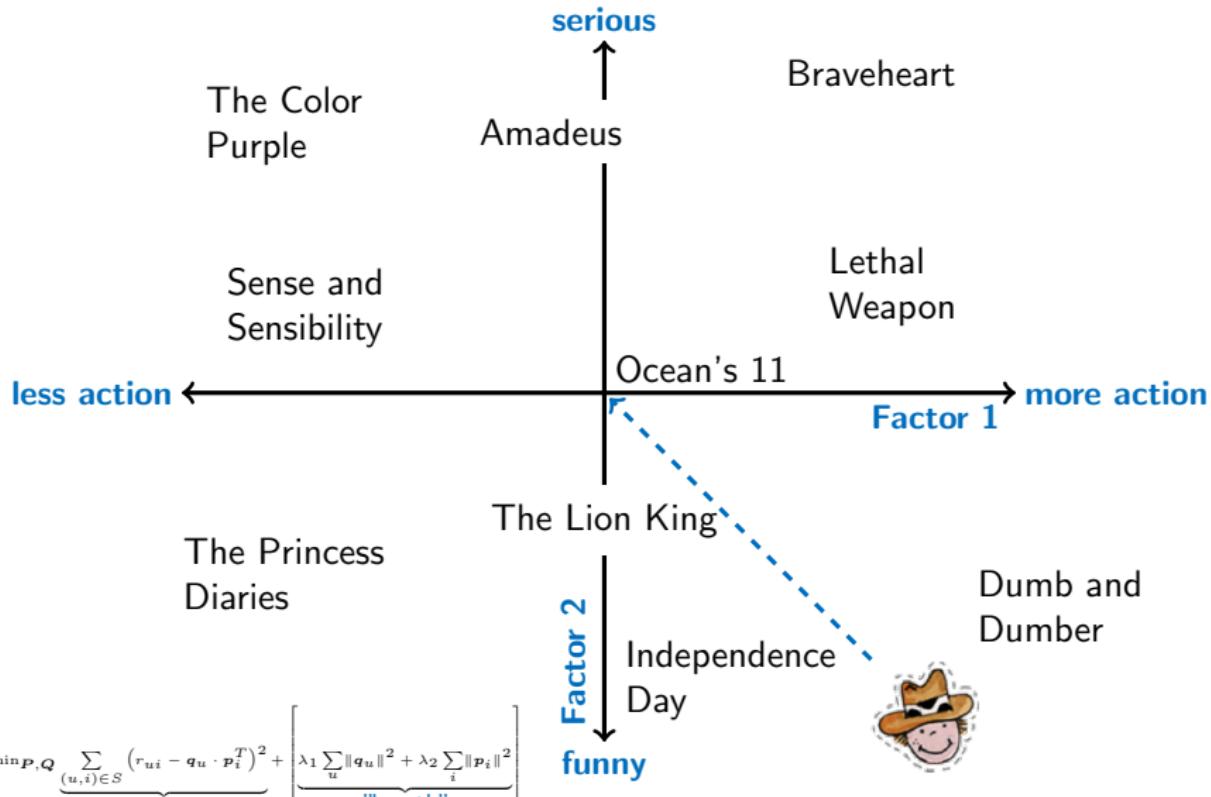
$$\min_{P, Q} \underbrace{\sum_{(u,i) \in S} \left(r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T \right)^2}_{\text{"error"}} + \left[\underbrace{\lambda_1 \sum_u \|\mathbf{q}_u\|^2}_{\text{"length"}^{\text{"}}} + \underbrace{\lambda_2 \sum_i \|\mathbf{p}_i\|^2}_{\text{"length"}^{\text{"}}} \right]$$

- λ_1 and λ_2 are user-defined regularization parameters
- Note: We do not care about the actual value of the objective function, but we care about the P, Q that achieve the minimum of the objective

