# 기계학습 (2022년도 2학기)

## Matrix Factorizations

### 전북대학교 컴퓨터공학부

# Overview

- Recall PCA: project data onto a low-dimensional subspace defined by the top eigenvalues of the data covariance

- We saw that PCA could be viewed as a linear autoencoder, which let us generalize to nonlinear autoencoders

- Today we consider another generalization, matrix factorizations

- view PCA as a matrix factorization problem

- extend to matrix completion, where the data matrix is only partially observed

- extend to other matrix factorization models, which place different kinds of structure on the factors

# PCA as Matrix Factorization

- Recall: each input vector $\mathbf{x}^{(i)}$ is approximated as $\mathbf{U}\mathbf{z}^{(i)}$, where $\mathbf{U}$ is the orthogonal basis for the principal subspace, and $\mathbf{z}^{(i)}$ is the code vector.

- Write this in matrix form: $\mathbf{X}$ and $\mathbf{Z}$ are matrices with one column per data point

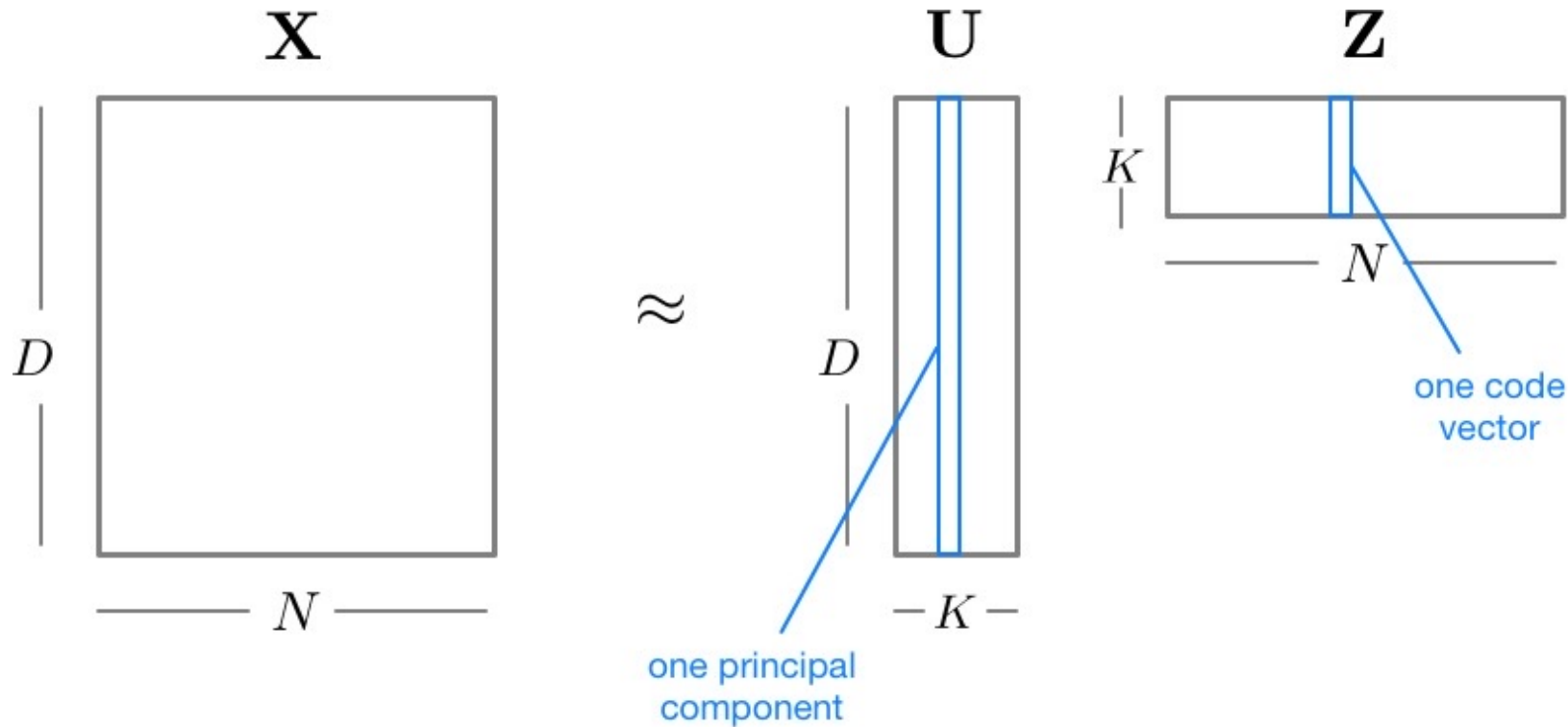- We transpose our usual convention for data matrices (for some parts of this lecture).

