

CS 189/289

Today's lecture outline:

1. Finish Transformers
2. Unsupervised learning, dimensionality reduction
3. PCA

Assigned reading: 16.1 (PCA)

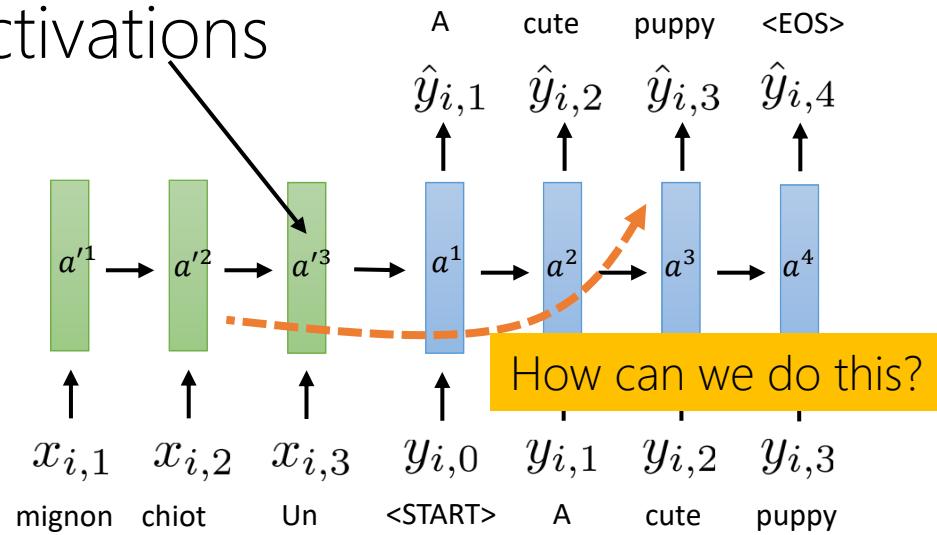
CS 189/289

Today's lecture outline

1. Finish Transformers
2. Unsupervised learning, dimensionality reduction
3. PCA

Recall: RNN bottleneck problem

all information about the conditioned sequence is contained in these activations



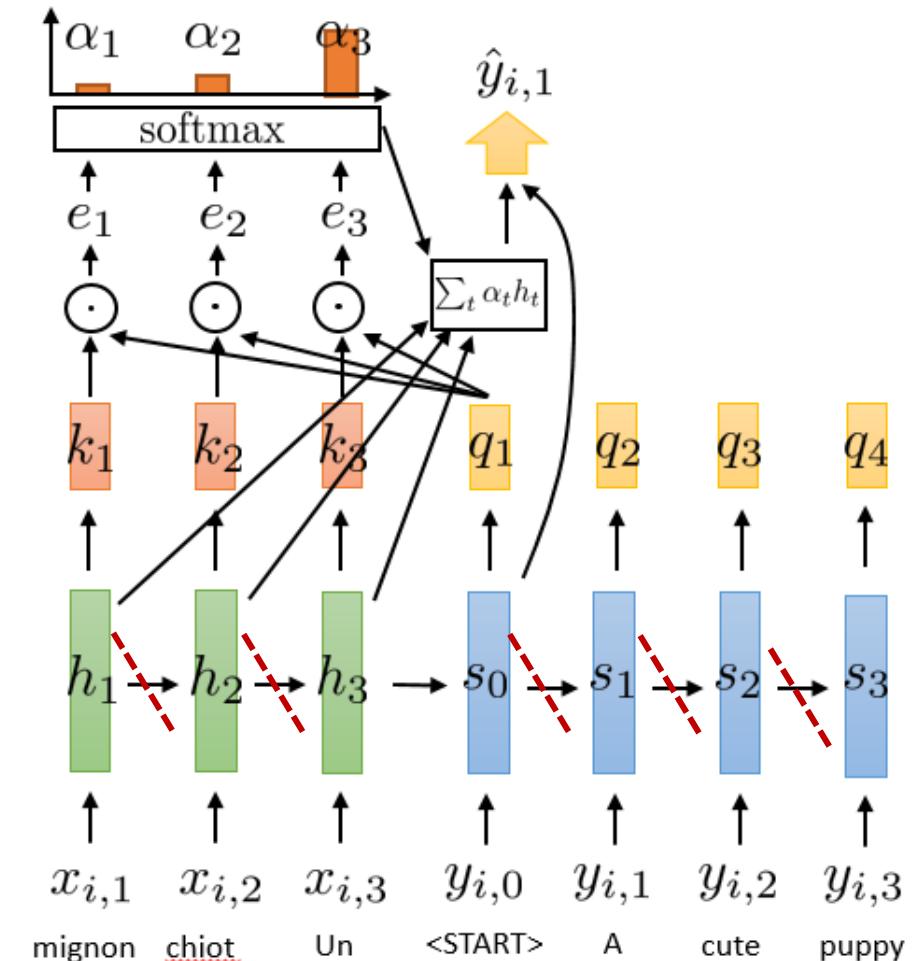
Idea: what if we could somehow “peek” at the source sentence while decoding?
Attention to the rescue!

Recall: Is Attention All We Need?

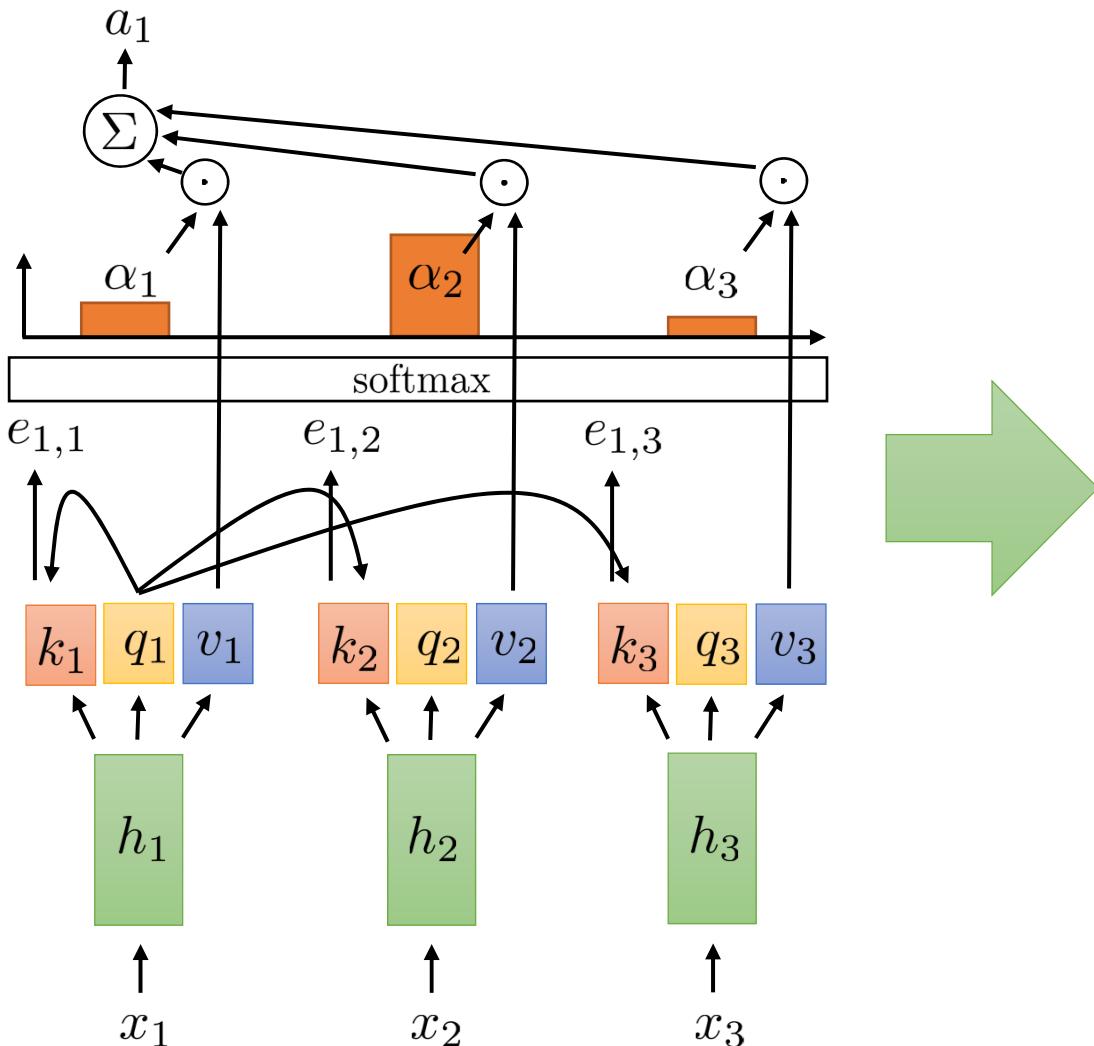
- If we have attention, do we even need recurrent connections?
- Can we transform our RNN into a purely attention-based model?

