

Roadmap

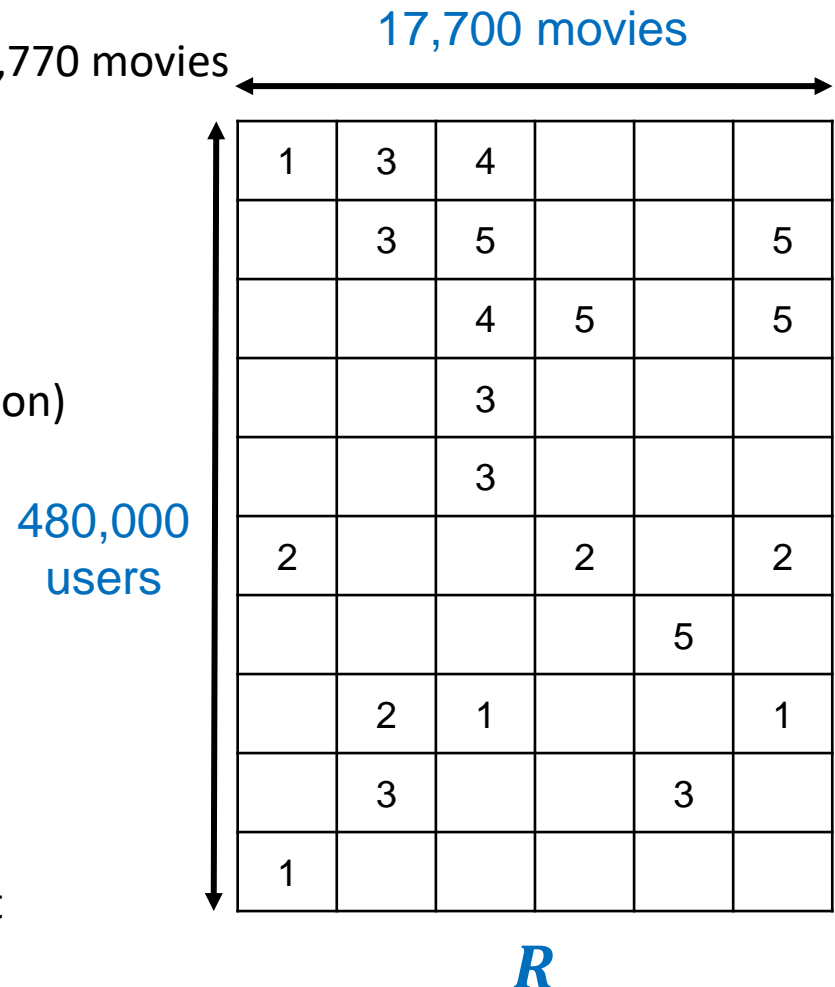
- Chapter: Dimensionality Reduction & Matrix Factorization
 1. Introduction
 2. Principal Component Analysis (PCA)
 3. Singular Value Decomposition (SVD)
 4. **Matrix Factorization**
 - **Motivation & Approach**
 - Regularization & Sparsity
 - Further Factorization Models
 5. Neighbor Graph Methods
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

- Competition
 - 2,700+ teams
 - \$1 million prize for 10% improvement



Evaluating Recommender Systems

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

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

true rating of
user u for item i

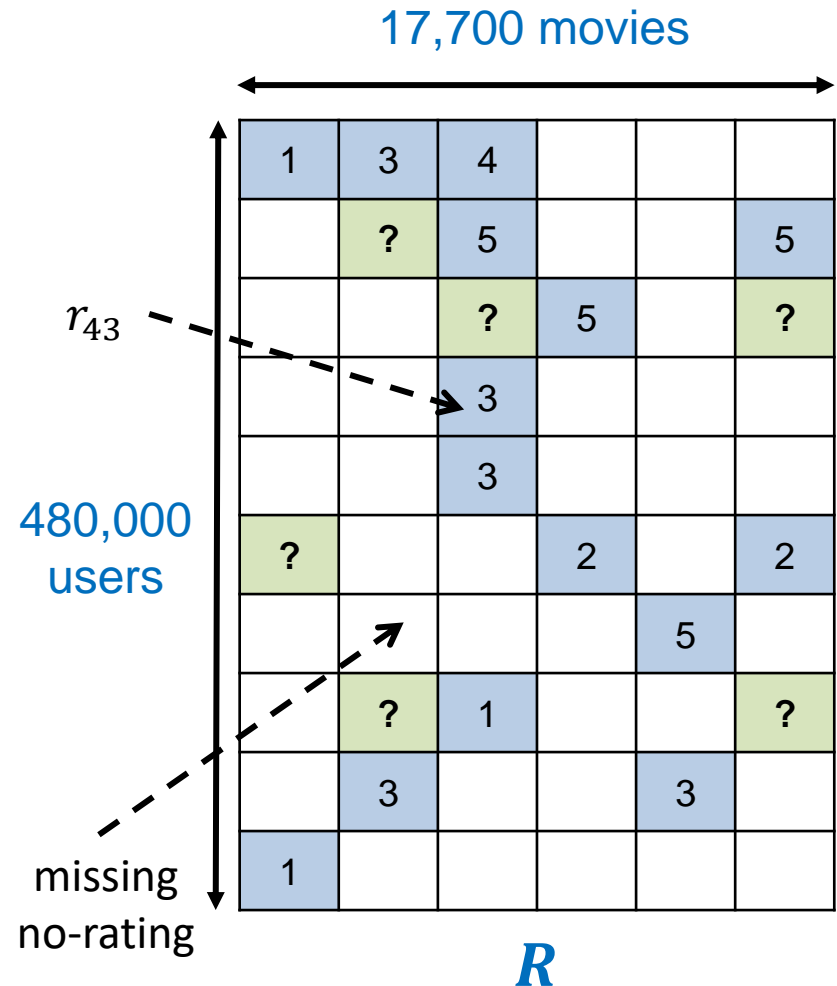
r_{ui}

predicted
rating

\hat{r}_{ui}


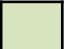

- Legend:

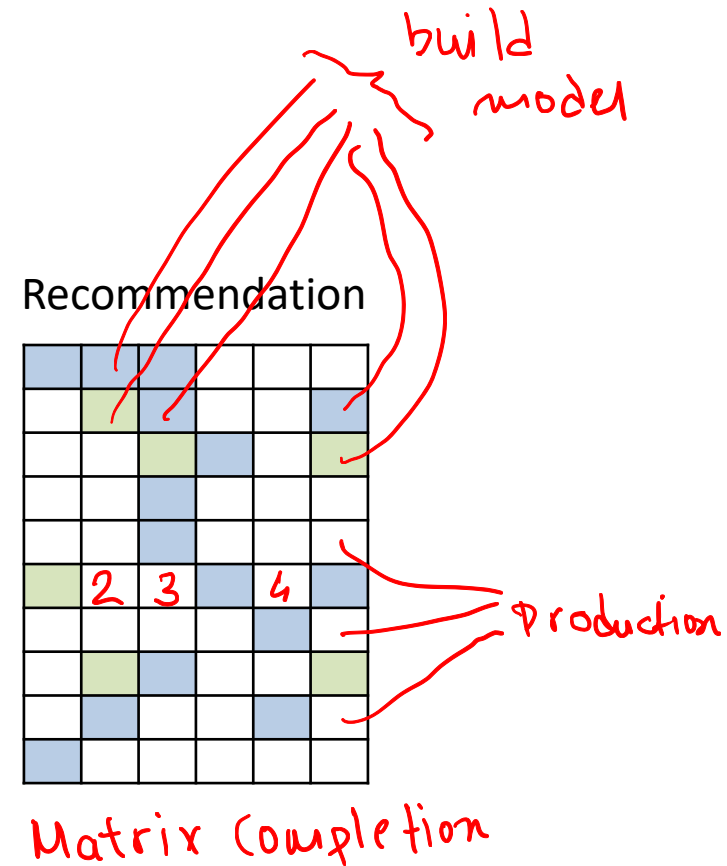
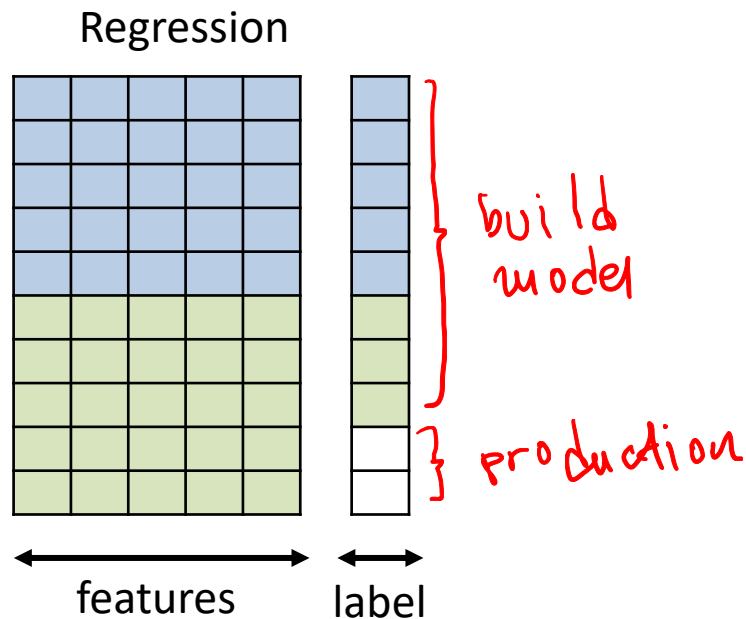
- Training and validation data
- Test data
- Missing data