- Writing the squared error in matrix form

$$\sum_{i=1}^{N} \|\mathbf{x}^{(i)} - \mathbf{U}\mathbf{z}^{(i)}\|^2 = \|\mathbf{X} - \mathbf{U}\mathbf{Z}\|_F^2$$

- Recall that the Frobenius norm is defined as $\|\mathbf{A}\|_F^2 = \sum_{i,j} a_{ij}^2$.

# PCA as Matrix Factorization

- So PCA is approximating $\mathbf{X} \approx \mathbf{UZ}$.



- Based on the sizes of the matrices, this is a rank-*K* approximation.
- Since **U** was chosen to minimize reconstruction error, this is the optimal rank-*K* approximation, in terms of $\|\mathbf{X} - \mathbf{UZ}\|_F^2$.

# PCA vs. SVD

- This has a close relationship to the Singular Value Decomposition (SVD) of **X**. This is a factorization
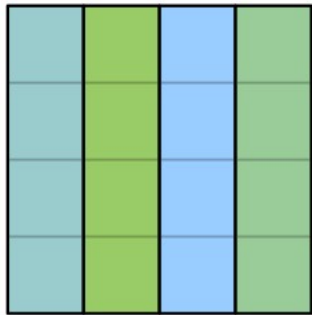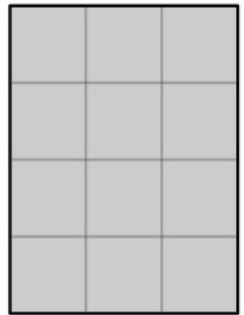
$$\mathbf{X} = \mathbf{U\Sigma V}^{\mathrm{T}}$$

Properties:

- $\mathbf{U}$, $\mathbf{\Sigma}$, and $\mathbf{V}^{\mathrm{T}}$ provide a real-valued matrix factorization of **X**, an $m \times n$ matrix.

- $\mathbf{U}$ is a $m \times m$ matrix with orthonormal columns $\mathbf{U}^{\top}\mathbf{U} = \mathbf{I}_{\mathrm{m}}$, where $\mathbf{I}_{\mathrm{m}}$ is the $m \times m$ identity matrix.

- $\mathbf{V}$ is an orthonormal $n \times n$ matrix, $\mathbf{V}^{\top}\mathbf{V} = \mathbf{I}_{\mathrm{n}}$.

- $\mathbf{\Sigma}$ is a $m \times n$ diagonal matrix, with non-negative singular values, $\sigma_1, \sigma_2, \ldots, \sigma_{\min\{m,n\}}$, on the diagonal, where the singular values are conventionally ordered from largest to smallest.
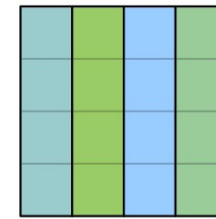
$$\mathbf{X} = \mathbf{U\Sigma V}^{\mathrm{T}} = \begin{pmatrix} | & | & & | \\ \vec{u}_1 & \vec{u}_2 & \cdots & \vec{u}_m \\ | & | & & | \end{pmatrix} \begin{pmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & 0 \\ & & \ddots & 0 \\ & & \sigma_m & 0 \end{pmatrix} \begin{pmatrix} - & \vec{v}_1^T & - \\ - & \vec{v}_2^T & - \\ & \vdots & \\ - & \vec{v}_n^T & - \end{pmatrix}$$
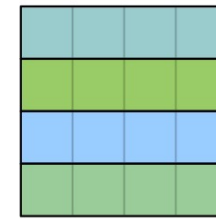
# PCA vs. SVD

