

## Week 10

↳ Midterms

↳ Anonymous eval!

- generally, reflect effort in class (OH + study sessions + psets)

- goal is learning:

↳ redo one long question for up to 1/2 pts lost

↳ I'll replace 1st midterm grade w/ 2nd, if better

- let's chart in terms!

↳ Guideline of 12 hours per class per week

|       |                        |    |               |         |   |
|-------|------------------------|----|---------------|---------|---|
| 2.5   | 1.5                    | 2  | 2             | 2       | 2 |
| class | prep attempt           | OH | study session | attempt |   |
|       | • reading #1<br>• pset |    |               | #2      |   |

## Beyond Supervised Learning

Q : Lots of unlabelled data (images, text),  
how do we use it?

This week, we will explore pretraining

## Auto encoders

Idea: Use data as its own label

$$f_0: \mathbb{R}^d \rightarrow \mathbb{R}^k$$

"Encoder"

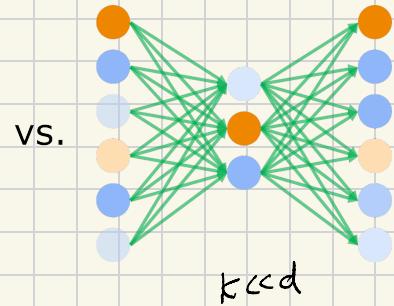
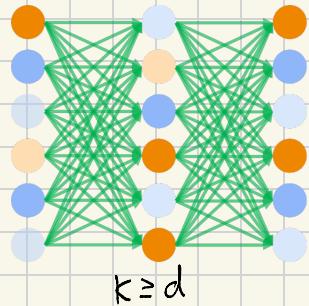
$$f_1: \mathbb{R}^k \rightarrow \mathbb{R}^d$$

"Decoder"

$z = f_0(x)$  is latent representation

$\tilde{x} = f_1(z)$  is reconstruction

$$\mathcal{L}_{\text{recon}} = \frac{1}{n} \sum_{i=1}^n \|x^{(i)} - \tilde{x}^{(i)}\|_2^2$$

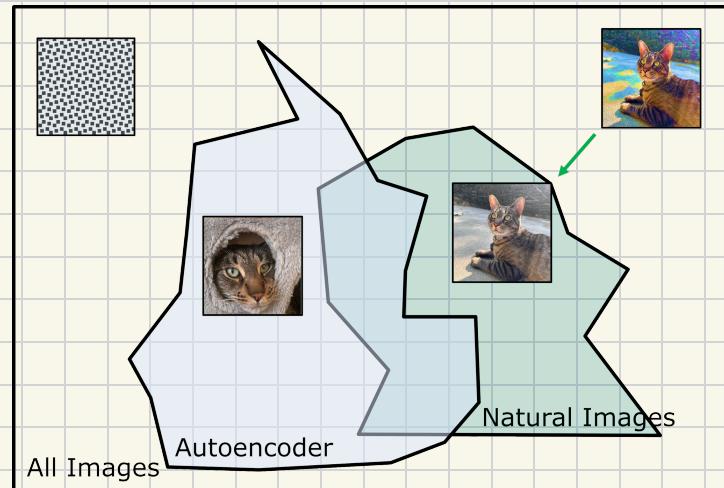


Applications include:

- data compression
- denoising
- inpainting
- representation learning

## Data Manifold

Necessarily losing information, okay because we only want to represent meaningful images



Prefer working with latent rep.

1. more efficient in lower dim
2. most meaningful features

Examples include:

- classification
- clustering
- generation

## Variational Autoencoders

Goal: make latent space "nicely behaved"

Treat latent as random sample

$$z \sim N(\mu_x, \Sigma_x)$$

↑ mean      ↙ variance  
 from encoder      from decoder

"Nicely behaved"  $\Leftrightarrow N(\mu_x, \Sigma_x) \approx N(0, I)$

Distribution  $P \approx$  Distribution  $Q$

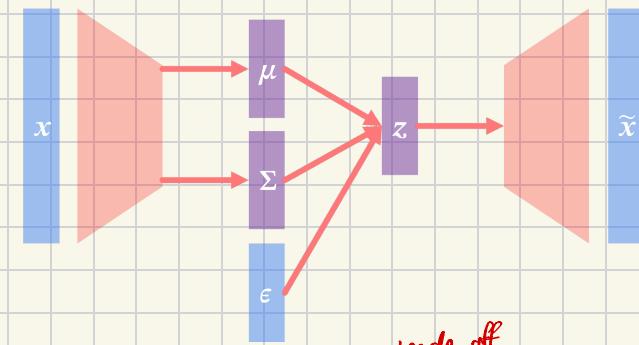
Loss: KL divergence

$$D_{KL}(P||Q) = E_{z \sim P} \left[ \log \frac{P(z)}{Q(z)} \right]$$

Example in  $\mathbb{R}$ :  $P = N(\mu_A, I)$ ,  $Q = N(\mu_B, I)$

How do we backpropagate through a sample?