Regression vs. Recommendation

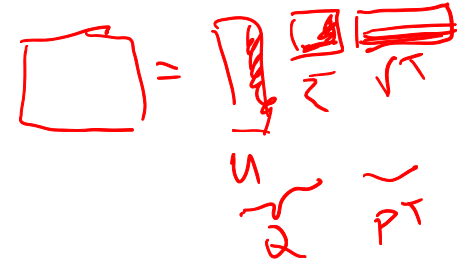
- Legend:

-  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: $R = U \Sigma V^T \rightarrow Q = U \Sigma, P = V$



items

1	0	3	0	0	5	0	0	5	0	4	0
0	0	5	4	0	0	4	0	0	2	1	3
2	4	0	1	2	0	3	0	4	3	5	0
0	2	4	0	5	0	0	4	0	0	2	0
0	0	4	3	4	2	0	0	0	0	2	5
1	0	3	0	3	0	0	2	0	0	4	0

R

users

users

factors

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

Q

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

factors

factors

u

2

~

P^T

SVD on Rating Data

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

$$\min_{U, \Sigma, V} \sum_{\substack{u=1 \dots n \\ i=1 \dots d}} \left(\underline{R_{ui}} - \underbrace{[U \Sigma V^T]_{ui}}_B \right)^2$$

$\min_B \|R - B\|_F^2$
 $\text{rank}(B) = k$

- SSE and RMSE are monotonically related:

- $\text{RMSE} = \frac{1}{\text{const.}} \sqrt{\text{SSE}}$

- Great news: SVD is minimizing RMSE

$UU^T = I$
 $VV^T = I$

- **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?

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:

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

\mathbf{q}_u $1 \times k$

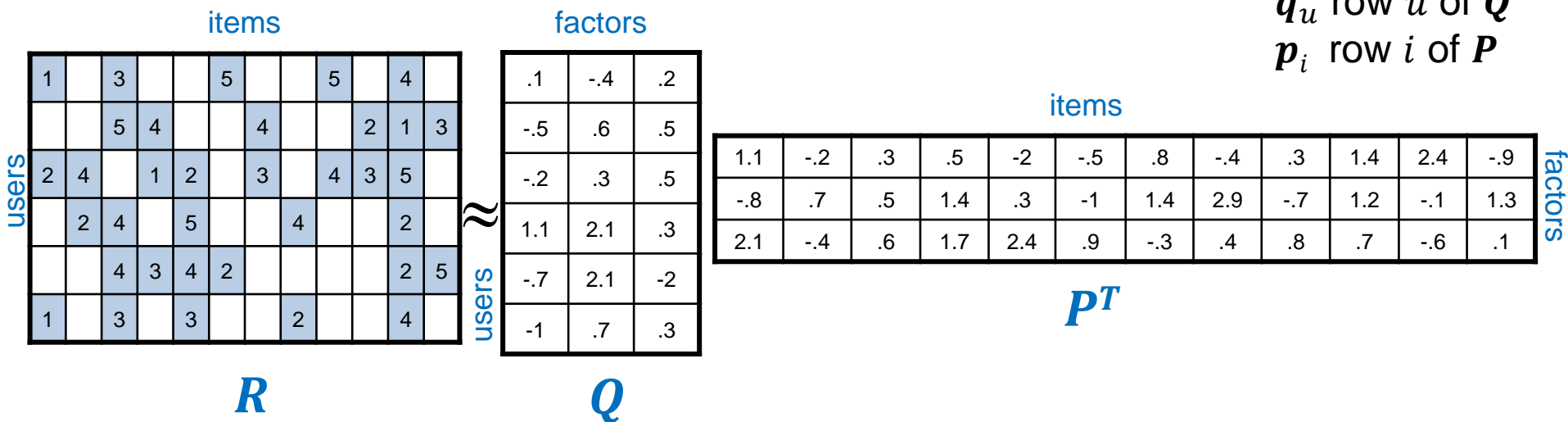
\mathbf{p}_i $1 \times k$

$\mathbf{q}_u \cdot \mathbf{p}_i^T$ $1 \times k (k \times k)^T = 1 \times 1$

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

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

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



Finding the Latent Factors (I)

- Goal: $\min_{\mathbf{P}, \mathbf{Q}} \sum_{(u, i) \in \mathcal{S}} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2 =: \min_{\mathbf{P}, \mathbf{Q}} f(\mathbf{P}, \mathbf{Q})$
- One approach: **Alternating Optimization** // a.k.a. block coordinate minimization
 - Pick initial values for \mathbf{P} and \mathbf{Q}
 - Alternatingly 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. $\mathbf{P}^{(t+1)} = \operatorname{argmin}_{\mathbf{P}} f(\mathbf{P}, \mathbf{Q}^{(t)})$
 3. $\mathbf{Q}^{(t+1)} = \operatorname{argmin}_{\mathbf{Q}} f(\mathbf{P}^{(t+1)}, \mathbf{Q})$
 4. $t = t + 1$
 5. goto 2 until convergence



Finding the Latent Factors (II)

- Initialization of \mathbf{P} and \mathbf{Q} :

- Use, e.g., SVD where missing entries are replaced by 0 (or overall mean)

- How to solve $\mathbf{P}^{(t+1)} = \underset{\mathbf{P}}{\operatorname{argmin}} f(\mathbf{P}, \mathbf{Q}^{(t)}) = \underset{\mathbf{P}}{\operatorname{argmin}} \sum_{(u,i) \in S} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2$

- Observation 1: Since \mathbf{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 $S_{*,i} = \{u \mid (u, i) \in S\}$ // = only users u who have rated item i

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

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**

$$X = \begin{bmatrix} \text{---} \\ \text{---} \mathbf{x}_j^T \text{---} \\ \text{---} \end{bmatrix}$$

$$\mathbf{y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n]$$

- $\min_{\mathbf{w}} \sum_{j=1}^n (y_j - \mathbf{w}^T \mathbf{x}_j)^2$ // y_j is scalar, \mathbf{x}_j and \mathbf{w} (column) vectors

- Optimal solution in closed form: $\mathbf{w} = \left(\frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^T \right)^{-1} \cdot \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j y_j = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

$$D \times D = (D \times 1) (1 \times D)$$

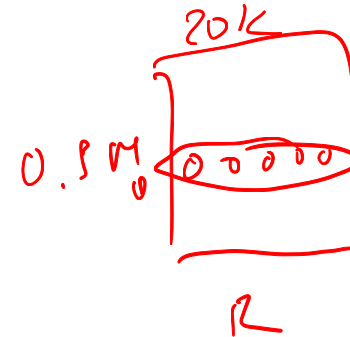
- In our case: $\mathbf{p}_i^T = \left(\frac{1}{|S_{*,i}|} \sum_{u \in S_{*,i}} \mathbf{q}_u^T \mathbf{q}_u \right)^{-1} \cdot \frac{1}{|S_{*,i}|} \sum_{u \in S_{*,i}} \mathbf{q}_u^T r_{ui}$

- Equivalently for \mathbf{q}_u


- Computation of $\operatorname{argmin}_{\mathbf{P}} (\mathbf{P}, \mathbf{Q}^{(t)})$ 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 $\mathbf{q}_u \leftarrow \mathbf{q}_u - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{q}_u}$ $\mathbf{p}_i \leftarrow \mathbf{p}_i - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{p}_i}$

- 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		3			5			5		4	
		5	4	?		4			2	1	3
2	4		1	2		3		4	3	5	
	2	4		5			4			2	
		4	3	4	2					2	5
1		3		3			2			4	

users

R

\approx

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

\mathbf{q}_u row u of Q

\mathbf{p}_i row i of P

factors

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

users

Q

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

factors

P^T

Rating Prediction

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

items

users

R

1		3			5			5		4	
		5	4	?		4			2	1	3
2	4		1	2		3		4	3	5	
	2	4		5			4			2	
		4	3	4	2					2	5
1		3		3			2			4	

\approx

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

\mathbf{q}_u row u of Q

\mathbf{p}_i row i of P

factors

users

Q

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

items

factors

P^T

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

Rating Prediction

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

items

users

R

1		3			5			5		4	
		5	4	2.4		4			2	1	3
2	4		1	2		3		4	3	5	
	2	4		5			4			2	
		4	3	4	2					2	5
1		3		3			2			4	