$$X = U\Sigma V^T$$



$$\underset{m \times n}{X} = \underset{m \times m}{U} \quad \underset{m \times n}{\Sigma} \quad \underset{n \times n}{V^T}$$

$$U \quad U^T \quad = \quad I_m$$

$$V \quad V^T \quad = \quad I_n$$

$V$에서 n개의 singular value를 취하면 이는 PCA에서 n개의 principal component와 같다. 즉, $X = U\Sigma V^T = ZV^T$로 놓을 수 있다.

주의:   1) SVD에서 $X$는 PCA에서 $X$의 transpose된 형태로 표현됨
          2) SVD의 $U$와 PCA의 $U$는 다른 행렬임

# Rank-*K* Approximation

■ rank(A): 행렬 A의 일차독립인 행 또는 열의 최대 갯수

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^{\mathrm{T}} = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \cdots + \sigma_r u_r v_r^T$$

where $r = \min(m, n)$

■ Rank-*K* approximation:

- k개의 $\sigma_i u_i v_i^T$ 을 선택하여 행렬 A를 근사하는 것 $(1 \le k \le r)$.
- 어떤 항을 선택하느냐에 따라 A를 근사하는 정확도가 다름.

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^{\mathrm{T}} = \begin{pmatrix} | & | & & | \\ \vec{u}_1 & \vec{u}_2 & \cdots & \vec{u}_m \\ | & | & & | \end{pmatrix} \begin{pmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & 0 \\ & & \ddots & 0 \\ & & & \sigma_m & 0 \end{pmatrix} \begin{pmatrix} - & \vec{v}_1^T & - \\ - & \vec{v}_2^T & - \\ & \vdots & \\ - & \vec{v}_n^T & - \end{pmatrix}$$

# Matrix Completion

We just saw that PCA gives the optimal low-rank matrix factorization.

Two ways to generalize this:

- 1) Consider when **X** is only partially observed.
  - A sparse 1000 × 1000 matrix with 50,000 observations (only 5% observed).
  - A rank 5 approximation requires only 10,000 parameters, so it's reasonable to fit this.



  - Unfortunately, no closed form solution.

- 2) Impose structure on the factors. We can get lots of interesting models this way.

# Recommender systems: Why?

-  400 hours of video are uploaded to YouTube every

-  353 million products and 310 million users

-  83 million paying subscribers and streams about 35 million songs

Who cares about all these videos, products and songs? People may care only about a few → Personalization: Connect users with content they may use/enjoy.

Recommender systems suggest items of interest and enjoyment to people based on their preferences

# Some recommender systems in action

# Some recommender systems in action



- Ideally recommendations should combine global and session interests, look at your history if available, should adapt with time, be coherent and diverse, etc.

# The Netflix problem

■ Movie recommendation: Users watch movies and rate them as good or bad.

| User | Movie | Rating |
|------|-------|--------|
| 🥷 | Thor | ★ ☆ ☆ ☆ ☆ |
| 🥷 | Chained | ★ ★ ☆ ☆ ☆ |
| 🥷 | Frozen | ★ ★ ★ ☆ ☆ |
| 🐱 | Chained | ★ ★ ★ ★ ☆ |
| 🐱 | Bambi | ★ ★ ★ ★ ★ |
| 😇 | Titanic | ★ ★ ★ ☆ ☆ |
| 😇 | Goodfellas | ★ ★ ★ ★ ★ |
| 😇 | Dumbo | ★ ★ ★ ★ ★ |
| 😇 | Twilight | ★ ★ ☆ ☆ ☆ |
| 😛 | Frozen | ★ ★ ★ ★ ★ |
| 😀 | Tangled | ★ ☆ ☆ ☆ ☆ |

Because users only rate a few items, one would like to infer their preference for unrated items

# Matrix completion problem

Matrix completion problem: Transform the table into a *N* users by *M* movies matrix **R**


Rating matrix

- Data: Users rate some movies. $R_{user,movie}$. Very sparse

- Task: Finding missing data, e.g. for recommending new movies to users. Fill in the question marks

- Algorithms: Alternating Least Square method, Gradient Descent, Non-negative Matrix Factorization, low rank matrix Completion, etc.