This has a few issues we must overcome:

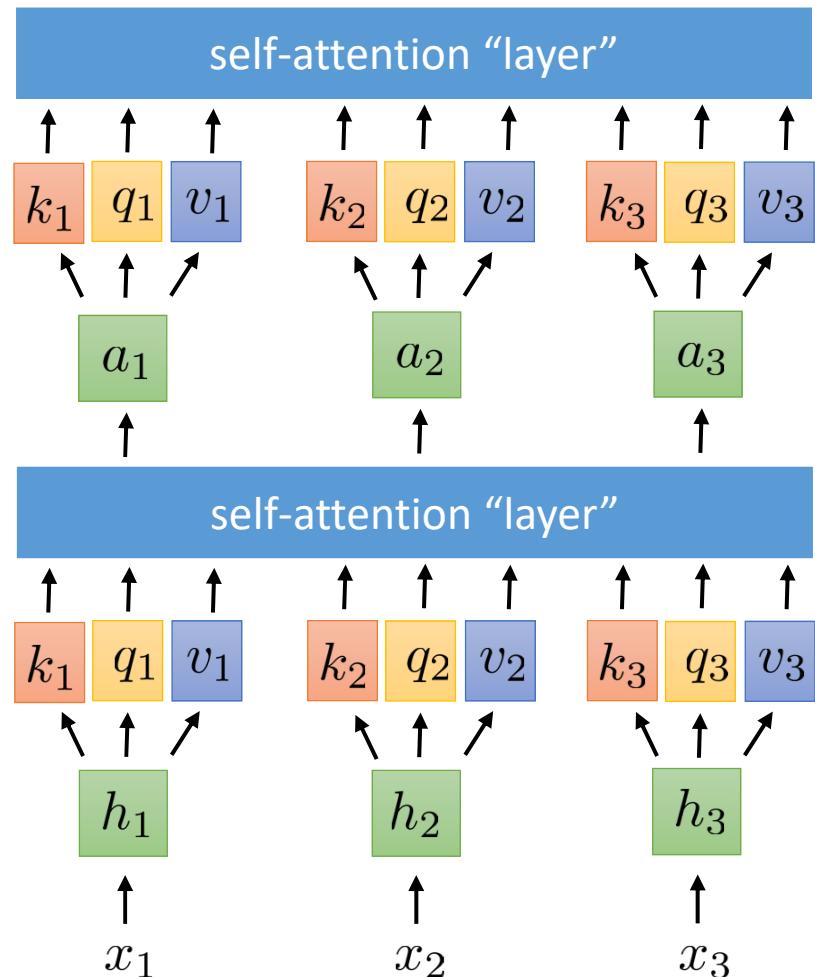
- Decoding position 3 can't access s_1 or s_0 .
- Solution: self-attention.



Recall: Self-Attention



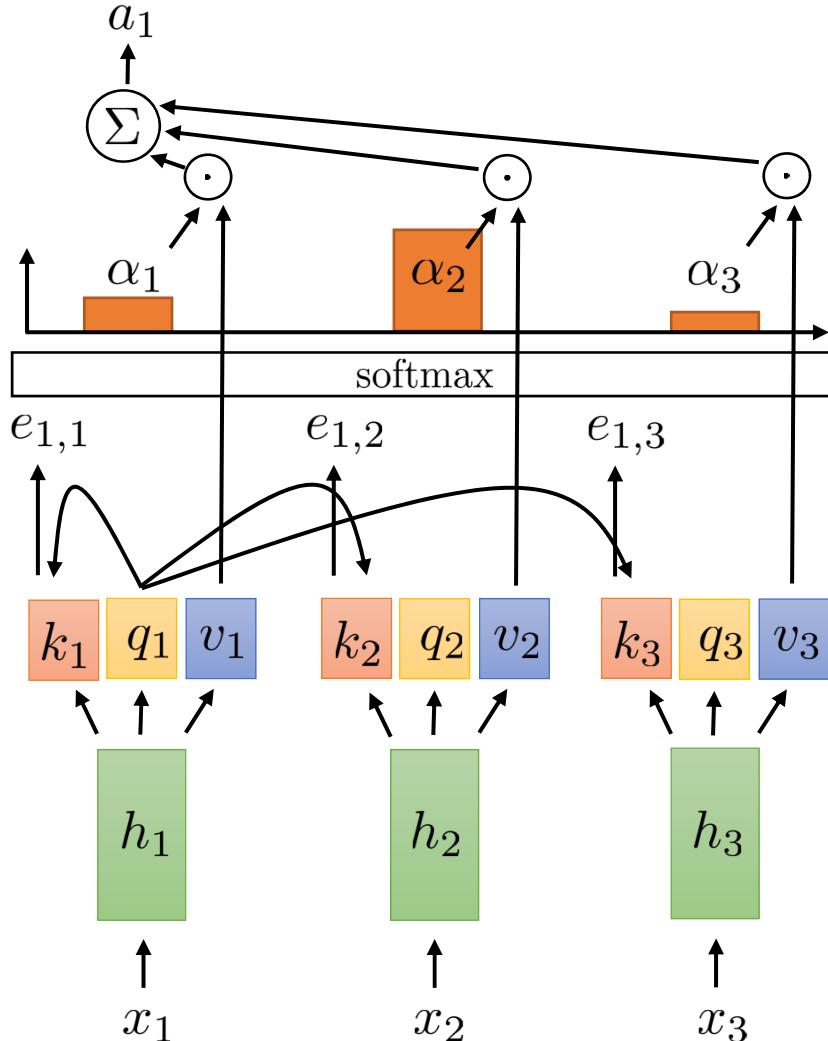
keep repeating until we've
processed this enough
then hand off to next part of overall
model