\approx

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

\mathbf{q}_u row u of Q

\mathbf{p}_i row i of P

factors

users

Q

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

factors

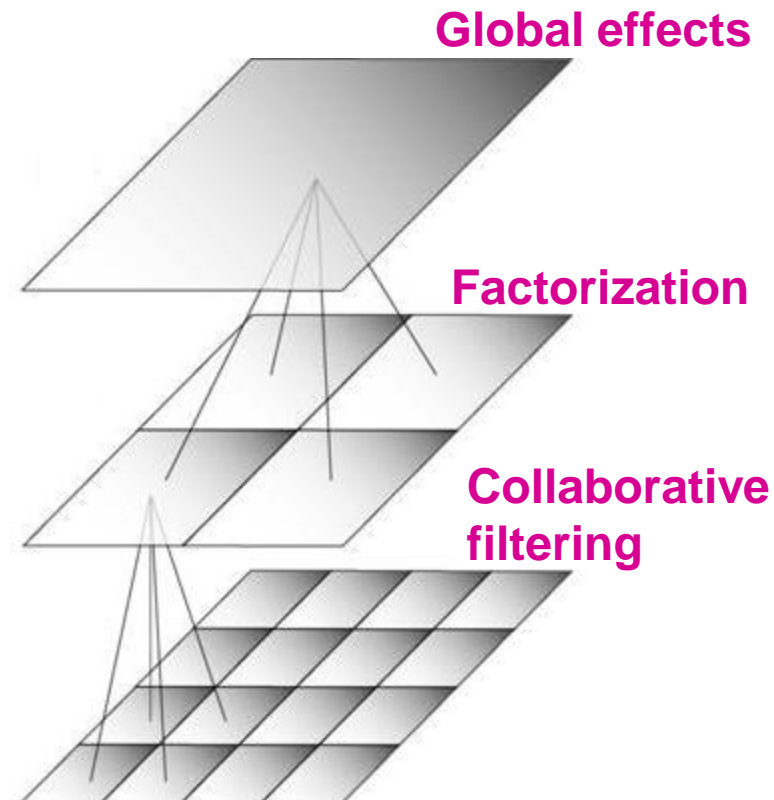
items

P^T

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

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~~^{item} 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

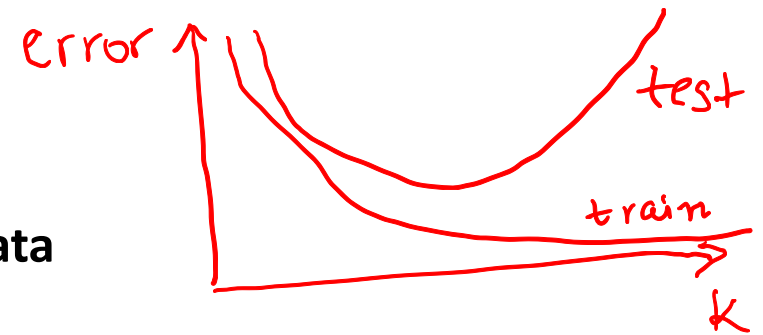
$$\min_{\substack{\mathbf{P}, \mathbf{Q} \\ b_u, b_i, b}} \sum_{(u, i) \in S} (r_{ui} - (\mathbf{q}_u \cdot \mathbf{p}_i^T + b_u + b_i + b))^2$$

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

Roadmap

- 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 (Non-linear 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