# Latent factor models

- In our current setting, latent factor models attempt to explain the ratings by characterizing both movies and users on a number of factors $K$ inferred from the ratings patterns.

- That is, we seek a representation movies and users as vectors in $\mathbb{R}^K$

- For simplicity, we can associate these factors (i.e. the dimensions of the vectors) with idealized concepts like
  - comedy
  - drama
  - action
  - But also uninterpretable dimensions

Can we use the sparse ratings matrix **R** to find these latent factors automatically?

# Alternating least squares

- Let the representation of user $n$ in the $K$-dimensional space be $\mathbf{u}_n$ and the representation of movie $m$ be $\mathbf{z}_m$

- Assume the rating user $n$ gives to movie $m$ is given by a dot product:

$$R_{nm} \approx \mathbf{u}_n^T \mathbf{z}_m$$

- In matrix form, if:

$$\mathbf{U} = \begin{bmatrix} - & \mathbf{u}_1^T & - \\ & \vdots & \\ - & \mathbf{u}_N^T & - \end{bmatrix} \text{ and } \mathbf{Z} = \begin{bmatrix} | & & | \\ \mathbf{z}_1 & \cdots & \mathbf{z}_M \\ | & & | \end{bmatrix}$$

then: $\mathbf{R} \approx \mathbf{UZ}$

- This is a matrix factorization problem!

# Approach: Matrix factorization methods



$$\mathbf{R} \approx \mathbf{U} \quad \mathbf{Z}$$

# Alternating least squares

- Let $O = \{(n, m) : \text{entry } (n, m) \text{ of matrix } R \text{ is observed}\}$

- Using the squared error loss, a matrix factorization corresponds to solving

$$\min_{\mathbf{U}, \mathbf{Z}} \frac{1}{2} \sum_{(n,m) \in O} \left( R_{nm} - \mathbf{u}_n^\top \mathbf{z}_m \right)^2$$

- The objective is non-convex and in fact it's NP-hard to optimize. (See Low-Rank Matrix Approximation with Weights or Missing Data is NP-hard by Gillis and Glineur, 2011)

- As a function of either **U** or **Z** individually, the problem is convex and easy to optimize. We can use coordinate descent, just like with K-means and mixture models!

Alternating Least Squares (ALS): fix **Z** and optimize **U**, followed by fix **U** and optimize **Z**, and so on until convergence.

# Alternating least squares

$R \approx UZ$

- ALS for Matrix Completion algorithm

$$U = \begin{bmatrix} — & \mathbf{u}_1^\top & — \\ & \vdots & \\ — & \mathbf{u}_N^\top & — \end{bmatrix} \text{ and } Z = \begin{bmatrix} | & & | \\ \mathbf{z}_1 & \cdots & \mathbf{z}_M \\ | & & | \end{bmatrix}$$

1. Initialize $\mathbf{U}$ and $\mathbf{Z}$ randomly

2. repeat

3.     **for** $n = 1, .., N$ **do**

$$\begin{aligned} R_{nm} &= \mathbf{u}_n^\mathrm{T}\mathbf{z}_m \\ \Leftrightarrow R_{nm}\mathbf{z}_m^T &= \mathbf{u}_n^\mathrm{T}\mathbf{z}_m\mathbf{z}_m^T \\ \Leftrightarrow R_{nm}\mathbf{z}_m^T(\mathbf{z}_m\mathbf{z}_m^T)^{-1} &= \mathbf{u}_n^\mathrm{T} \\ \Leftrightarrow (\mathbf{z}_m\mathbf{z}_m^T)^{-1}R_{nm}\mathbf{z}_m &= \mathbf{u}_n \end{aligned}$$