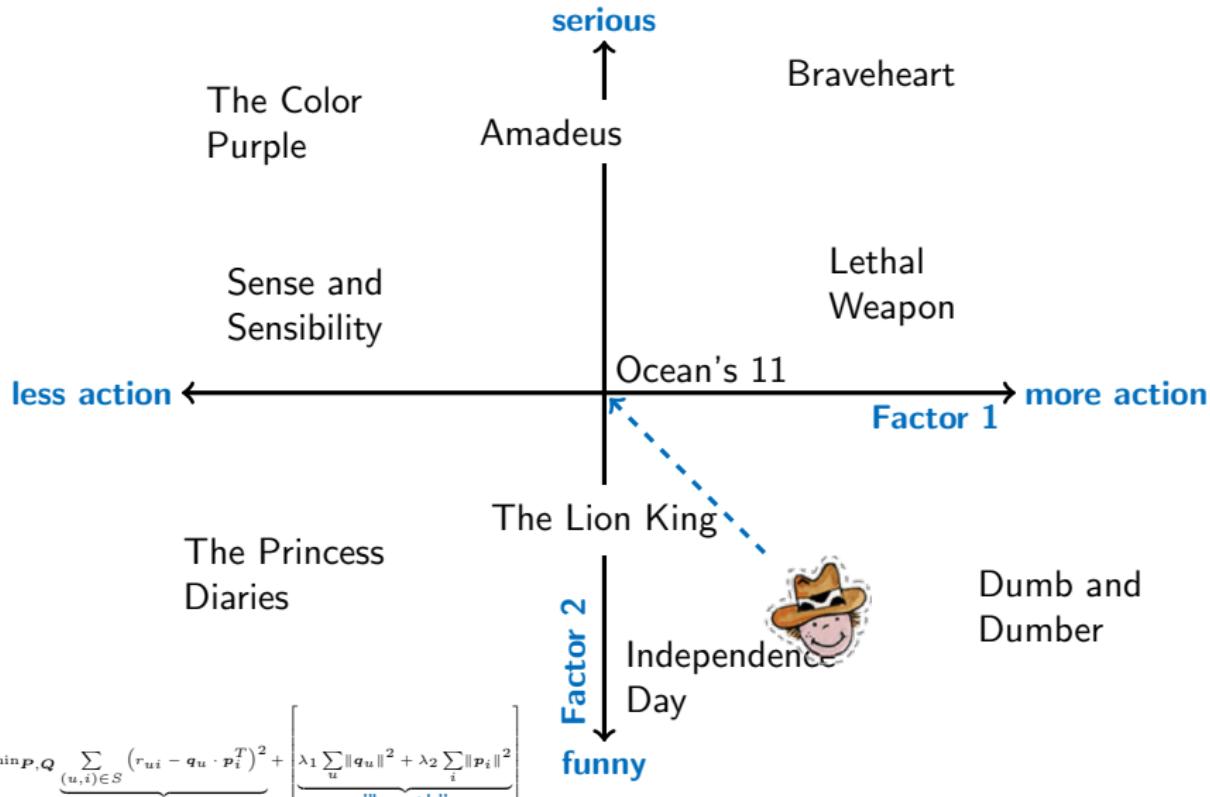
The Effect of Regularization



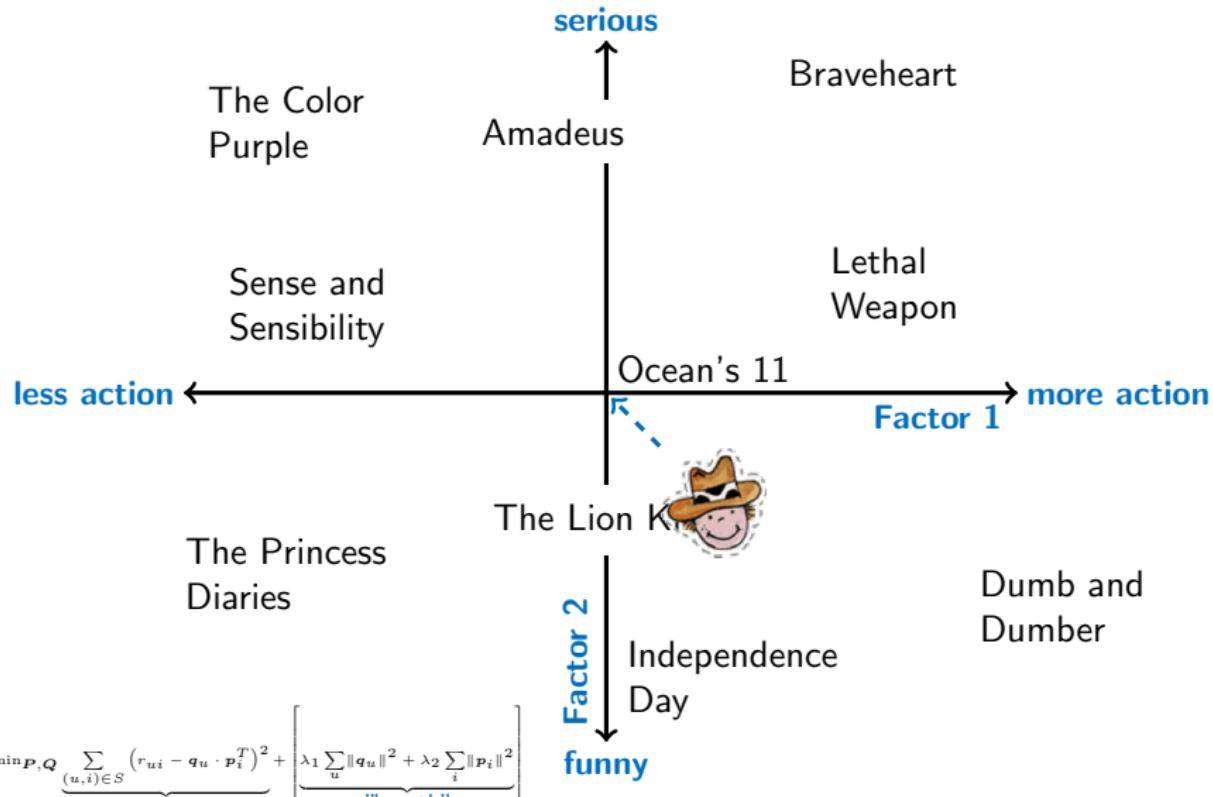
The Effect of Regularization



The Effect of Regularization



The Effect of Regularization



Regularization in Our Use Case

- Regularized Problem:

$$\min_{P, Q} \sum_{(u,i) \in S} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2 + \left[\lambda_1 \sum_u \|\mathbf{q}_u\|^2 + \lambda_2 \sum_i \|\mathbf{p}_i\|^2 \right]$$

- Recap of unregularized problem: In each iteration of the alternating optimization we solved an OLS regression problem
- For the regularized version this becomes ridge regression
 - $\min_{\mathbf{p}_i} \sum_{i \in S_{*,i}} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2 + \lambda_2 \|\mathbf{p}_i\|^2$
 - Effect: parameter values are forced to become smaller
 - Large values that only capture noise are avoided
- You know how to solve this!
 - Closed form solution (see linear regression slides)
 - Gradient decent

$$(\mathbf{x}^T \mathbf{x} + \lambda I)^{-1} \mathbf{x}^T \mathbf{y}$$