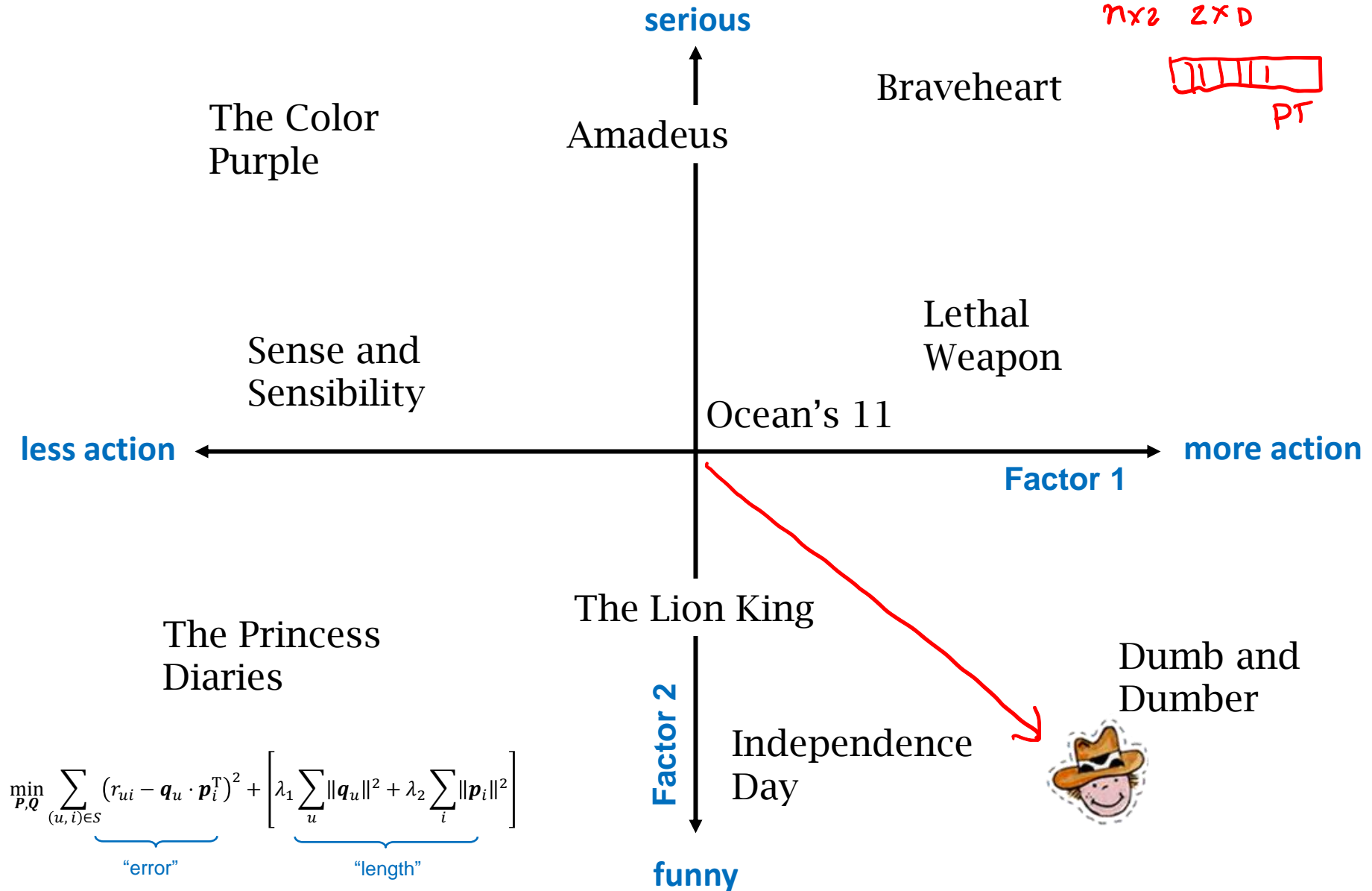
Regularization

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

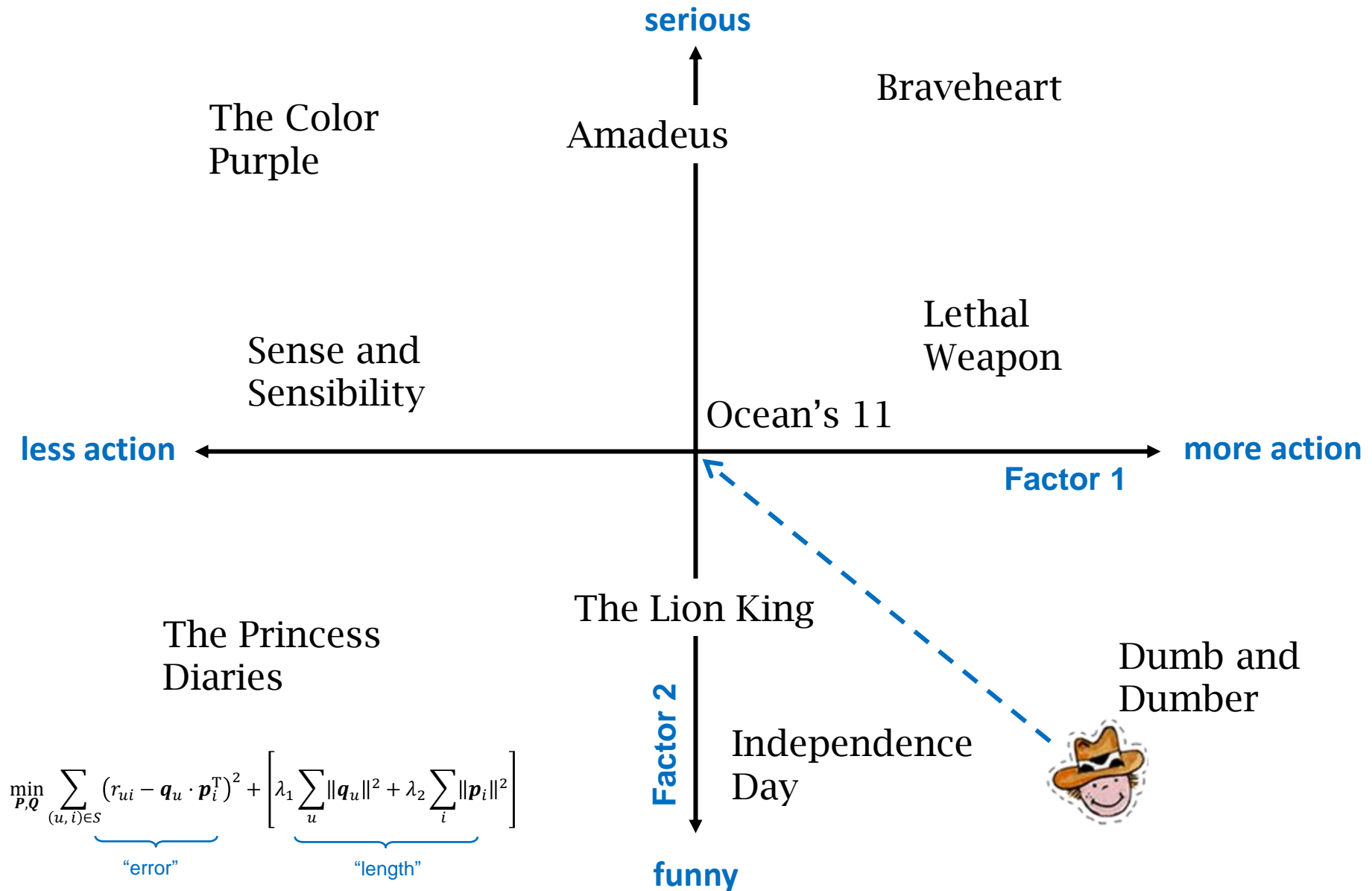
$$\min_{\mathbf{P}, \mathbf{Q}} \underbrace{\sum_{(u,i) \in S} (r_{ui} - \mathbf{q}_u \cdot \mathbf{p}_i^T)^2}_{\text{“error”}} + \underbrace{\left[\lambda_1 \sum_u \|\mathbf{q}_u\|^2 + \lambda_2 \sum_i \|\mathbf{p}_i\|^2 \right]}_{\text{“length”}}$$

- 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 \mathbf{P}, \mathbf{Q} that achieve the minimum of the objective