4.        $\mathbf{u}_n = \left(\sum_{m:(n,m)\in O} \mathbf{z}_m\mathbf{z}_m^\top\right)^{-1} \sum_{m:(n,m)\in O} R_{nm}\mathbf{z}_m$

5.     **for** $m = 1, .., M$ **do**

6.        $\mathbf{z}_m = \left(\sum_{n:(n,m)\in O} \mathbf{u}_n\mathbf{u}_n^\top\right)^{-1} \sum_{n:(n,m)\in O} R_{nm}\mathbf{u}_n$

7. until convergence   $\cdots\cdots\to$   $\min_{\mathbf{U},\mathbf{Z}} \frac{1}{2} \sum_{(n,m)\in O} \left(R_{nm} - \mathbf{u}_n^\top\mathbf{z}_m\right)^2$

# K-Means

- It's possible to view K-means as a matrix factorization.

- Stack the indicator vectors $\mathbf{u}_n^T$ for assignments into a $N{\times}K$ matrix $\mathbf{u}$, and stack the cluster centers $\mathbf{z}_k$ into a matrix $K{\times}M$ matrix $\mathbf{Z}$.



예) k는 영화 장르, $\mathbf{u}_n^T$는 유저가 선호하는 장르의 representation, $\mathbf{z}_m$은 개별 영화의 representation, $\mathbf{z}_k$는 영화들이 장르 k와의 연관성(유사성)에 대한 representation

- "Reconstruction" of the data (replace each point with its cluster center) is given by **UZ**.

- K-means distortion function in matrix form:

영화 장르들을 cluster center로 놓고, 유저들을 각 장르에 clustering.
유저는 하나 이상의 cluster에(장르를) 할당될 수 (선호할 수) 있음.

$$\sum_{n=1}^{N}\sum_{k=1}^{K} r_k^{(n)}||\mathbf{m}_k - \mathbf{x}^{(n)}||^2 = ||\mathbf{X} - \boldsymbol{UZ}||_F^2$$

19

# K-Means

- Can sort by cluster for visualization:



$$\approx$$

# Co-clustering (optional)

- We can take this a step further.

- Idea: feature dimensions can be redundant, and some feature dimensions cluster together.

- Co-clustering clusters both the rows and columns of a data matrix, giving a block structure.

- We can represent this as the indicator matrix for rows, times the matrix of means for each block, times the indicator matrix for columns

# Sparse Coding

- Efficient coding hypothesis: the structure of our visual system is adapted to represent the visual world in an efficient way
  - E.g., be able to represent sensory signals with only a small fraction of neurons having to fire (e.g. to save energy)

- Olshausen and Field fit a sparse coding model to natural images to try to determine what's the most efficient representation.

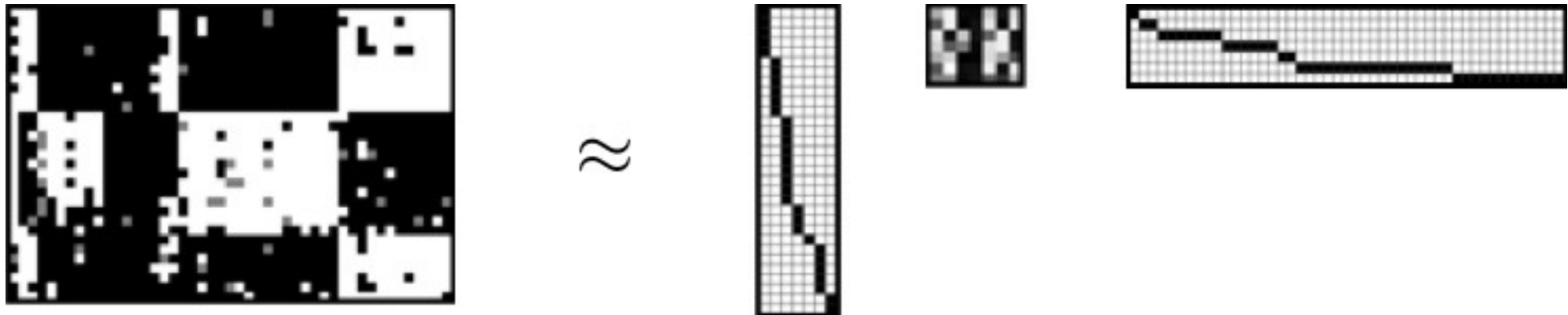- They didn't encode anything specific about the brain into their model, but the learned representations bore a striking resemblance to the representations in the primary visual cortex