From Self-Attention to Transformers

- The basic concept of **self-attention** can be used to develop a very powerful type of sequence model, called a **transformer**
 - But to make this actually work, we need to develop a few additional components to address some fundamental limitations
- | | |
|---------------------------|--|
| 1. Positional encoding | addresses lack of sequence information |
| 2. Multi-headed attention | allows querying multiple positions at each layer |
| 3. Adding nonlinearities | so far, each successive layer is <i>linear</i> in the previous one |
| 4. Masked decoding | how to prevent attention lookups into the future? |

Self-Attention is Linear



$$k_t = W_k h_t \quad q_t = W_q h_t \quad v_t = W_v h_t$$

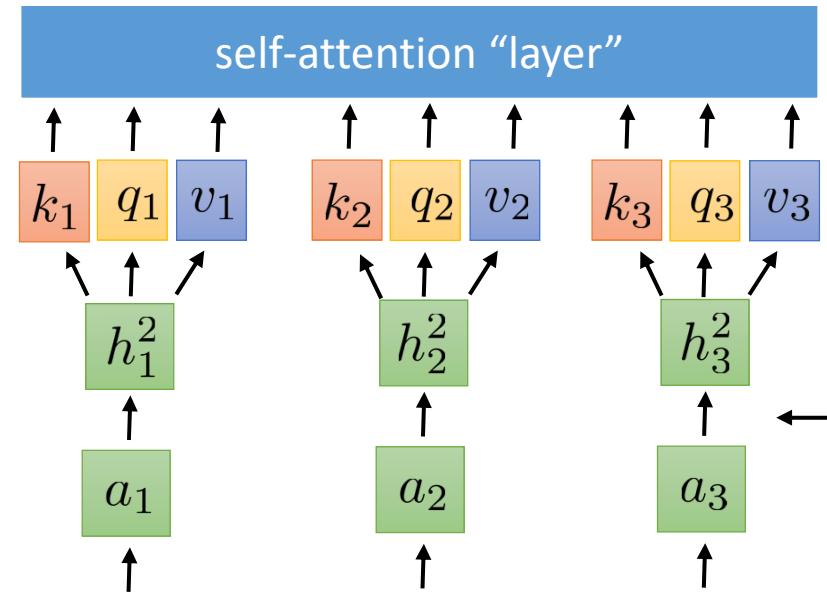
$$\alpha_{l,t} = \exp(e_{l,t}) / \sum_{t'} \exp(e_{l,t'})$$

$$e_{l,t} = q_l \cdot k_t$$

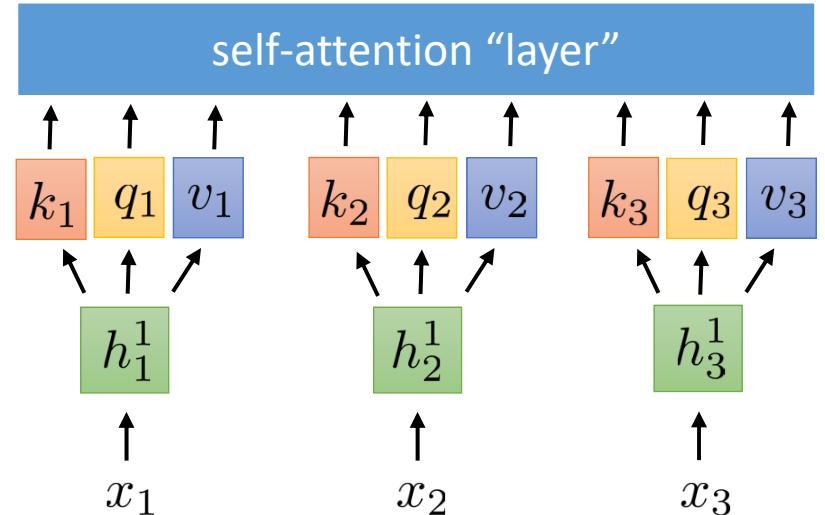
$$a_l = \sum_t \alpha_{l,t} v_t = \sum_t \alpha_{l,t} W_v h_t = W_v \sum_t \alpha_{l,t} h_t$$

Every self-attention “layer” is a linear transformation of the previous layer

Alternating self-attention & non-linearity



some non-linear (learned) function
e.g., $h_t^\ell = \sigma(W^\ell a_t^\ell + b^\ell)$

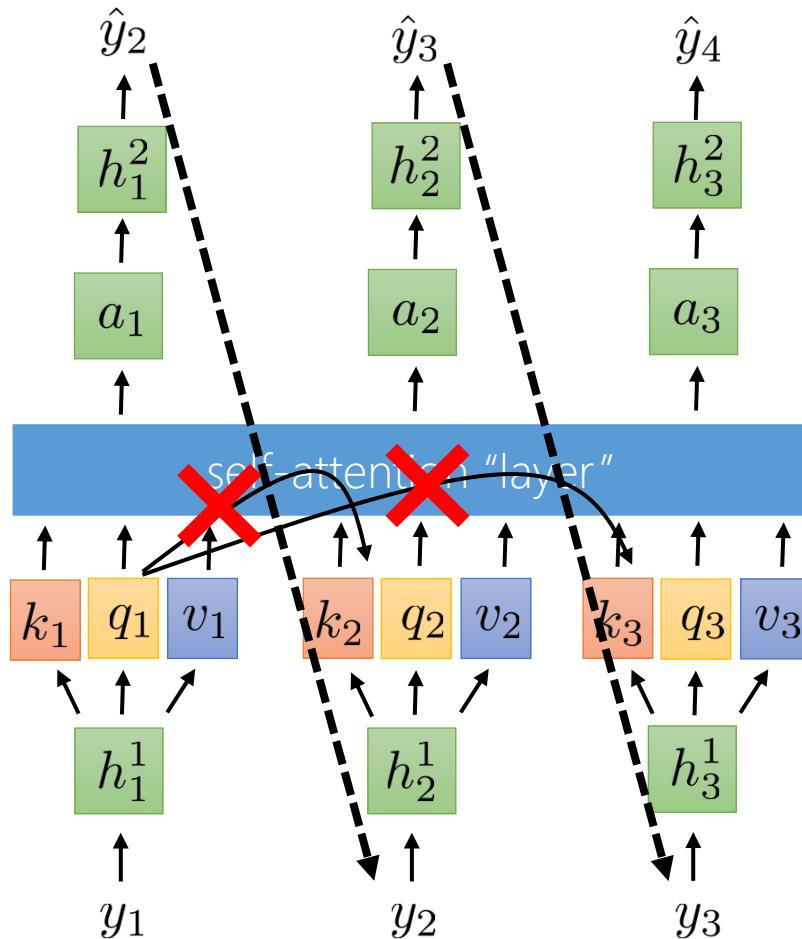


just a simple neural net applied at every position after every self-attention layer

From Self-Attention to Transformers

- The basic concept of **self-attention** can be used to develop a very powerful type of sequence model, called a **transformer**
- But to make this actually work, we need to develop a few additional components to address some fundamental limitations
 - 1. Positional encoding addresses lack of sequence information
 - 2. Multi-headed attention allows querying multiple positions at each layer
 - 3. Adding nonlinearities so far, each successive layer is *linear* in the previous one
 - 4. Masked decoding how to prevent attention lookups into the future?