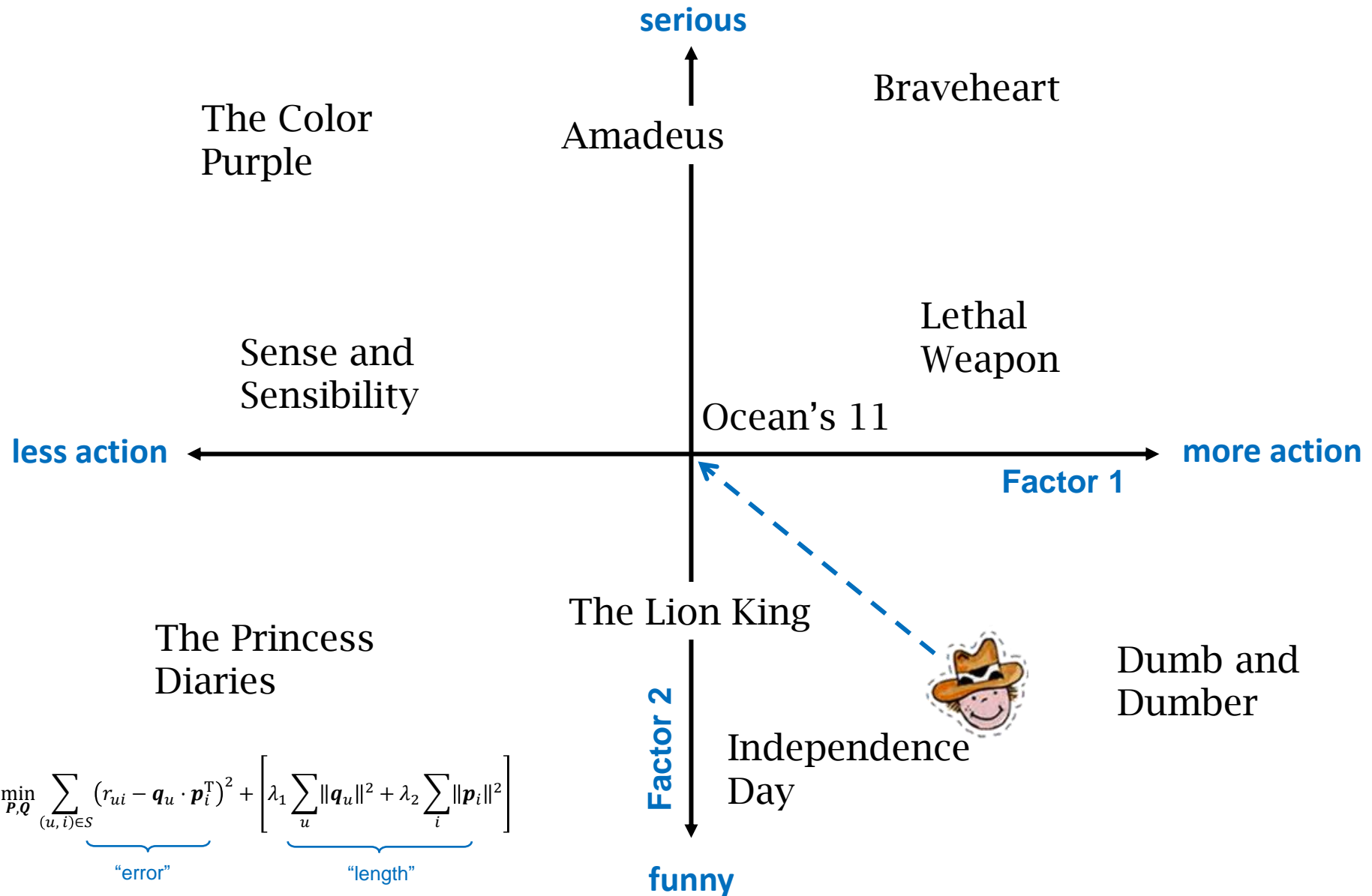
The Effect of Regularization



The Effect of Regularization

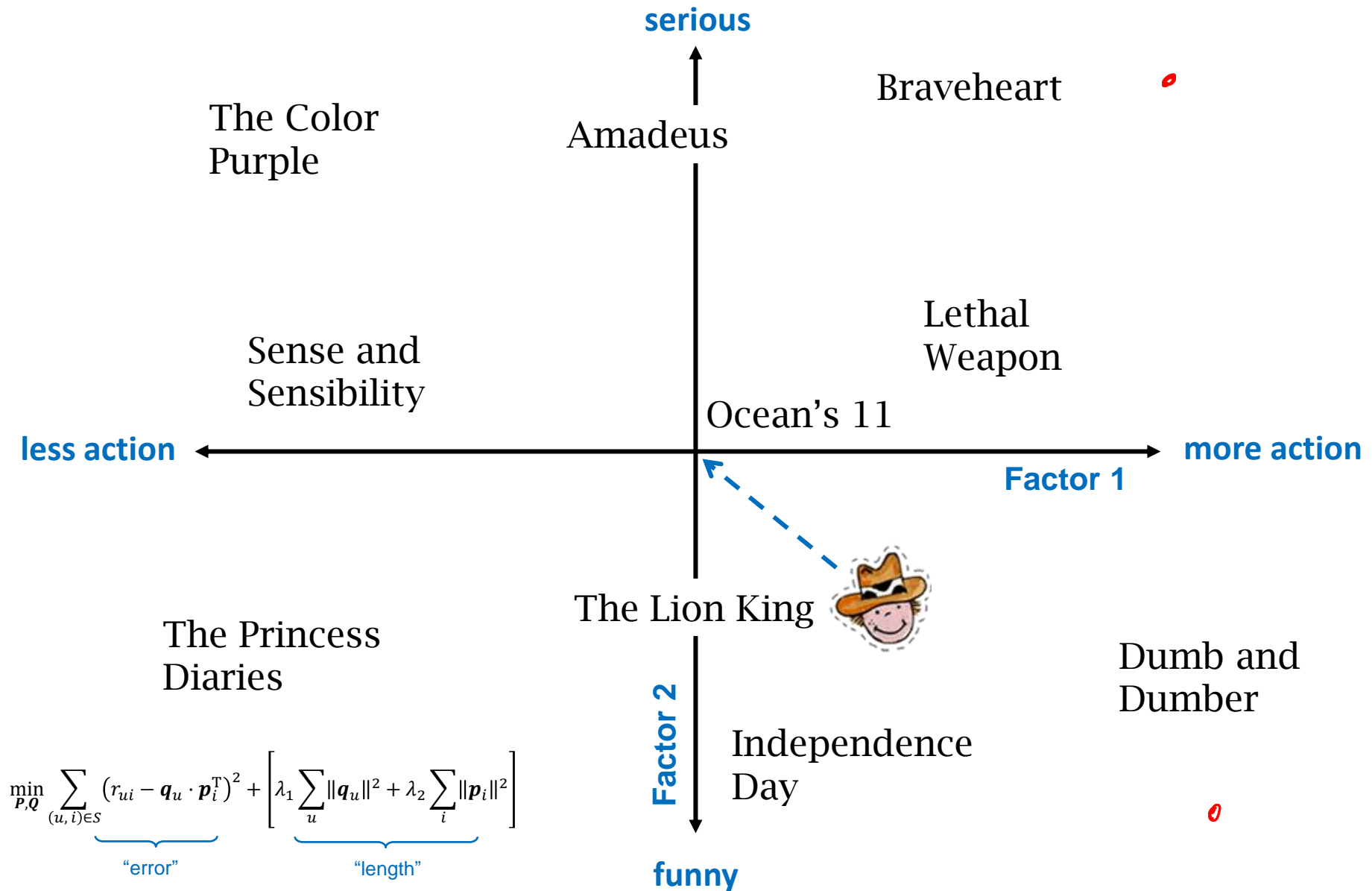


The Effect of Regularization



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

The Effect of Regularization



Regularization in Our Use Case

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

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

SGD for MF with Regularization and Biases

- SGD on the objective

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

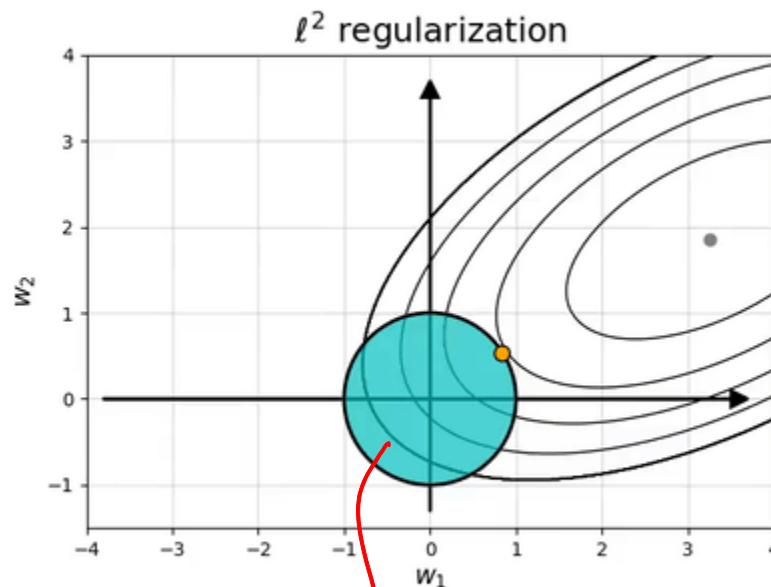
- Pick a random user u and a random item i with rating r_{ui} (batch size 1)
- The gradient update step is very similar to before
 - $e_{ui} \leftarrow r_{ui} - (\mathbf{q}_u \cdot \mathbf{p}_i^T + b_u + b_i + b)$ helper variable, the current error
 - $\mathbf{q}_u \leftarrow \mathbf{q}_u + 2\eta(e_{ui}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

L2 vs. L1 Regularization

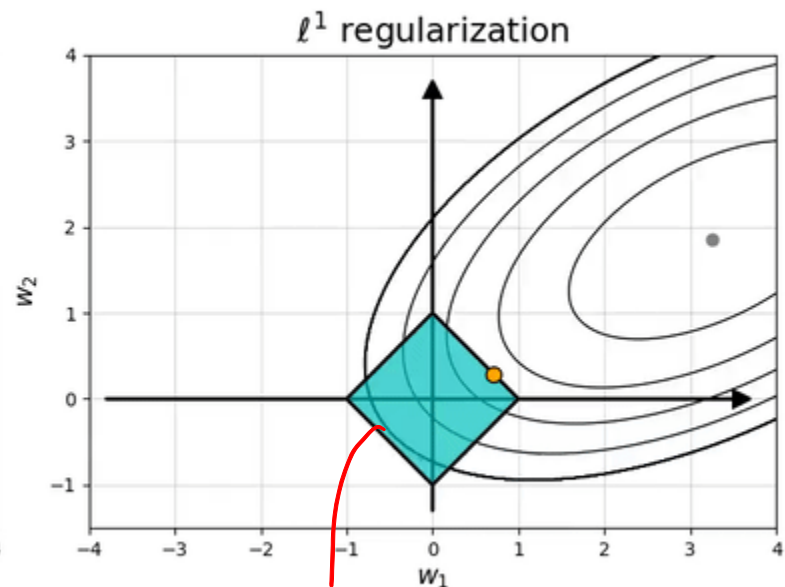
- Comparison: L2 vs. L1 regularization
 - L2 tries to "shrink" the parameter vector \mathbf{w} **equally**
 - Large values are highly penalized due to the square in the L2 norm
 - Unlikely that any component will be **exactly 0**
 - L1 allows large values in individual components of \mathbf{w} by shrinking 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 L1 constraint promotes sparsity



$$\|w\|_2 \leq 1$$



$$\|w\|_1 \leq 1$$

by @itayevron

$$\min_w f(w) \quad \text{subject to} \quad \|w\| \leq 1 \quad \rightarrow \quad \min_w f(w) + \lambda (\|w\| - 1)$$