- Similar to, but more general than, PCA

# Sparse Coding

- This algorithm works on small (e.g. 20 × 20) image patches, which we reshape into vectors (i.e. ignore the spatial structure)

- Suppose we have a dictionary of basis functions $\{\boldsymbol{a}_k\}_{k=1}^{K}$ which can be combined to model each patch

- Each patch is approximated as a linear combination of a small number of basis functions:

$$x = \sum_{k=1}^{K} s_k \mathbf{a}_k = \mathbf{As}$$

- This is an overcomplete representation, in that typically $K > D$ (e.g. more basis functions than pixels)

- The requirement that **s** is sparse makes things interesting

# Sparse Coding



$$\mathbf{x} \approx \sum_{k=1}^{K} s_k \mathbf{a}_k = \mathbf{As}$$

Since we use only a few basis functions, **s** is a sparse vector.

# Sparse Coding

- We'd like choose **s** to accurately reconstruct the image, but encourage sparsity in **s**.

- What cost function should we use?

- Inference in the sparse coding model:

$$\min_{\mathbf{s}} \|\mathbf{x} - \mathbf{As}\|^2 + \beta\|\mathbf{s}\|_1$$

- Here, $\beta$ is a hyperparameter that trades off reconstruction error vs. sparsity.

- There are efficient algorithms for minimizing this cost function (beyond the scope of this class)

# Sparse Coding: Learning the Dictionary

- We can learn a dictionary by optimizing both **A** and $\{s_i\}_{i=1}^{N}$ to trade off reconstruction error and sparsity

$$\min_{\{\mathbf{s}_i\},\mathbf{A}} \quad \sum_{i=1}^{N} \|\mathbf{x} - \mathbf{A}\mathbf{s}_i\|^2 + \beta\|\mathbf{s}_i\|_1$$
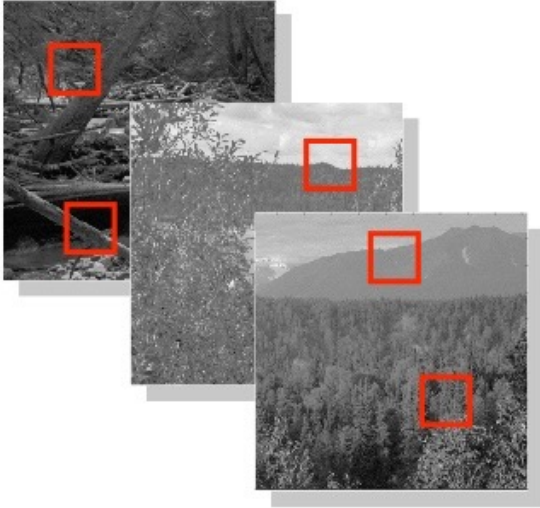
$$\text{subject to} \quad \|\mathbf{a}_k\|^2 \leq 1 \text{ for all } k$$

- Why is the normalization constraint on $\boldsymbol{a}_k$ needed?
  - Regularization term 때문에 $\mathbf{s}_i$는 작아지면서 $\|\mathbf{x} - \mathbf{A}\mathbf{s}_i\|^2$ 를 줄이기 위해 $\boldsymbol{a}_k$가 커짐

- Reconstruction term can be written in matrix form as $\|\mathbf{X} - \mathbf{A}\mathbf{S}\|_F^2$ , where **S** combines the $\boldsymbol{s}_i$ as columns