Self-attention can see the future!



Problem:

- Each position looks at future values.
- Also cyclic: output 1 depends on input 2 which depends on output 1.

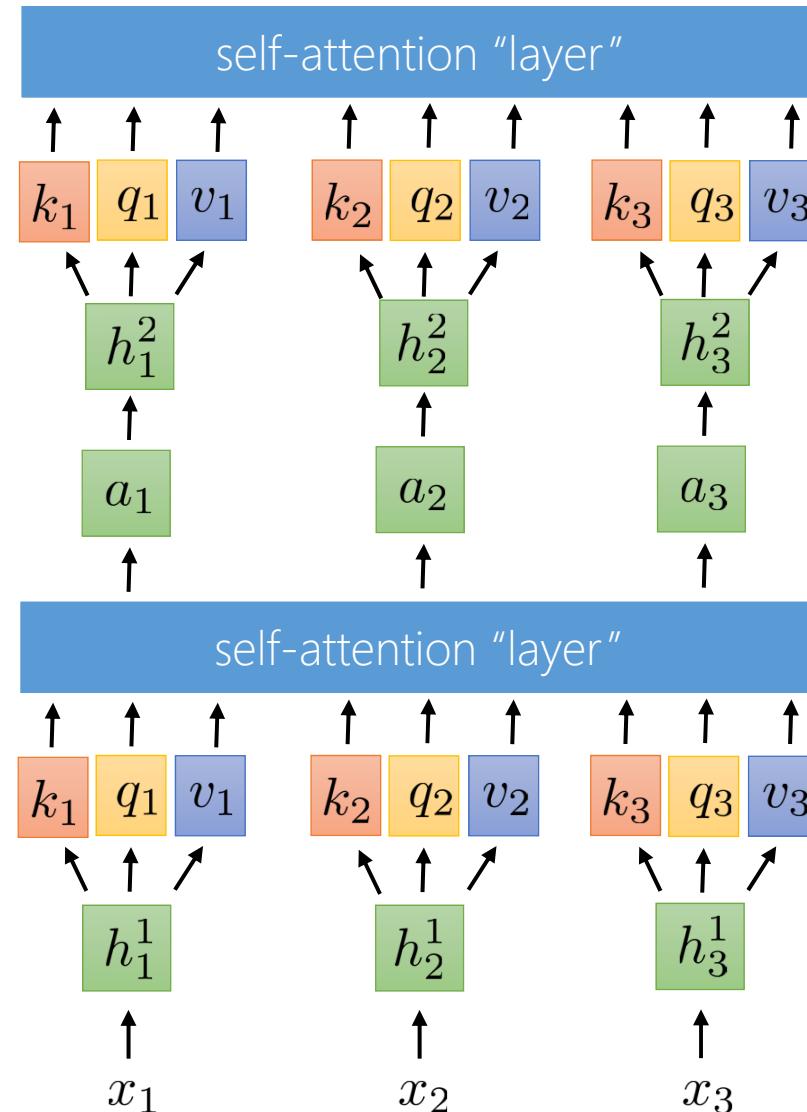
Solution ("masked attention"):

$$e_{l,t} = \begin{cases} q_l \cdot k_t & \text{if } t \leq l \\ -\infty & \text{otherwise} \end{cases}$$

$$a_l = \sum_t \alpha_{l,t} v_t$$
$$\alpha_{l,t} = \exp(e_{l,t}) / \sum_{t'} \exp(e_{l,t'})$$

Now we are read for
The Transformer!