$$z = \mu_x + \sum_x^{1/2} \epsilon \in N(0, I) \text{ draw}$$



trade off

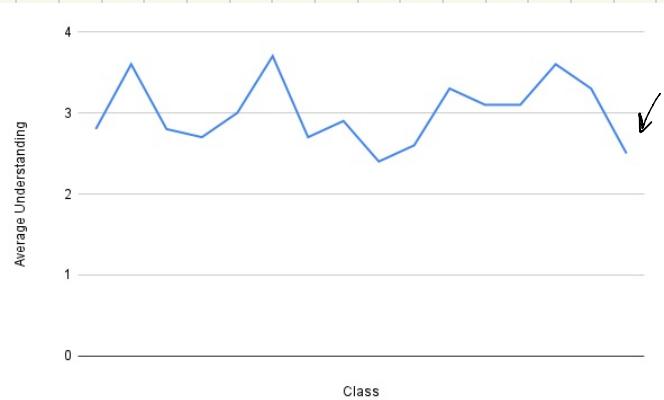
$$\mathcal{L}_{VAE} = \mathcal{L}_{\text{recon}} + \lambda D_{KL}(P||Q)$$

↑  
 z encodes  
 meaning of x  
 ↑  
 z distributed  
 similarly

## Logistics

- Midterm Q Redo due 10/31

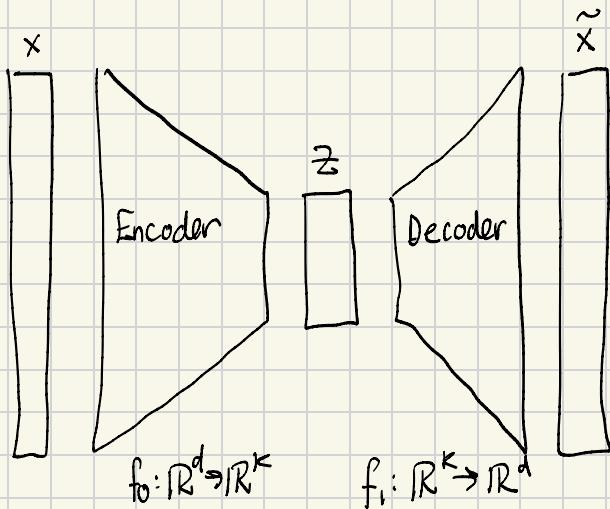
↳ instructions on discord



first section  
had worse  
pacing + explanation,  
I'm sorry :-)

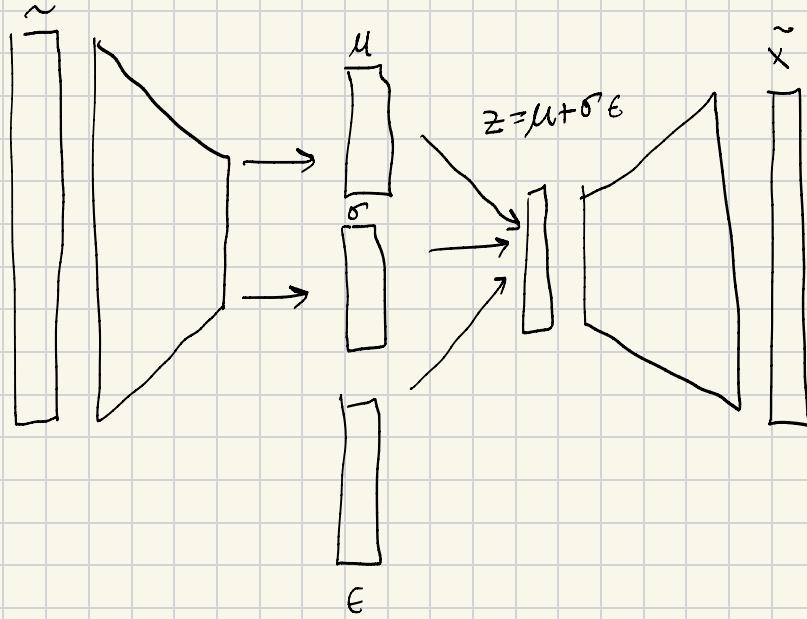
Please do reading  
and come to office hours!