- Can fit using an alternating minimization scheme over **A** and **S**, just like K-means, EM, low-rank matrix completion, etc.
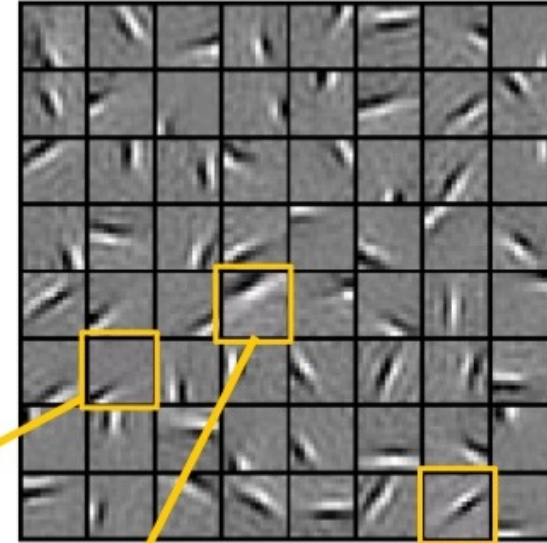
# Sparse Coding: Learning the Dictionary
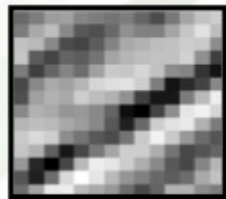
- Basis functions learned from natural images:



Natural Images

Learned bases $(\phi_1, \phi_2, \ldots, \phi_{64})$: "Edges"
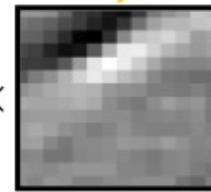
Test Example