SGD for MF with Regularization and Biases

- SGD on the objective

$$\mathcal{L} := \min_{\substack{\mathbf{P}, \mathbf{Q} \\ b_u, b_i, b}} \sum_{(u,i) \in S} \left(r_{ui} - (\mathbf{q}_u \cdot \mathbf{p}_i^T + b_u + b_i + b) \right)^2 + \left[\lambda_1 \sum_u \|\mathbf{q}_u\|^2 + \lambda_2 \sum_u \|\mathbf{p}_u\|^2 \right]$$

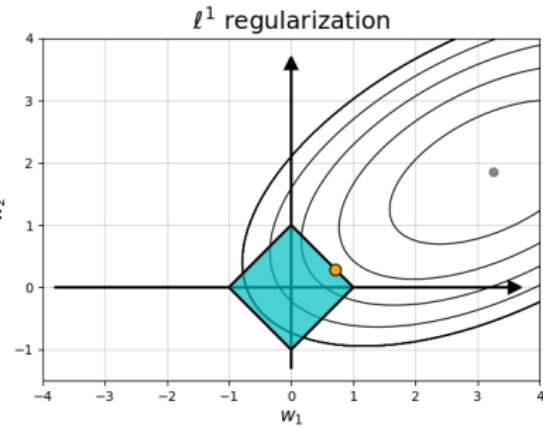
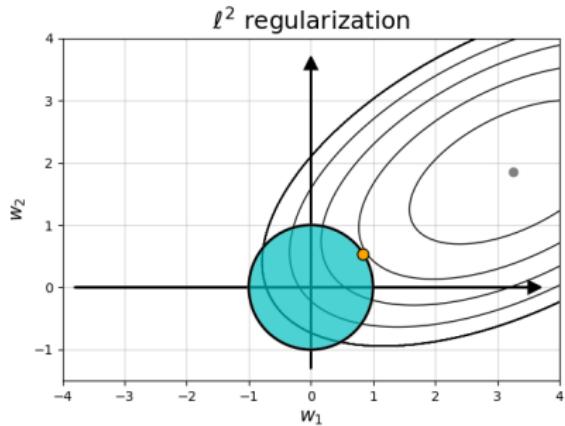
- Pick a random user u and a random item i with rating r_{ui} (batch size 1)
- The gradient update step is very similar to before
 - $e_{ui} \leftarrow r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T$ \\ \\\backslash helper variable, the current error
 - $\mathbf{q}_u \leftarrow \mathbf{q}_u + 2\eta (e_{ui} \mathbf{p}_i - \lambda_1 \mathbf{q}_u)$
 - $\mathbf{p}_i \leftarrow \mathbf{p}_i + 2\eta (e_{ui} \mathbf{q}_u - \lambda_2 \mathbf{p}_i)$
- $b_u \leftarrow b_u + 2\eta e_{ui}$
- $b_i \leftarrow b_i + 2\eta e_{ui}$
- Usually directly set b to the global mean $b = \frac{1}{|S|} \sum_{(u,i) \in S} r_{ui}$