## Autoencoder



## Variational Autoencoder

Distribution of latent is "nice"



$$\mathcal{L}_{\text{recon}} = \sum_{i=1}^n \|x^{(i)} - \tilde{x}^{(i)}\|_2^2$$

$$\mathcal{L}_{\text{dist}} = \sum_{i=1}^n D_{\text{KL}}(N(\mu_{x^{(i)}}, \sigma_{x^{(i)}}^2) || N(0, 1))$$

$$\mathcal{L} = \mathcal{L}_{\text{recon}} + \lambda \mathcal{L}_{\text{dist}}$$

## Principal Component Analysis

The "linear regression" of autoencoders

$$f_0 = W_0 \in \mathbb{R}^{k \times d} \quad f_1 = W_1 \in \mathbb{R}^{d \times k}$$

$$z^T = x^T W_0$$

$$\boxed{\phantom{0}} = \boxed{\phantom{0}} \boxed{W_0}$$

$$\tilde{x}^T = x^T W_0 \quad W_1 = z^T W_1$$

$$\boxed{\phantom{0}} = \boxed{\phantom{0}} \boxed{W_1}$$

$$= \boxed{\phantom{0}} \boxed{W_0} \boxed{W_1}$$

$$x \in \mathbb{R}^{n \times d}$$

$$\begin{matrix} x^{(1)T} \\ \vdots \\ x^{(n)T} \\ x \end{matrix} \boxed{W_0} \boxed{W_1} = \begin{matrix} \tilde{x}^{(1)T} \\ \vdots \\ \tilde{x}^{(n)T} \\ \tilde{x} \end{matrix}$$

$$\begin{aligned} \mathcal{L}_{\text{recon}} &= \sum_{i=1}^n \|x^{(i)} - \tilde{x}^{(i)}\|_2^2 \\ &= \sum_{i=1}^n \sum_{j=1}^d (x_j^{(i)} - \tilde{x}_j^{(i)})^2 \\ &= \|x - \tilde{x}\|_F^2 \end{aligned}$$

sum of  
every entry  
squared

Q: What is the best  $W_0$  and  $W_1$ ?

## Singular Value Decomposition

$$X = U \Sigma V^T \quad \checkmark \quad r = \text{rank}$$

$n \times d \quad n \times r \quad r \times r \quad r \times d$

$$X = \begin{matrix} u^{(i)} \\ \Sigma \\ v^{(i)T} \end{matrix}$$

Diagram illustrating the decomposition:

- $X$  is a rectangle.
- $U$  is a vertical column of rectangles, each labeled  $u^{(i)}$ .
- $\Sigma$  is a square matrix with diagonal entries labeled  $\sigma_i$ .
- $V^T$  is a vertical column of rectangles, each labeled  $v^{(i)T}$ .

$$X = \sum_{i=1}^r u_i \sigma_i v_i^T$$

$$U_i \text{ orthonormal} \quad U^T U = I_{r \times r}$$

$$V_i \text{ orthonormal} \quad V^T V = I_{r \times r}$$

$$\text{Fact: } \|X\|_F^2 = \sum_{i=1}^r \hat{\sigma}_i^2$$

Q: What is the best rank  $k$  (latent dim) approximation to  $X$ ?

Assume  $k \leq r$

$$X_k = \sum_{i=1}^k u_i \sigma_i v_i^T$$

$$\|X - X_k\|_F^2 = ?$$

## Back to PCA

$$\text{Goal: } \tilde{X} = X_k$$

Q: What are  $W_0$  and  $W_1$ ?