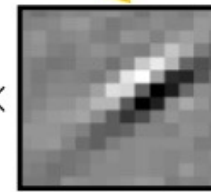$$\approx \quad 0.8 \times \underset{\phi_{36}}{\ } \quad + \quad 0.3 \times \underset{\phi_{42}}{\ } \quad + \quad 0.5 \times \underset{\phi_{63}}{\ }$$
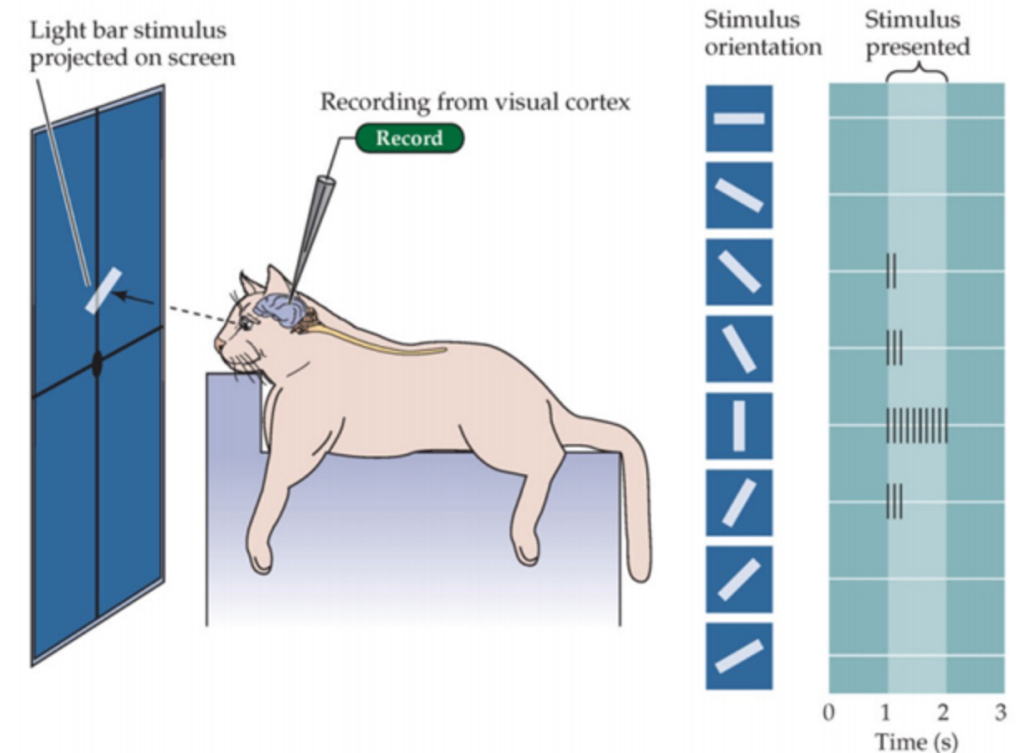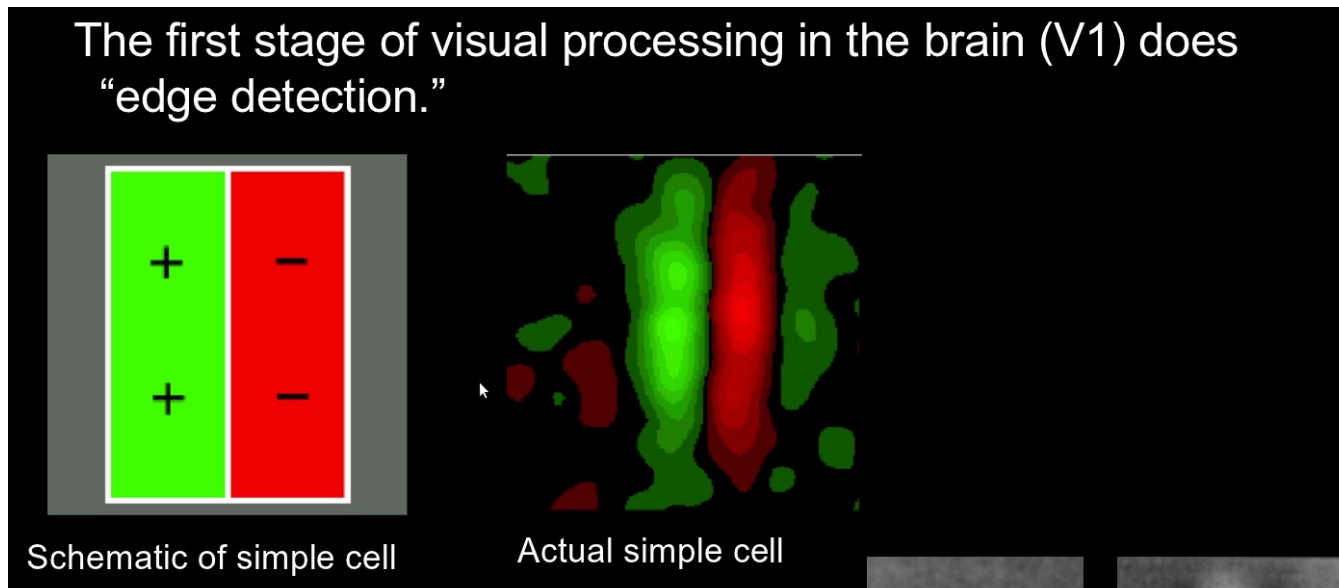
$$[\alpha_1, \ldots, \alpha_{64}] = [0, \ldots, 0.8, \ldots, 0.3, \ldots, 0.5, \ldots, 0]$$

# Sparse Coding: Learning the Dictionary

- The sparse components are oriented edges, similar to what a conv net learns

- But the learned dictionary is much more diverse than the first-layer conv net representations: tiles the space of location, frequency, and orientation in an efficient way

- Each basis function has similar response properties to cells in the primary visual cortex (the first stage of visual processing in the brain)



The first stage of visual processing in the brain (V1) does "edge detection."

Schematic of simple cell          Actual simple cell



Light bar stimulus projected on screen

Recording from visual cortex
Record

Stimulus orientation    Stimulus presented

0  1  2  3
Time (s)

28

# Example: Image Denoising

# Example: Image Restoration

# Summary

- PCA can be viewed as fitting the optimal low-rank approximation to a data matrix.

- Matrix completion is the setting where the data matrix is only partially observed
  - Solve using ALS, an alternating procedure analogous to EM

- PCA, K-means, co-clustering, sparse coding, and lots of other interesting models can be viewed as matrix factorizations, with different kinds of structure imposed on the factors.