L2 vs. L1 Regularization

- Comparison: L2 vs. L1 regularization
 - L2 tries to "shrink" the parameter vector w **equally**
 - Large values are highly penalized due to the square in the L2 norm
 - Unlikely that any component will be **exactly 0**
 - L1 allows large values in individual components of w by shrinking other components to 0
 - L1 is suited to enforce **sparsity** of the parameter vector
- Why sparsity
 - Better interpretation
 - Only few values contribute to result
 - Unintuitive that **sparse input data** is generated based on **dense signal**
 - Less storage, faster processing

L2 vs. L1 Intuition

- An L 1 constraint promotes sparsity



by @itayevron

Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. Principal Component Analysis (PCA)
3. Singular Value Decomposition (SVD)
4. **Matrix Factorization**
 - Motivation & Approach
 - Regularization & Sparsity
 - **Further Factorization Models**
5. Neighbor Graph Methods
6. Autoencoders (Nonlinear Dimensionality Reduction)

Further Factorization Models

- Matrix factorization is extremely powerful
 - Dimensionality reduction, data analysis/data understanding, prediction of missing values, ...
- Various extensions have been proposed
 - Enforcing different constraints or operating on different data types
 - Important goal: better interpretation of result
- Non Negative Matrix Factorization (next slides)
- Boolean Matrix Factorization $0, 1$
 - Factorize Boolean A in Boolean Q and P

Non-Negative Matrix Factorization

- Often data is given in form of non negative values
 - Rating values between 1 to 5; income, age, ... of persons; number of words in a document; grayscale images; etc.
- However: SVD (and the other approaches presented before) might lead to factors containing **negative values**
 - Difficult to interpret: non negative data is "generated" based on negative factors; what do these negative factors mean?
 - Predicted values might also become negative
- Solution: Non Negative Matrix Factorization

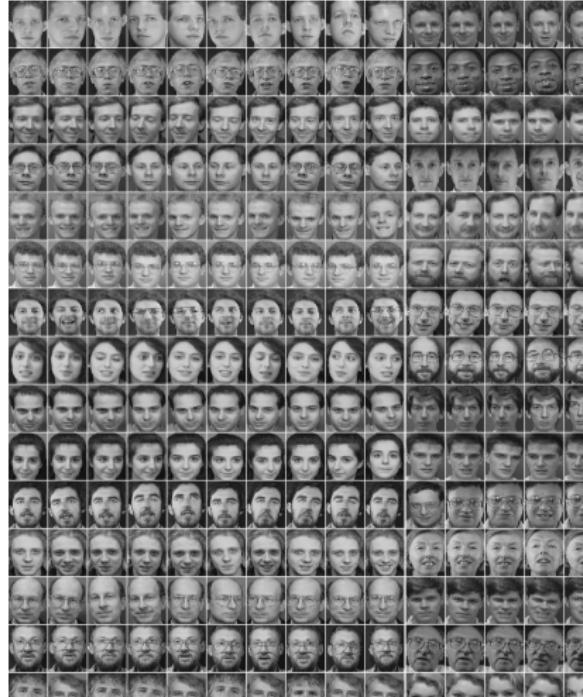
Non-Negative Matrix Factorization

- Task: Factorize non negative \mathbf{A} in non negative \mathbf{Q} and \mathbf{P} , i.e.
$$\mathbf{A} \approx \mathbf{Q} \cdot \mathbf{P}^T$$
- Formally:
 - Given $\mathbf{A} \in \mathbb{R}^{n \times d}$ with $A_{ij} \geq 0$ and integer k , find $\mathbf{P} \in \mathbb{R}^{n \times k}$,
 $\mathbf{Q} \in \mathbb{R}^{k \times d}$ such that $\|\mathbf{A} - \mathbf{Q} \cdot \mathbf{P}^T\|_F$ is minimized subject to $\mathbf{Q} \geq \mathbf{0}$ and $\mathbf{P} \geq \mathbf{0}$