Sequence-to-sequence with self-attention

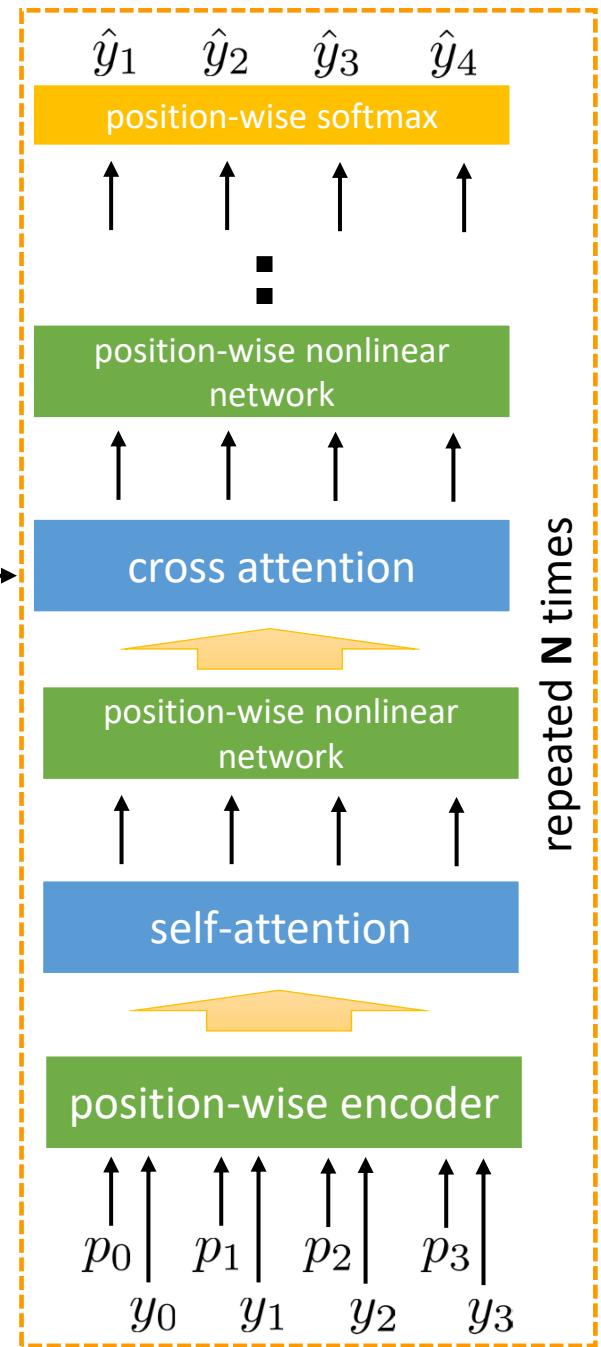
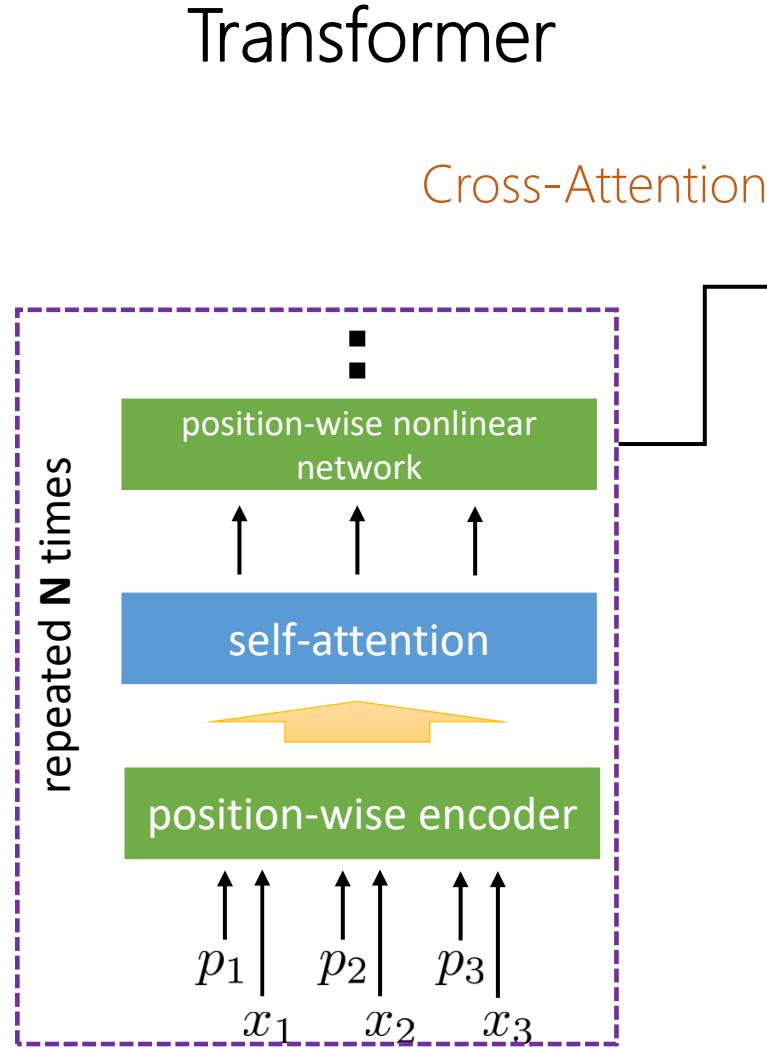
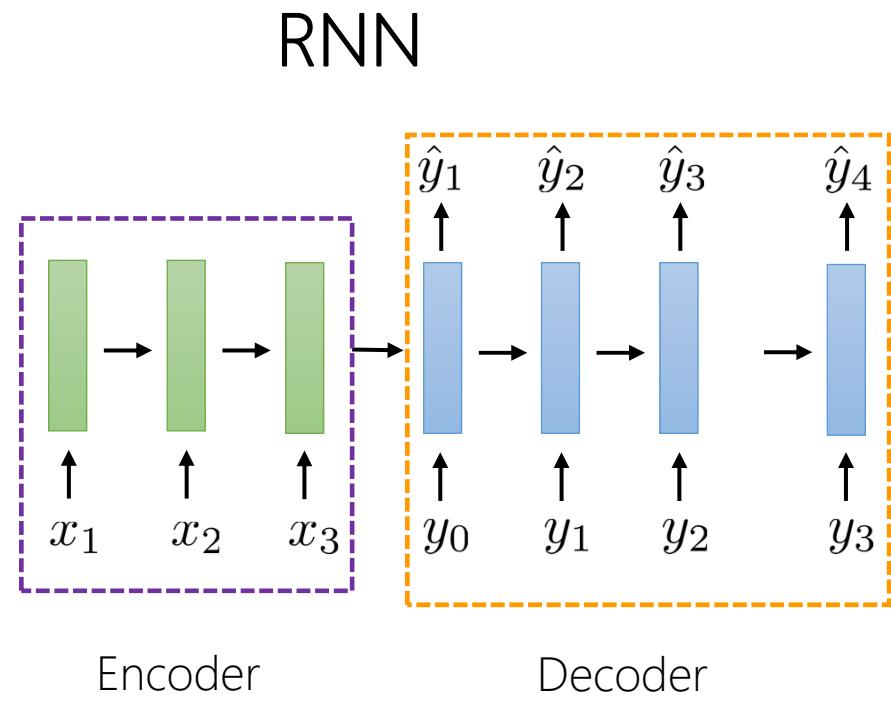


"*Transformer*" architecture:

- Stacked self-attention layers with position-wise nonlinearities.
- *Transform* one sequence into another at **each** layer.
- For sequence data.

[Vaswani et al. Attention Is All You Need. 2017]

Encoder-Decoder Transformer



Transformers pros and cons

Downsides:

- Attention computations are theoretically* $O(n^2)$.
- Slightly more complex to implement (positional encodings, etc.)

Benefits:

- + Better long-range connections (compared to RNN).
- + *Much easier to parallelize.
- + In practice, can make it much deeper (more layers) than RNN.

- Benefits often vastly outweigh the downsides.
- Transformers work much better than RNNs in general.
- One of the most important sequence modeling improvements of the past decade.
- Can use just encoder, or just decoder, or both.

CS 189/289

Today's lecture outline

1. Finish Transformers
2. Unsupervised learning, dimensionality reduction
3. PCA

Unsupervised learning

So far: *supervised learning*, $\{(x_i, y_i)\}$ for $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$ or $y \in \mathbb{Z}$.

Often model just $\{x_i\}$: *unsupervised learning*, includes:

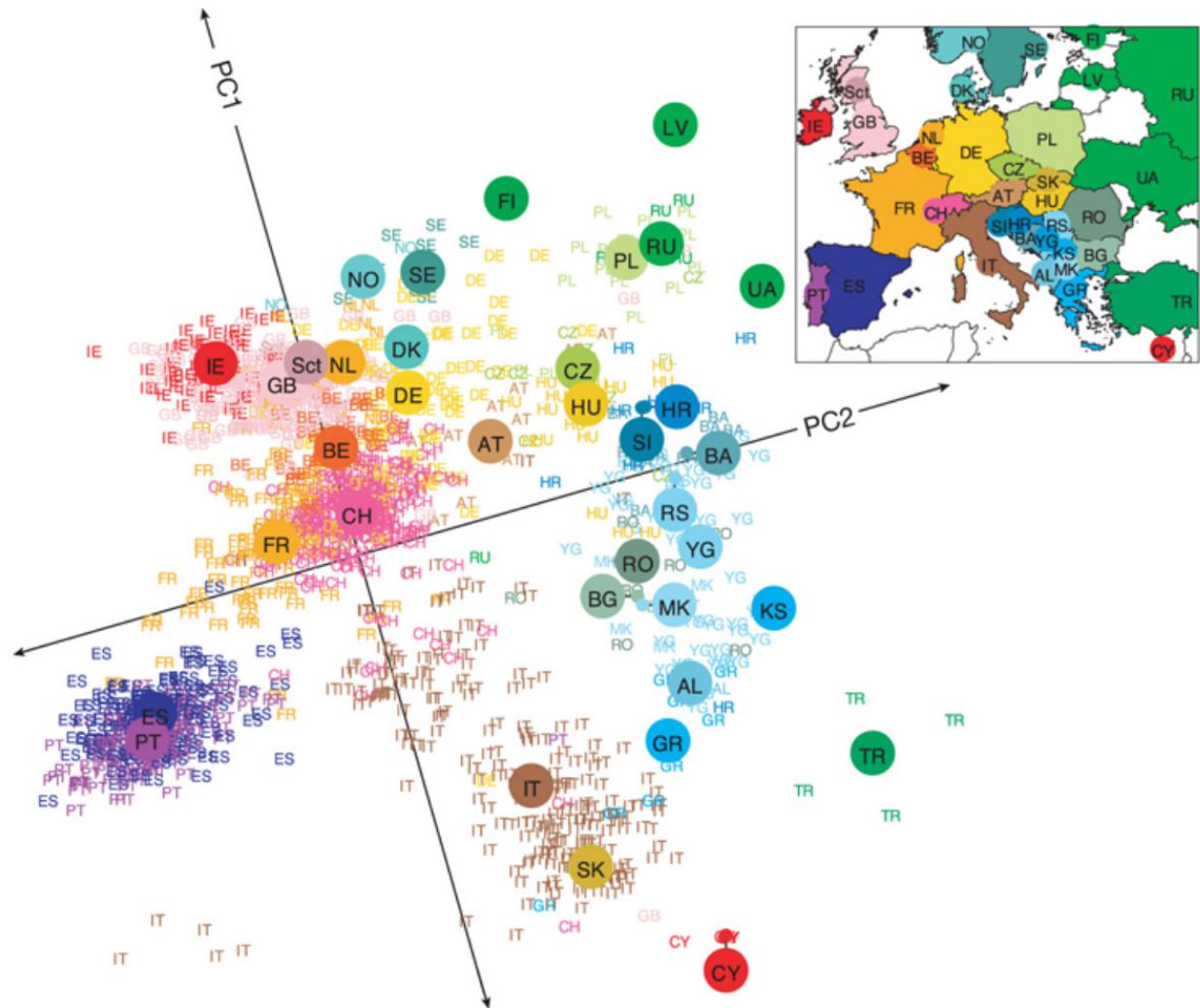
- i. *Dimensionality reduction*, $z \in \mathbb{R}^m = f_\theta(x \in \mathbb{R}^d)$, $m \ll d$.
- ii. *Clustering*, for each x_i , assign cluster label, $z_i \in \{1,2,3 \dots K\}$
- iii. *Representation learning*, $z \in \mathbb{R}^m$, $z = f_\theta(x)$, or $z \sim p_\theta(x)$.
- iv. *Density estimation*, evaluate $p_\theta(x)$.
- v. "Generative" modeling, $x \sim p_\theta(x)$