Roadmap

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

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

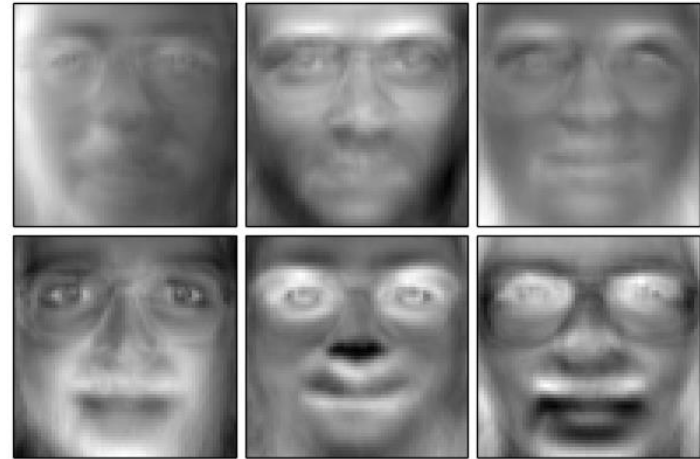
- Constrained optimization

NNMF Example: Olivetti Faces Data

$$X = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$



Eigenfaces - RandomizedPCA - Train time 1.4s



Non-negative components - NMF - Train time 2.9s



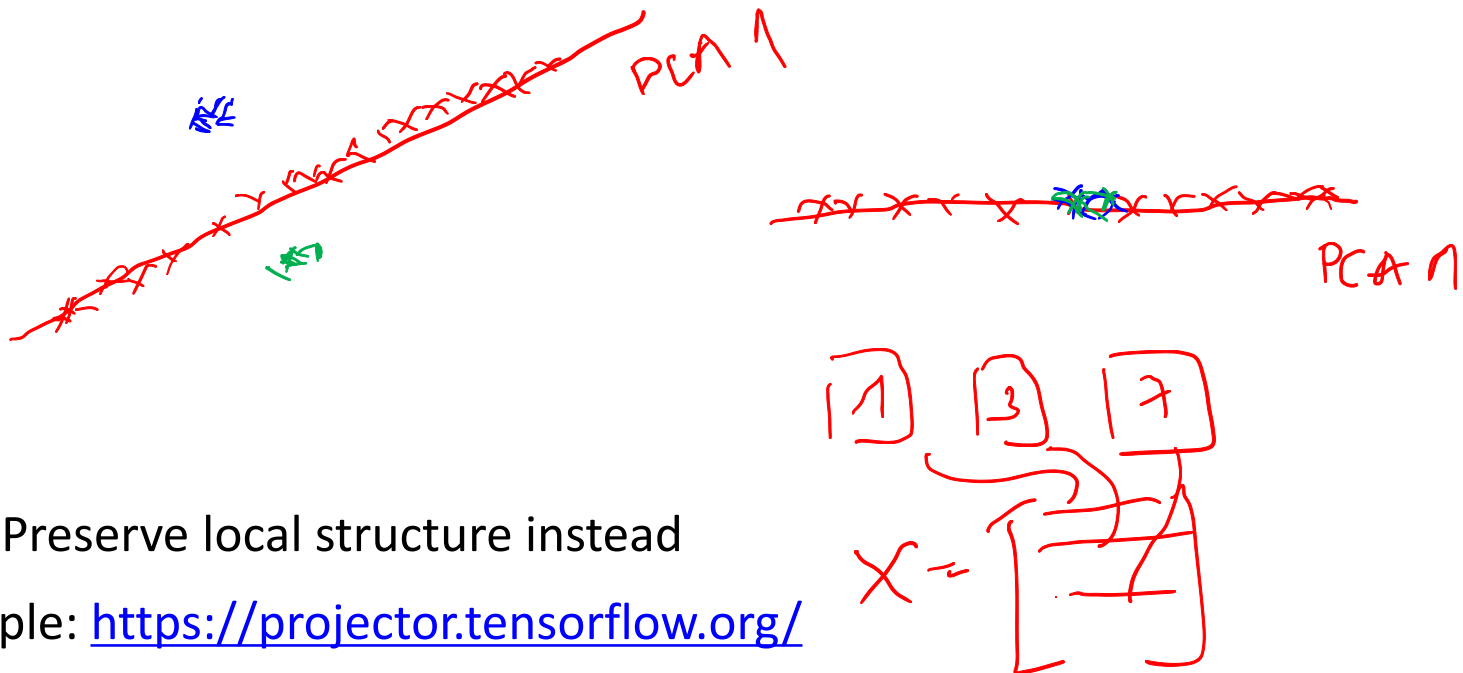
results according to scikit-learn

Roadmap

- 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 (Non-linear 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):

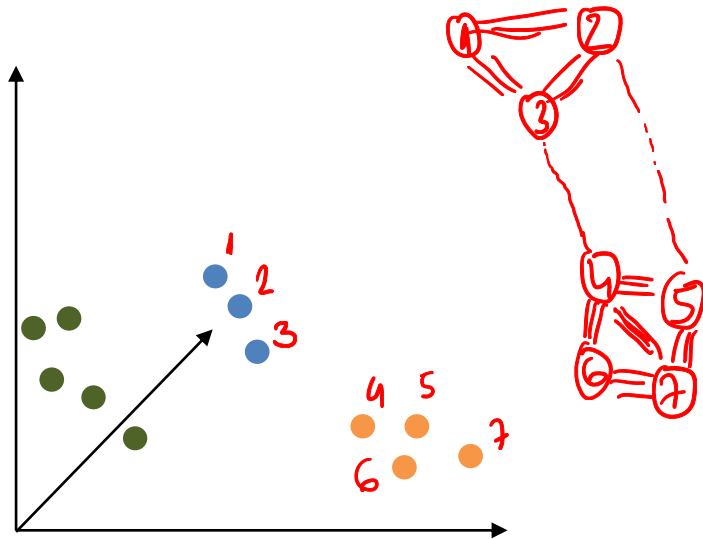


- Idea: Preserve local structure instead
- Example: <https://projector.tensorflow.org/>

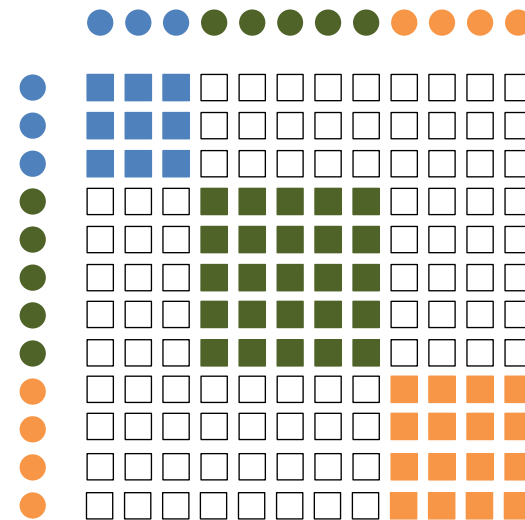
Neighbor Graph Methods

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

Highdimensional data:





Neighbor / similarity graph:

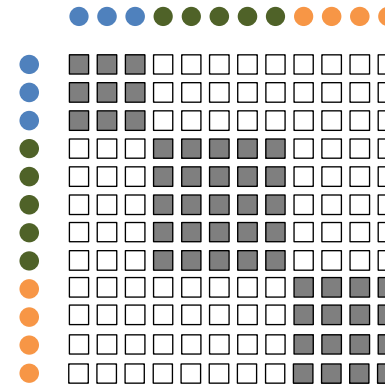
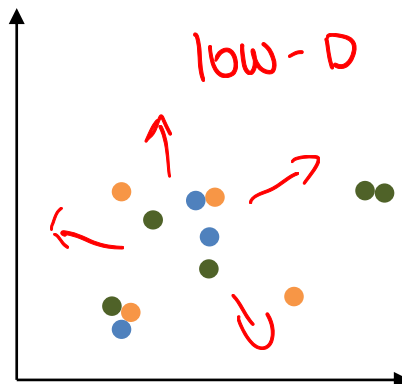
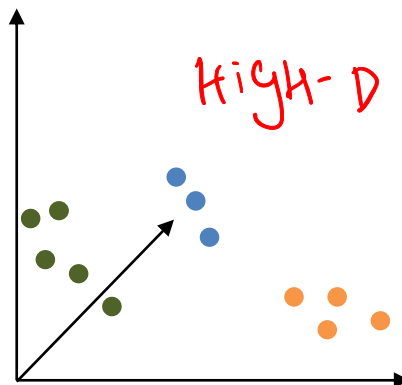


- Idea: Low dimensional neighborhood similar to the original neighborhood

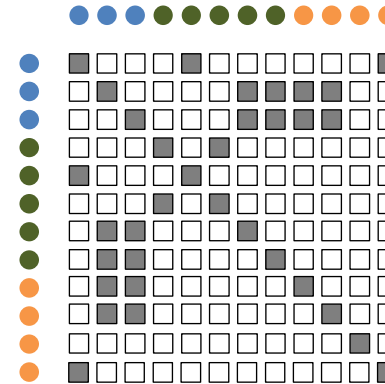
Neighbor Graph Methods

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

 - high similarity
 - low similarity



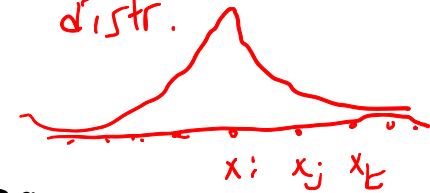
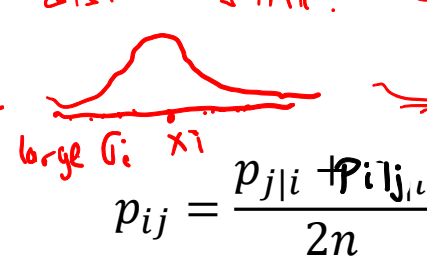
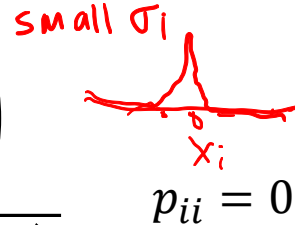
adjacency matrix
of the graph /
pairwise similarity
matrix /
pairwise distance
matrix



t-SNE: t-Distributed Stochastic Neighbor Embedding

- High-dimensional similarities for input \mathbf{x}_i :

$$p_{j|i} = \frac{\exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_k\|^2}{2\sigma_i^2}\right)}$$



- σ_i chosen to achieve fixed perplexity (effective number of neighbors)

$x_i = [0.1 \ 2 \ 3 \ 4 \ 5]$
 $x_j = [\dots \dots \dots]$

- Low-dimensional similarities for (to be optimized) parameters \mathbf{y}_i :

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$

$q_{ii} = 0$

Variables to be optimized

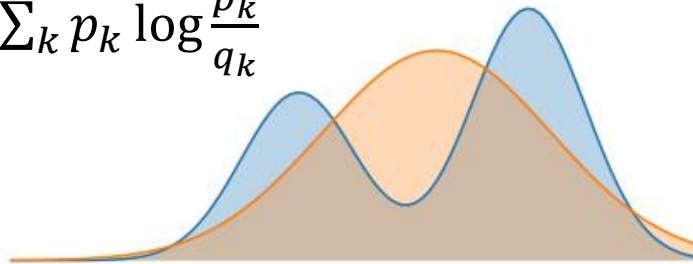
- Minimize the KL divergence: $\min_{\mathbf{y}_i} KL(P||Q) = \sum_i \sum_{j \neq i} p_{ij} \log \frac{p_{ij}}{q_{ij}}$

Note on the KL Divergence

$$P(\theta) = \mathcal{P} = [p_1, p_2, \dots, p_k]$$
$$Q = [q_1, q_2, \dots, q_k]$$