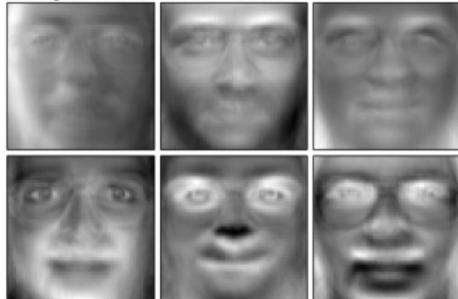
$$\min_{\mathbf{P} \geq \mathbf{0}, \mathbf{Q} \geq \mathbf{0}} \|\mathbf{A} - \mathbf{Q} \cdot \mathbf{P}^T\|_F$$

- Constrained optimization

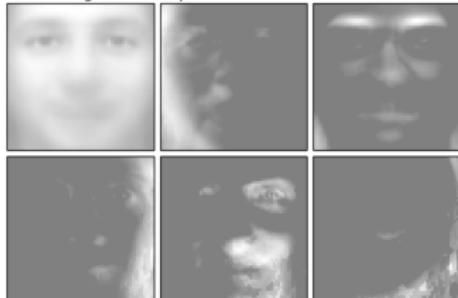
NNMF Example: Olivetti Faces Data



Eigenfaces - RandomizedPCA - Train time 1.4s



Non-negative components - NMF - Train time 2.9s



Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. Principal Component Analysis (PCA)
3. Singular Value Decomposition (SVD)
4. Matrix Factorization
5. **Neighbor Graph Methods**
6. Autoencoders (Nonlinear Dimensionality Reduction)

Preserving global vs. preserving local similarity

- PCA tries to find a global structure
 - Can lead to local inconsistencies
 - Far away point can become nearest neighbors
- Illustration (during lecture):



- Idea: Preserve local structure instead
- Example: <https://projector.tensorflow.org/>

Neighbor Graph Methods

- All methods we have seen so far are based on matrix factorizations
- An alternative class of methods based on **neighbor graphs**

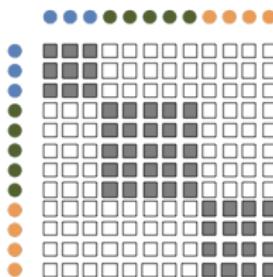
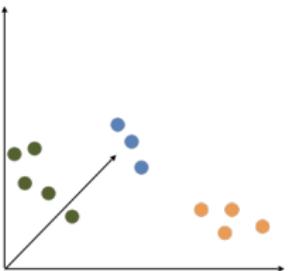


- Idea: Low dimensional neighborhood similar to the original neighborhood

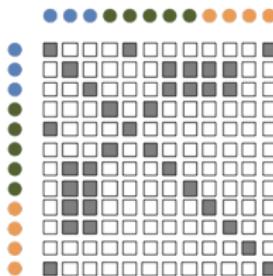
Neighbor Graph Methods

 - High Similarity / Low distance
 - low Similarity / High distance

1. Construct neighbor graph of high dimensional data
2. Initialize points in low dimensional space
3. Optimize coordinates in low dimensional space s.t. similarities align



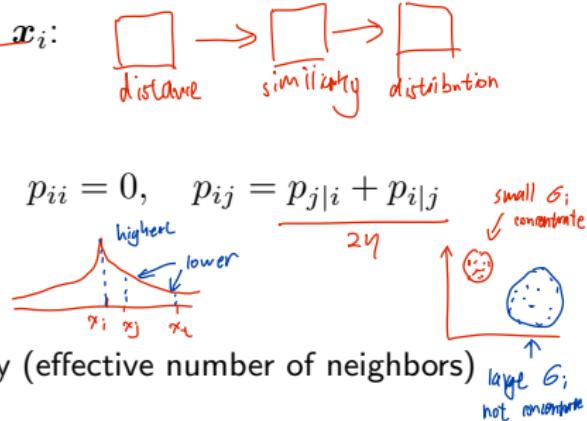
pairwise
Similarity matrix
pair wise distance matrix



tSNE : t Distributed Stochastic Neighbor Embedding

- High dimensional similarities for input x_i :

$$p_{j|i} = \frac{\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}\right)}$$



- σ_i chosen to achieve fixed perplexity (effective number of neighbors)
n. hyperparameter
- Low dimensional similarities for (to be optimized) parameters y_i :

$$q_{i,j} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|y_k - y_l\|^2)^{-1}} \quad q_{ii} = 0$$

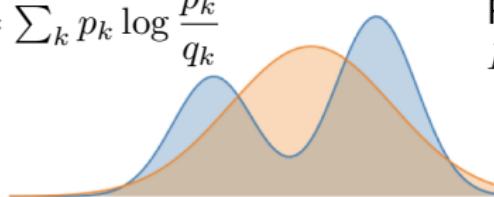
- Minimize the KL divergence:

$$\min_{y_i} KL(P||Q) = \sum_i \sum_{j \neq i} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