e.g. of Dimensionality reduction



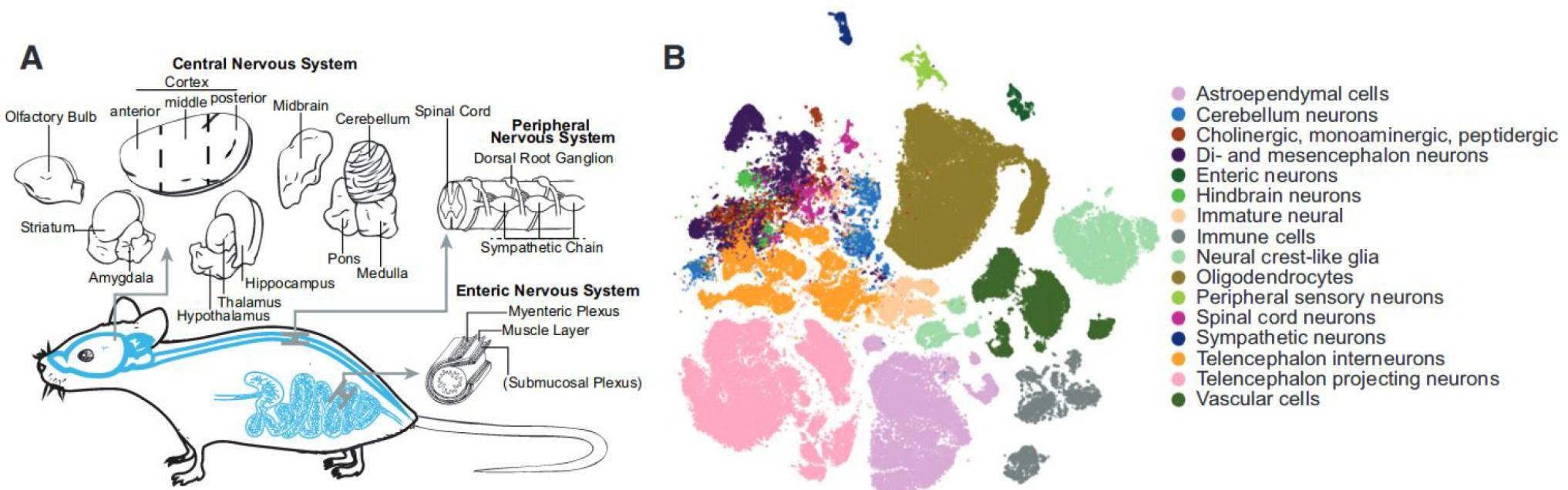
ATGTTGAATCTG
AAAGTGAAATGT
TATTATAACGAAG
AAGTATTTGCTA
GACCTCAAAAACC.
CTTCATCATAAC.

*DNA from set of
individuals*

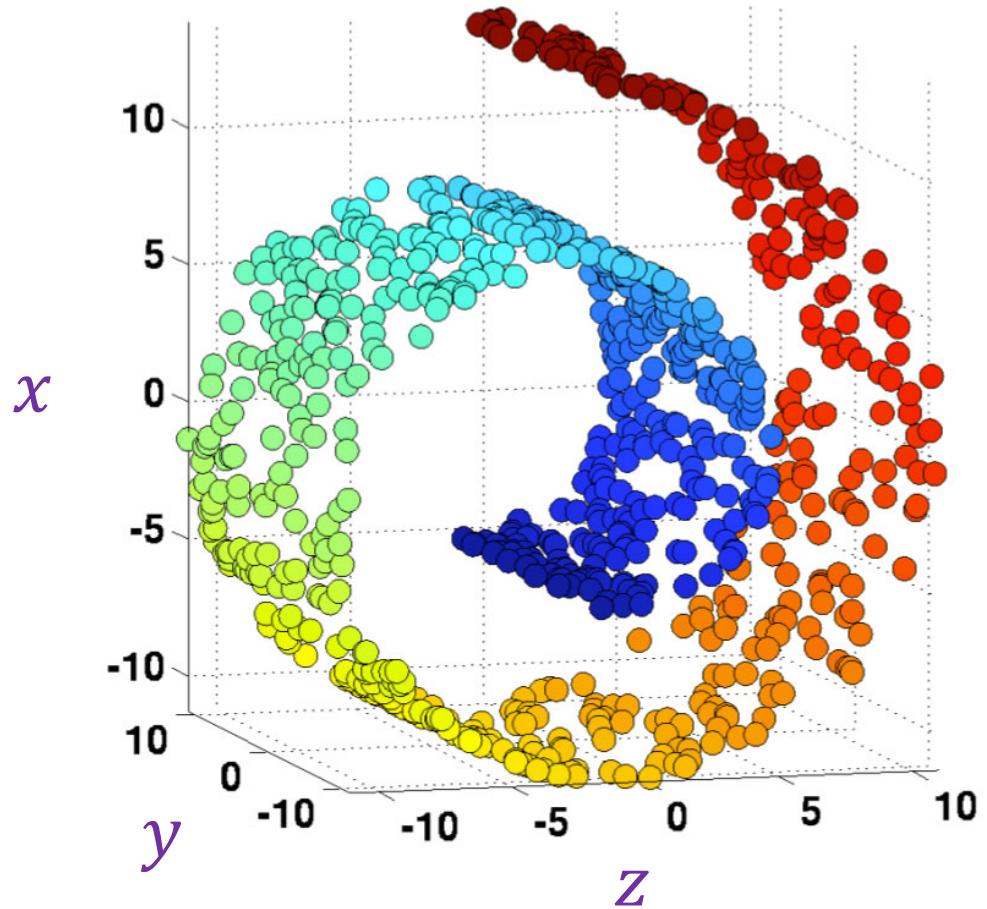


e.g. of dimensionality reduction+clustering

Single-cell transcriptomics (single-cell RNA sequencing): samples are cells, features are genes.