A:  $W_0 = V_k$      $W_1 = V_k^T$   
 $d \times k$  ✓     $k \times d$  ✓

$$\tilde{X} = X W_0 W_1^T$$

$$V_k V_k^T = \sum_{i=1}^k v_i v_i^T$$

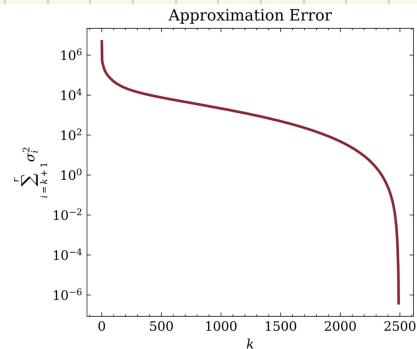
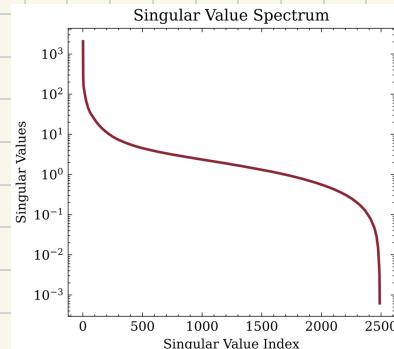
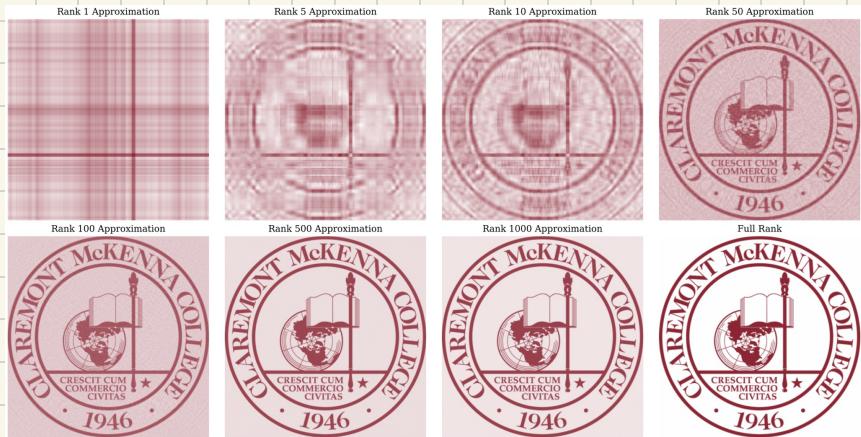
$$\begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad \begin{bmatrix} -v_i^T \\ 1 \end{bmatrix}$$

$$\tilde{X} = \sum_{j=1}^r u_j \sigma_j v_j^T \sum_{i=1}^k v_i v_i^T$$

$$= \sum_{j=1}^k u_j \sigma_j v_j^T = X_k$$

$$X W_0 W_1^T = X_k$$

reconstruction

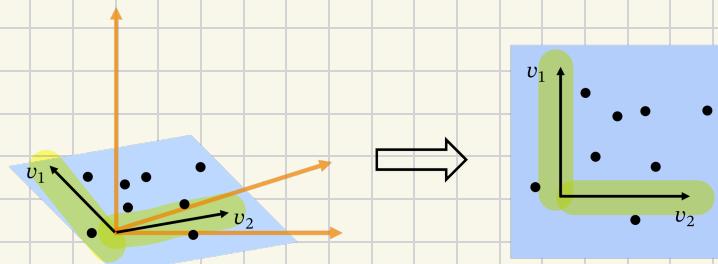


## Latent Representation

$$X W_o = X V_k \quad \text{latent representation}$$

$X$ :  $n \times d$     $W_o$ :  $d \times k$     $V_k$ :  $n \times k$

$$z^{(i)} = x^{(i)T} V_k = \sum_{j=1}^k \langle x^{(i)}, v_j \rangle$$



$V_k$

principal components

## Semantic Embeddings

$X \in \mathbb{R}^{n \times d}$  "document-word" matrix

$$X_{ij} = \mathbb{I}[ \text{word } j \text{ appears in doc } i ]$$

|                  | cat | dog | kitten | puppy | great | good |
|------------------|-----|-----|--------|-------|-------|------|
| doc <sub>1</sub> | 1   |     | 1      | 1     |       | 1    |
| doc <sub>2</sub> |     | 1   |        |       | 1     |      |
| doc <sub>3</sub> |     | 1   |        | 1     |       |      |
| :                |     |     |        |       |       |      |

$z^T$  latent representation of document

$v_i, v_j$  principal components of words

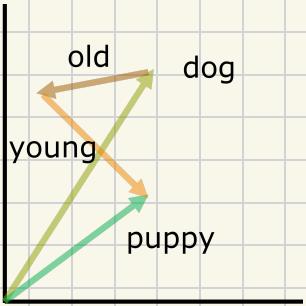
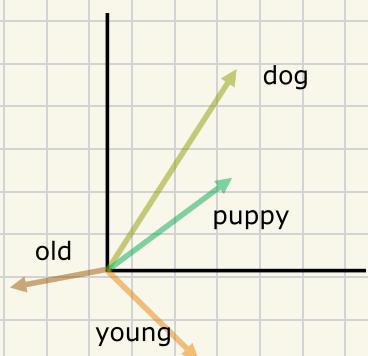
If  $\langle z, v_i \rangle \approx 1$  and  $\langle z, v_j \rangle \approx 1$ ,  $v_i$  and  $v_j$  are close!

If  $v_i$  and  $v_j$  are in the same doc,  $v_i$  and  $v_j$  are close!

[d words and n documents]

$$\begin{matrix} & i & j \\ & \vdots & \vdots \\ \approx & \boxed{\square \quad \square} & = & \boxed{\vdots \quad z^T \quad \vdots} & \boxed{v_k^T} \\ & x_k & & u_k \Sigma_k & k \times d \end{matrix}$$

## Word Math



$$\|Q \alpha\|_2^2 = \alpha^T Q^T Q \alpha$$

## Unsupervised Translation