Note on the KL Divergence

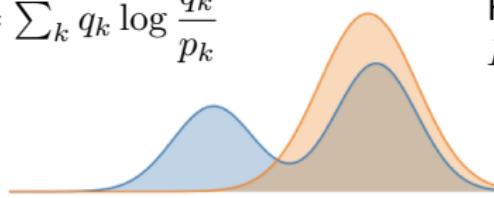
- $KL(P||Q) \geq 0$ for any P and Q , $KL(P||Q) = 0$ iff $P = Q$
- The KL Divergence is asymmetric $KL(P||Q) \neq KL(Q||P)$
 - Example: Given P we are optimizing over Q

$$\bullet \quad KL(P||Q) = \sum_k p_k \log \frac{p_k}{q_k}$$



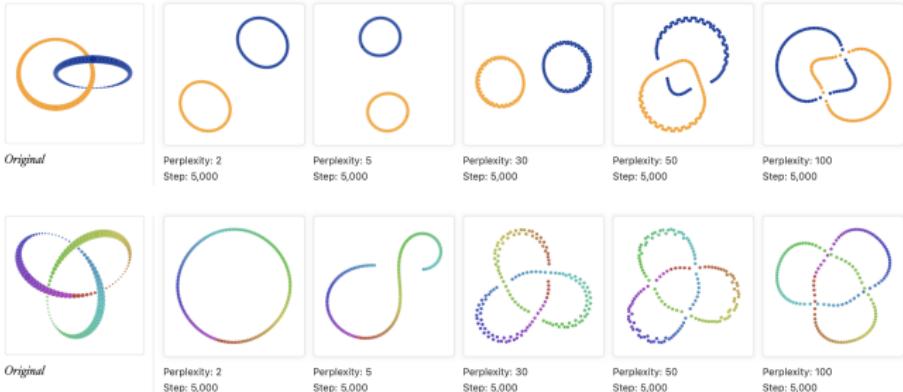
Forward KL is Mean Seeking
 $KL(P||Q)$

$$\bullet \quad KL(Q||P) = \sum_k q_k \log \frac{q_k}{p_k}$$



Reverse KL is Mode Seeking
 $KL(Q||P)$

Interactive tSNE



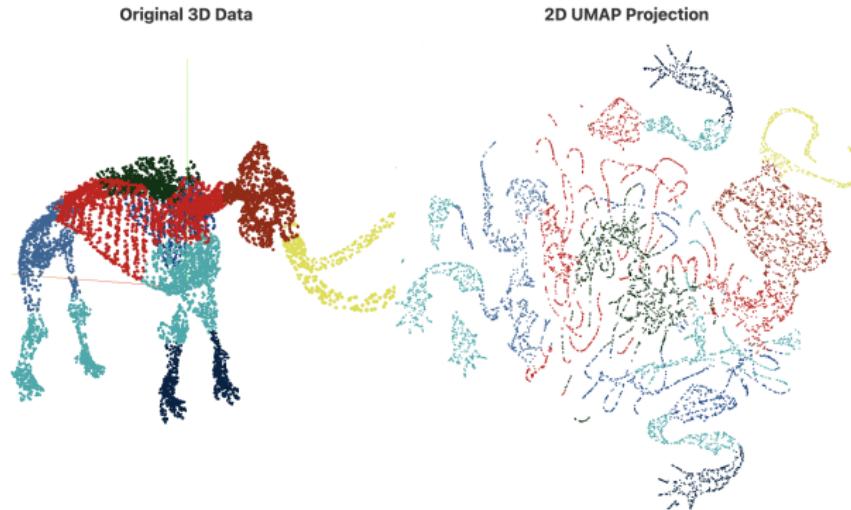
Source of illustration, for more details and interactive widgets

The Advantages and Disadvantages of t-SNE

- The current standard for visualizing high dimensional data
 - Helps understand "black box" algorithms like DNN
 - Reduced "crowding problem" with heavy tailed distribution
-
- t-SNE plots can sometimes be mysterious or misleading
 - Be careful with cluster sizes and cluster distances
 - Can be sensitive to hyperparameters
 - Random noise does not always look random
 - No easy way to compute the embedding of new data
 - Not great for more than 3 dimensions

Uniform Manifold Approximation and Projection

- Principled approach relying on topological spaces, category theory, ...