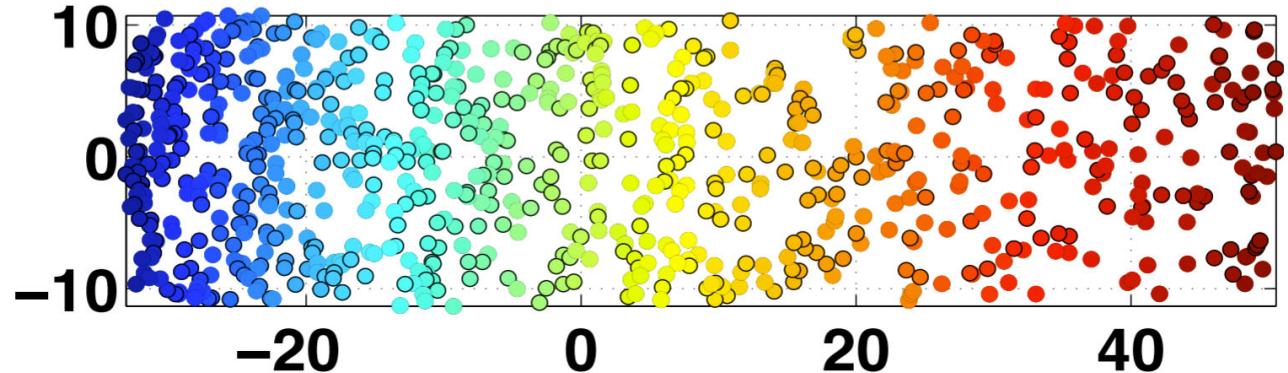


What dimensionality do these points live in?

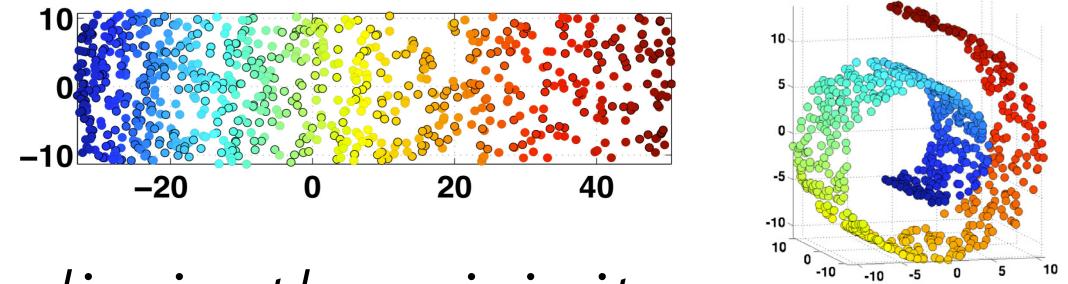


$$x \in \mathbb{R}^3 = [x, y, z].$$

But could uniquely describe each point with just 2 coordinates.



The “manifold hypothesis”



- *“High dimensional data tend to lie in the vicinity of a low dimensional manifold”* e.g. [Fefferman 2013]
- Manifold: roughly speaking, a space that locally feels like a Euclidean space.
- For us: a manifold is a lower dimensional part of the observation space in which the data tend to lie.
- *“embedding the data in a lower dimensional manifold”, or “an embedding of the data”.*

What dimensionality do these points live in?



- 5000 faces, $x_i \in \mathbb{R}^{32 \times 32=1024}$
- How low a dimension do you think we can go and still “keep” the image?
- Turns out we can go down to ~ 100 from the “ambient” **1024** dimensions!
- Trick: carefully create 100 special “basis” images.
- Principal Components Analysis (PCA) will yield the PC basis vectors.

What dimensionality do these points live in?



- 5000 faces
 - How low can we go?
 - Turns 1024 dimensions into 1024 “basis” images
 - Trick: consider the principal components of A) will tell us
-
- 22x22 1024
- The image shows a 5x5 grid of 25 grayscale images. Each image is a 22x22 pixel PCA basis image. The images are arranged in a single row. The first image in the top-left corner has the text "22x22 1024" above it.

PCA “basis” images, $x \in \mathbb{R}^{1024}$

Representation with new truncated basis

original faces, $x \in \mathbb{R}^{1024}$

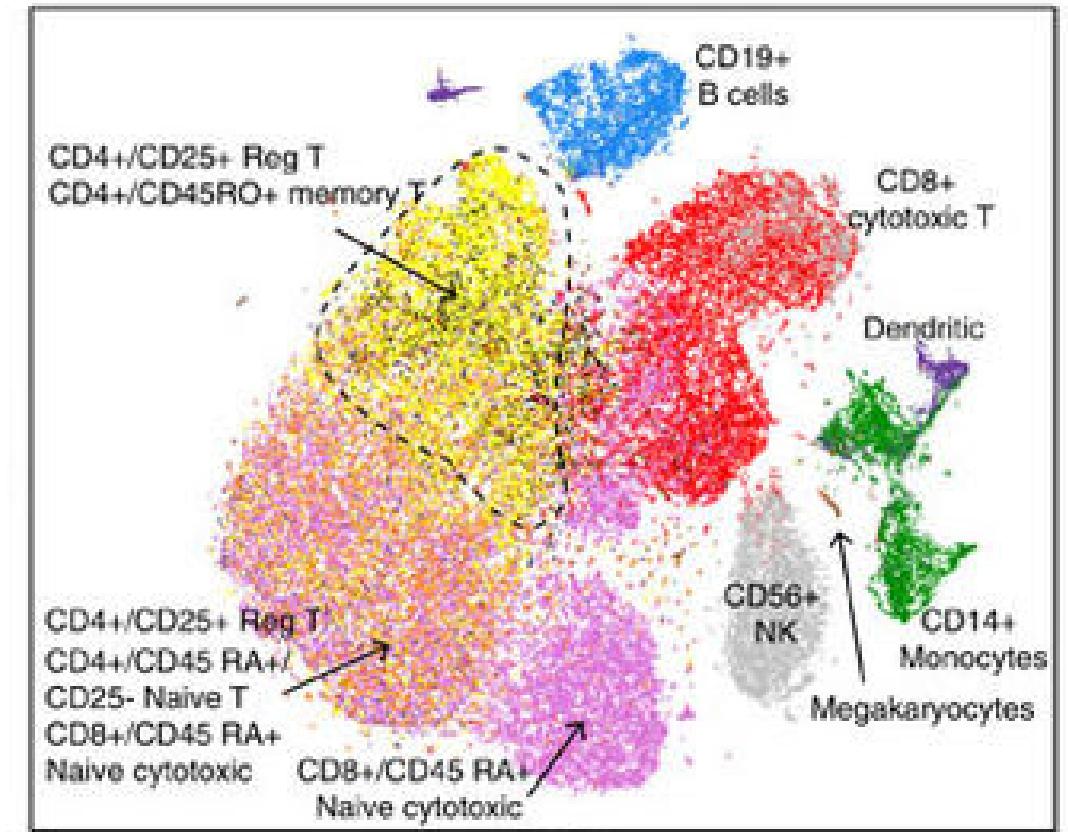


approximate faces, $x' \in \mathbb{R}^{100}$



Why might we want to reduce dimensionality?