- $KL(P||Q) \geq 0$ 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

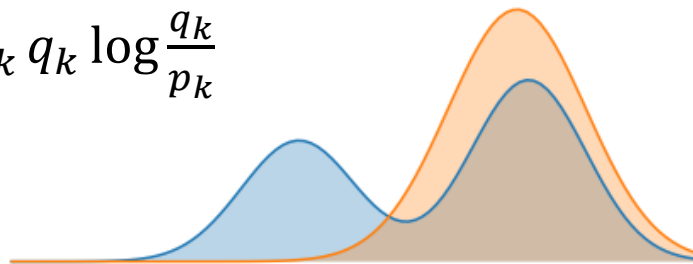
$$KL(P||Q) = \sum_k p_k \log \frac{p_k}{q_k}$$



Forward KL is Mean-Seeking

$$KL(P||Q)$$

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

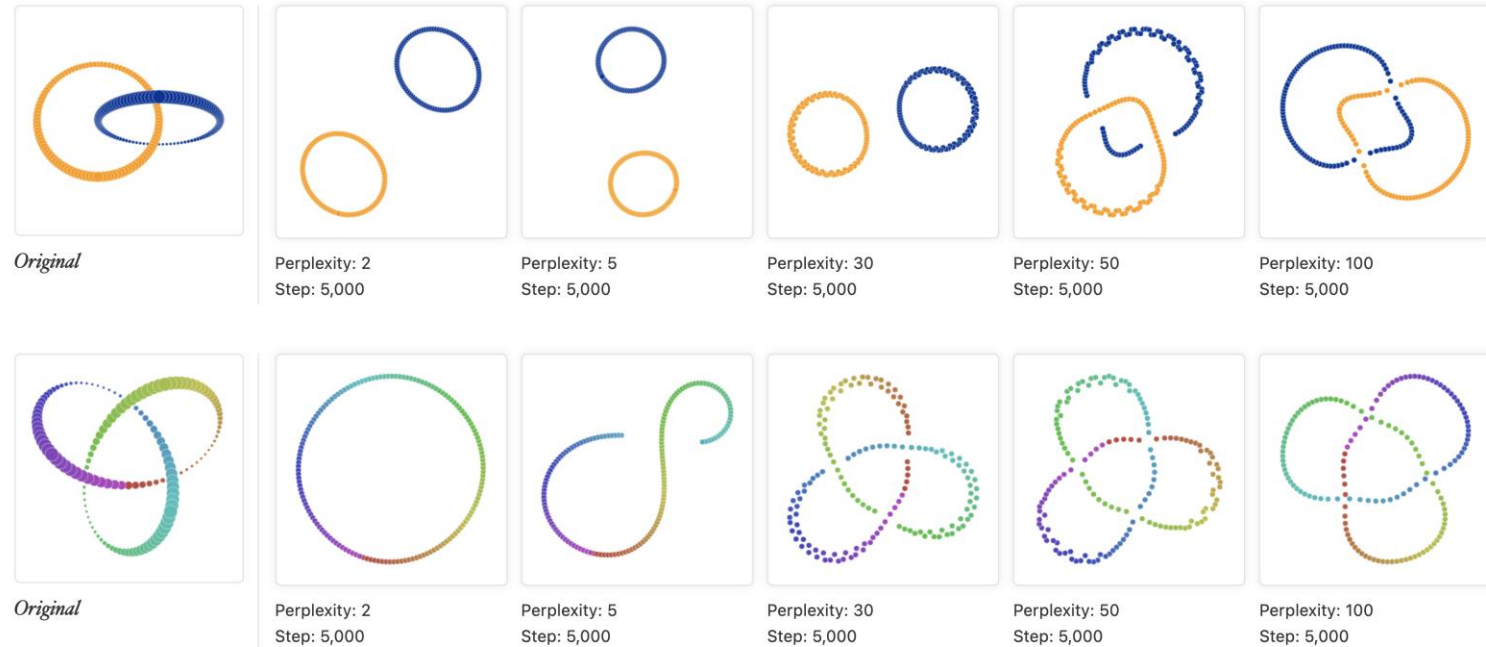


Reverse KL is Mode-Seeking

$$KL(Q||P)$$

Figures adapted from Dibya Ghosh.

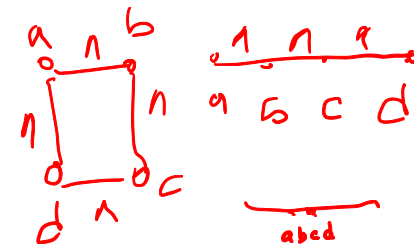
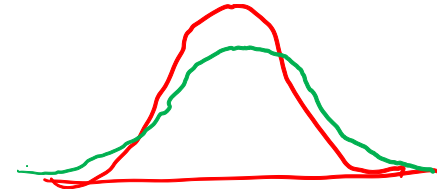
Interactive t-SNE



[Source of illustration, for more details and interactive widgets.](#)

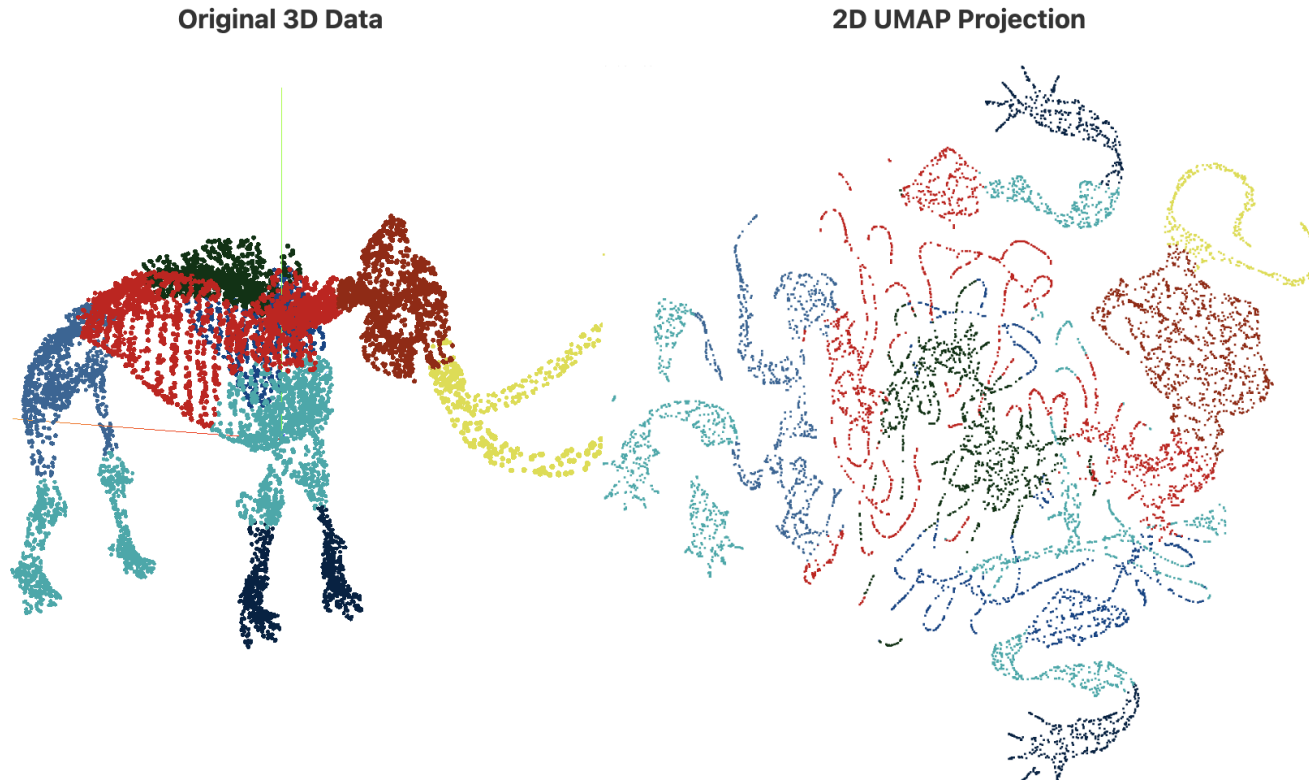
The Advantages and Disadvantages of t-SNE

- 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, ...