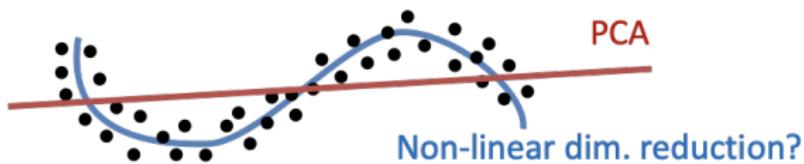
[Link for details, widgets, a comparison to T SNE, ...](#)

Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. Principal Component Analysis (PCA)
3. Singular Value Decomposition (SVD)
4. Matrix Factorization
5. Neighbor Graph Methods
6. **Autoencoders (Nonlinear Dimensionality Reduction)**

Motivation

- PCA / SVD can only capture linear structure (linear variations in the data)
 - Recall: transformed data is $\tilde{\mathbf{Y}} = \tilde{\mathbf{X}} \cdot \Gamma$
 - Linear projection by the eigenvectors Γ
- However, data often lies on a non linear low dimensional manifold



- Idea: find a non linear projection of the data

Autoencoders



- An **autoencoder** is a neural network that:
 - finds a compact representation of the data
 - by learning to reconstruct its input

$$f(\mathbf{x}, \mathbf{W}) := \hat{\mathbf{x}} \approx \mathbf{x}$$

- Goal: minimize the reconstruction error between $\hat{\mathbf{x}}$ and \mathbf{x}

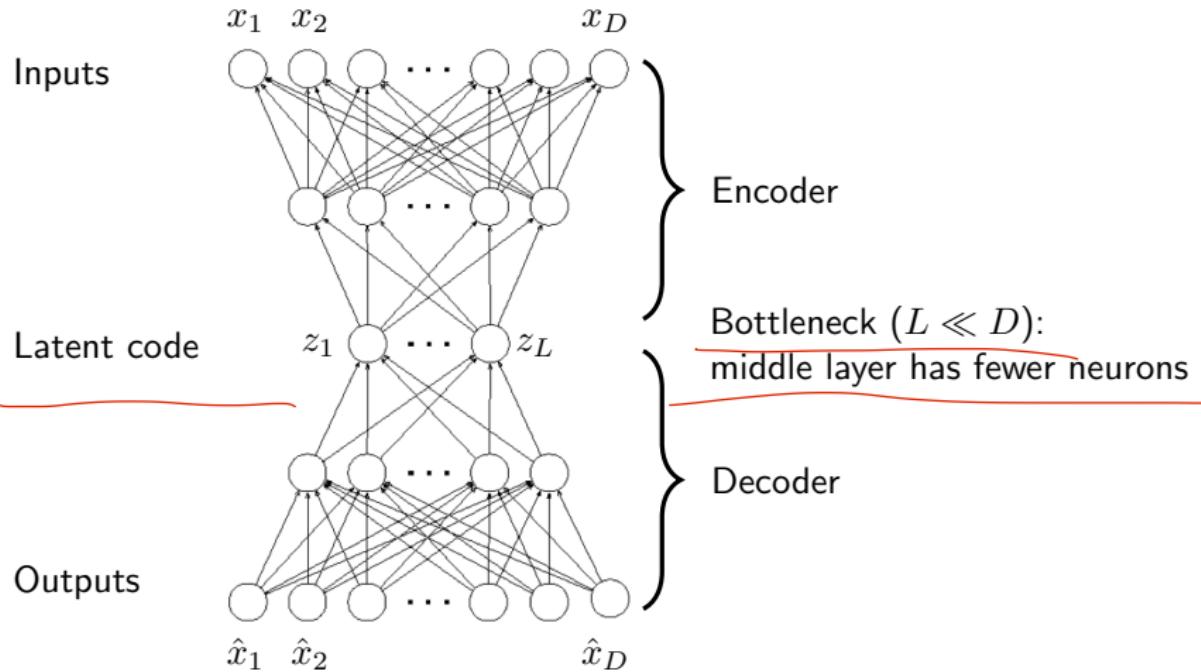
$$\min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N |f(\mathbf{x}_i, \mathbf{W}) - \mathbf{x}_i|^2$$

- Alternative view: find a **latent representation** $\mathbf{z} \in \mathbb{R}^L$ which is a compact code for $\mathbf{x} \in \mathbb{R}^D$ since $L \ll D$
 - $f_{enc}(\mathbf{x}) = \mathbf{z}$ # encoder: project the data to a lower dimension
 - $f_{dec}(\mathbf{z}) \approx \mathbf{x}$ # decoder: reconstruct the data from the latent code