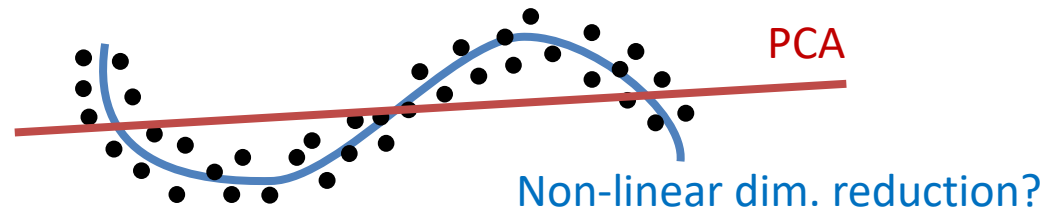
[Link for details, widgets, a comparison to T-SNE, ...](#)

Roadmap

- 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 (Non-linear Dimensionality Reduction)**

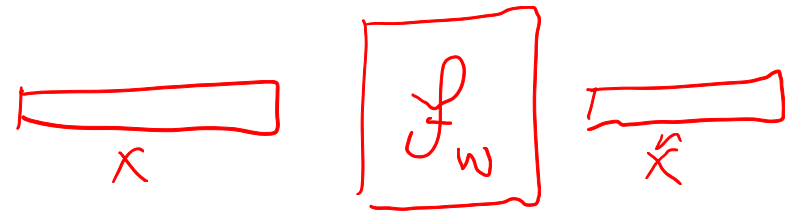
Motivation

- PCA / SVD can only capture linear structure (linear variations in the data)
 - Recall: transformed data is $\mathbf{Y} = \tilde{\mathbf{X}} \cdot \mathbf{\Gamma}$
 - Linear projection by the eigenvectors $\mathbf{\Gamma}$
- 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

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

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

$$\min_W \frac{1}{N} \sum_{i=1}^N \|f(x_i, W) - x_i\|^2$$

\uparrow *noisy* \uparrow *clean*

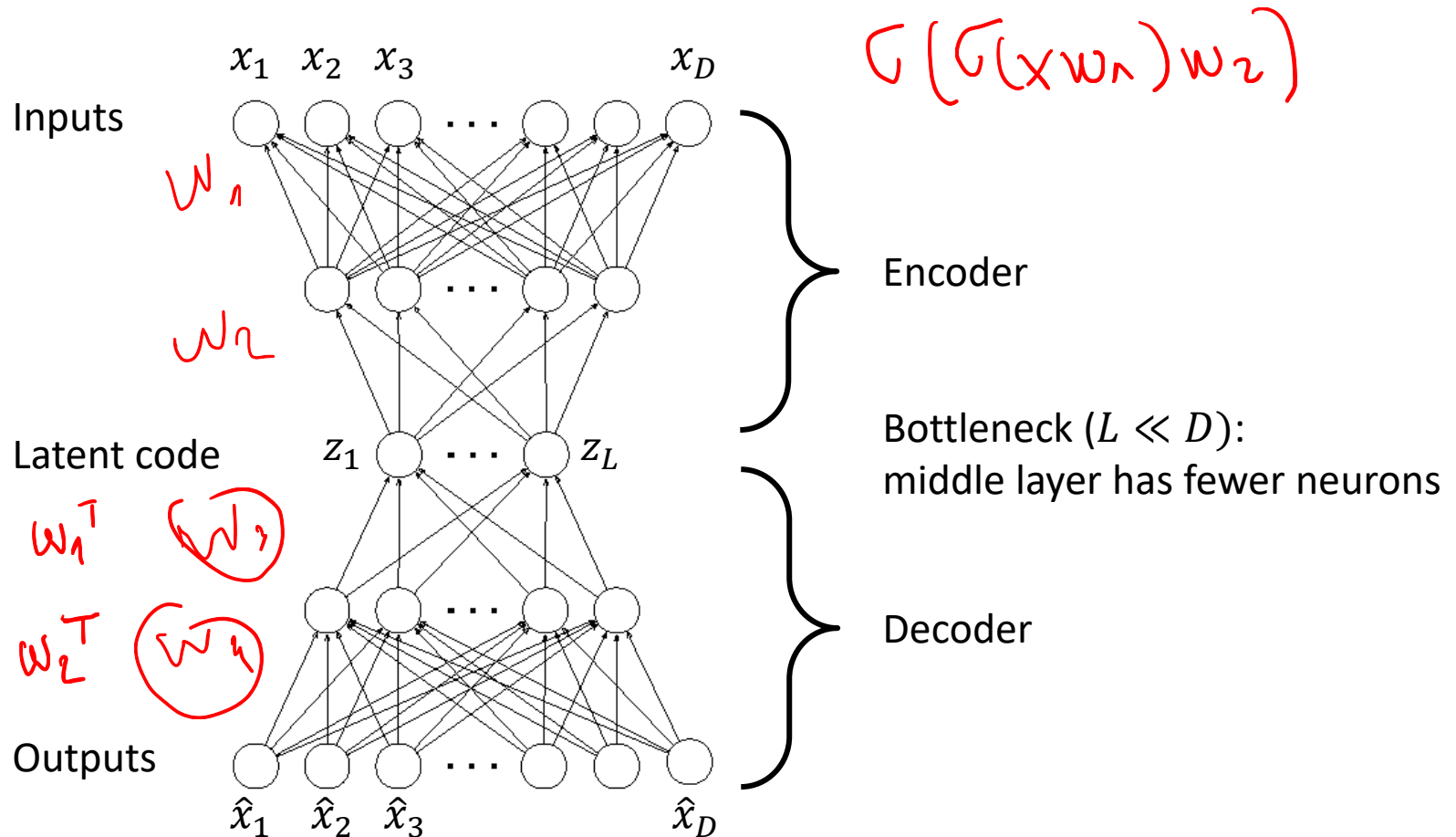
- Alternative view: find a **latent representation** $z \in \mathbb{R}^L$ which is a compact code for $x \in \mathbb{R}^D$ since $L \ll D$

- $f_{enc}(x) = z$ # encoder: project the data to a lower dimension
- $f_{dec}(z) \approx x$ # decoder: reconstruct the data from the latent code

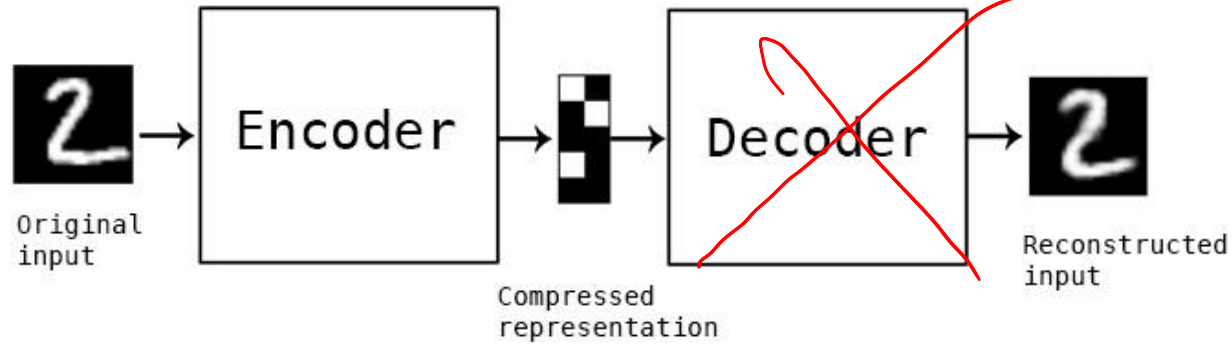
$$f_{dec}(f_{enc}(x)) \approx x$$

Autoencoders

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



Autoencoders: Example



Original image:



Reconstruction:



Autoencoders

- 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?

- $f_{enc}(\mathbf{x}, \mathbf{W}_1) = \mathbf{x}\mathbf{W}_1$ $\# \mathbf{W}_1 \in \mathbb{R}^{D \times L}$

- $f_{dec}(\mathbf{z}, \mathbf{W}_2) = \mathbf{z}\mathbf{W}_2$ $\# \mathbf{W}_2 \in \mathbb{R}^{L \times D}$

$$\mathbf{W} := \mathbf{W}_1 \mathbf{W}_2$$

$D \times D \quad D \times L \quad L \times D$

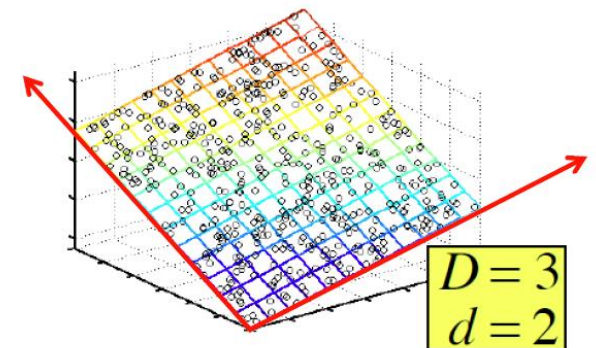
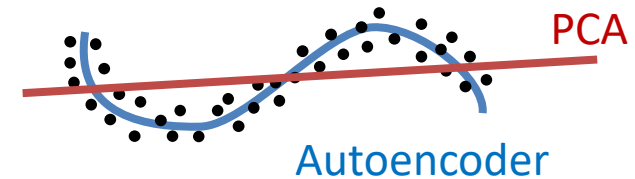
- We have: $f_{dec}(f_{enc}(\mathbf{x})) = \mathbf{x}\mathbf{W}_1\mathbf{W}_2$ and

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

- Optimal solution (assuming normalization):
 - PCA: $\mathbf{W}^* = \mathbf{\Gamma}$ where $\mathbf{\Gamma}$ 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

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](#)