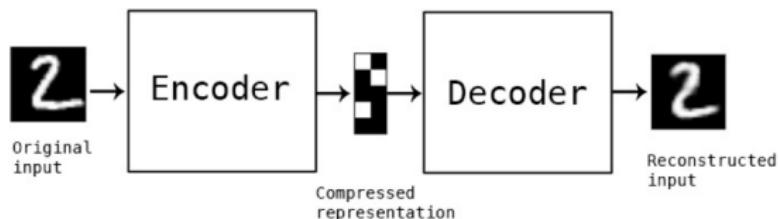
$$f_{dec}(f_{enc}(\mathbf{x})) \approx \mathbf{x}$$

Autoencoders

- f_{enc} and f_{dec} are non-linear functions implemented by a neural network



Autoencoders: Example



Original image:



Reconstruction:



- An autoencoder whose code dimension is less than the input dimension ($L \ll D$) is called **undercomplete**
 - Forces the autoencoder to capture the most salient features of the data
 - $L \geq D$ (**overcomplete**) the autoencoder can simply learn the identity function
- Training autoencoders in practice:
 - Add a regularization term to the SSE loss to prevent overfitting
 - Weight tying: f_{enc} and f_{dec} share the same weights
- Other extensions:
 - Denoising autoencoders (DAEs): receive a corrupted (noisy) training data point as input and predict the "clean" uncorrupted data point as output
 - Variational autoencoders (covered in our MLGS lecture)

(Linear) Autoencoders & PCA: Comparison

- What if f_{enc} and f_{dec} are linear functions?

- $f_{enc}(\mathbf{x}, \mathbf{W}_1) = \mathbf{x}\mathbf{W}_1$

$$\#\mathbf{W}_1 \in \mathbb{R}^{D \times L}$$

- $f_{dec}(\mathbf{z}, \mathbf{W}_2) = \mathbf{z}\mathbf{W}_2$

$$\#\mathbf{W}_2 \in \mathbb{R}^{L \times D}$$

- We have: $f_{dec}(f_{enc}(\mathbf{x})) = \mathbf{x}\mathbf{W}_1\mathbf{W}_2$ and

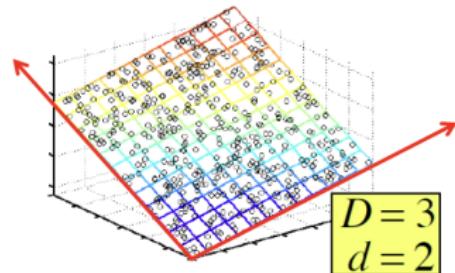
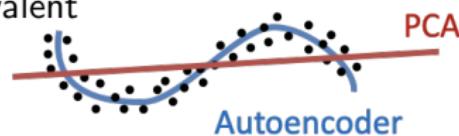
$$\min_{\mathbf{W}_1 \mathbf{W}_2} \frac{1}{N} \sum_{i=1}^N |f(\mathbf{x}_i, \mathbf{W}) - \mathbf{x}_i|^2 = \min_{\mathbf{W}_1 \mathbf{W}_2} \frac{1}{N} \sum_{i=1}^N |\mathbf{x}_i \mathbf{W}_1 \mathbf{W}_2 - \mathbf{x}_i|^2 =$$
$$\underbrace{\min_{\mathbf{W}} \frac{1}{N} \sum_{i=1}^N |\mathbf{x}_i \mathbf{W} - \mathbf{x}_i|^2}_{\text{Equivalent: } \mathbf{W}_1 \mathbf{W}_2 = \mathbf{W} \text{ s.t rank of } \mathbf{W} \text{ is } L}$$

- Optimal solution (assuming normalization):

- PCA: $\mathbf{W}^* = \boldsymbol{\Gamma}$ where $\boldsymbol{\Gamma}$ are the (top L) eigenvectors of $\mathbf{X}^T \mathbf{X}$

Summary

- Dimensionality Reduction and Matrix Factorization are highly related
 - And can be used for various purposes
- PCA, SVD (and linear Autoencoders) are equivalent
 - Optimal low rank approximation
- Matrix factorization for rating prediction
 - Very general formulation: allows to handle, e.g., missing entries
- Autoencoders and t-SNE perform non linear dimensionality reduction
- Why are these techniques useful?
 - Less storage required
 - More efficient processing of the data
 - Remove redundant and noisy features
 - Discover hidden correlations/topics/concepts
 - Interpretation and visualization



Reading material

- Bishop, chapters: 12.1 , 12.2.1
- Leskovec, Rajaraman, Ullman- Mining of Massive Datasets: chapter 11
- Goodfellow- Deep Learning: chapter 14
- t-SNE
- Understanding UMAP