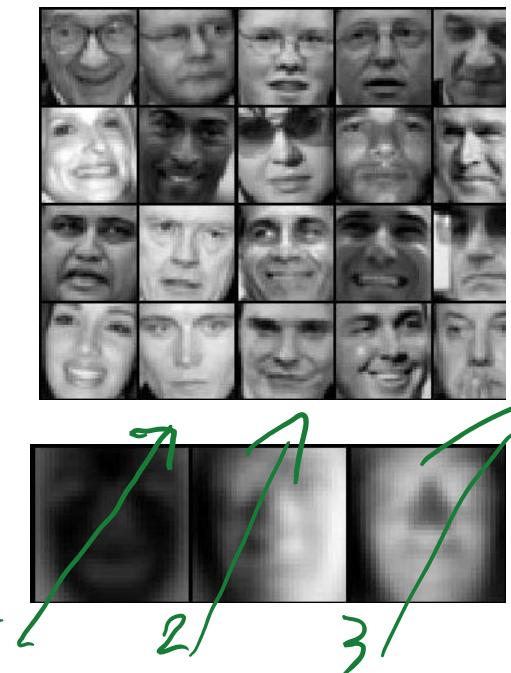
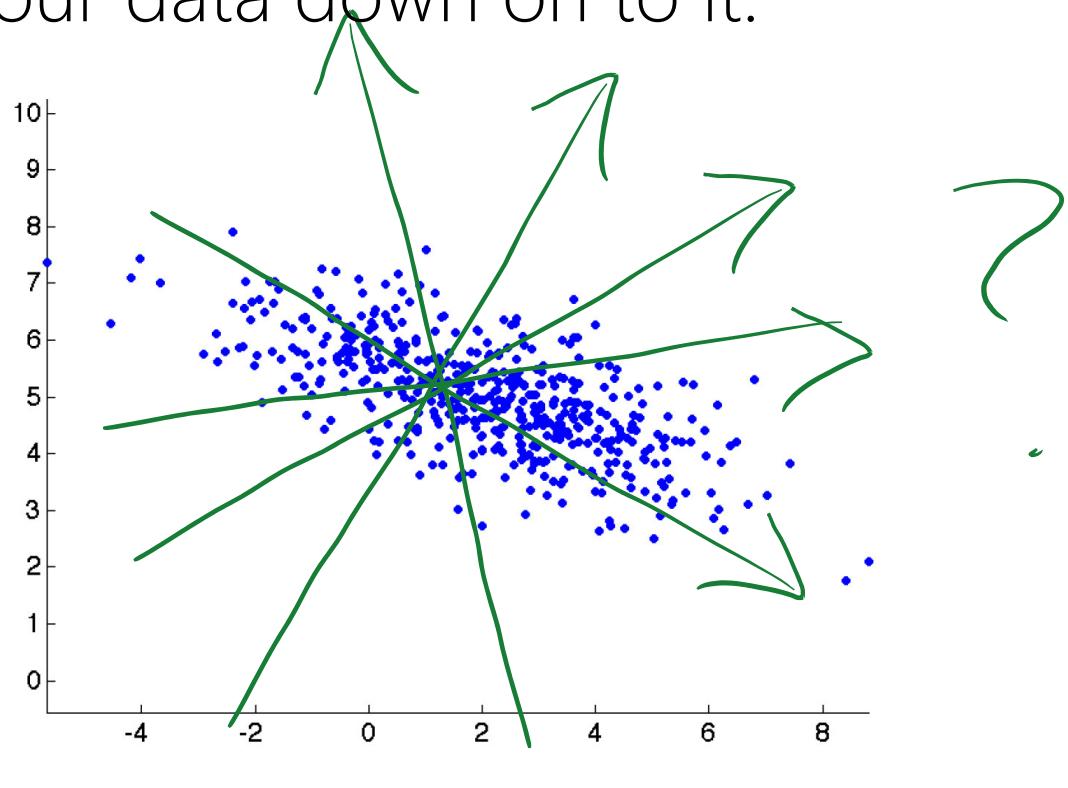
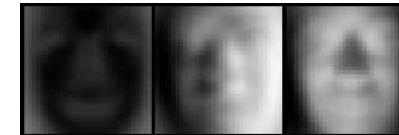
- Visualization, e.g. 2D plots.
- To denoise the data, or remove systematic artifacts (big one in biology).
- To compress the data (e.g. audio, images).
- To speed up supervised learning, or other analyses.



<https://www.nature.com/articles/ncomms14049/figures/3>

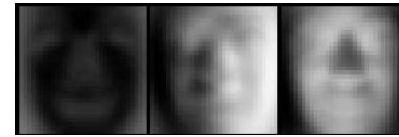
Principal Components Analysis (PCA)

- “eigenfaces” are the PCA basis
- Look for the **direction** in the original space that “retains most of the **“information”** if you project your data down on to it.

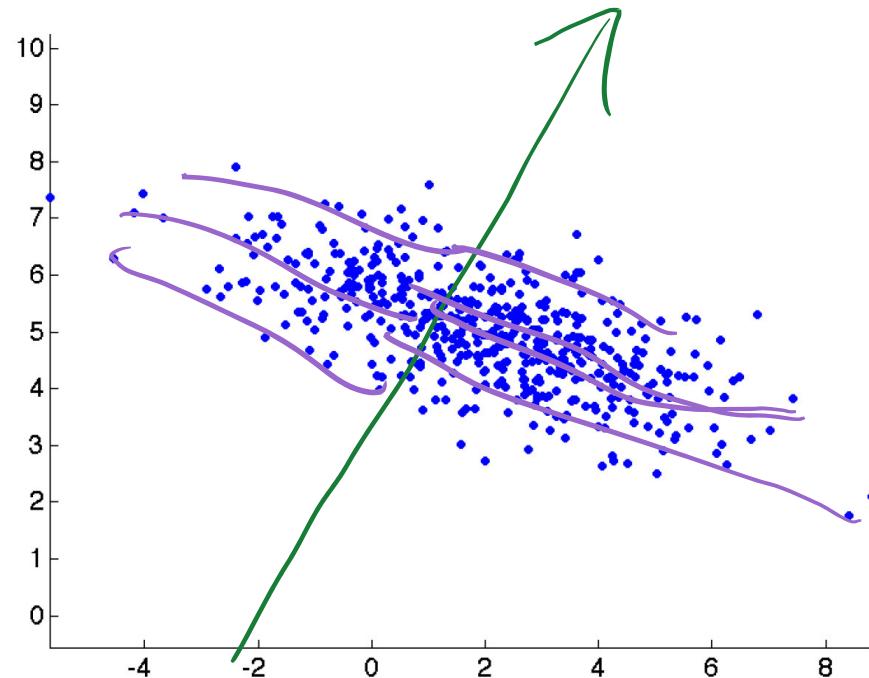


Principal Components Analysis (PCA)

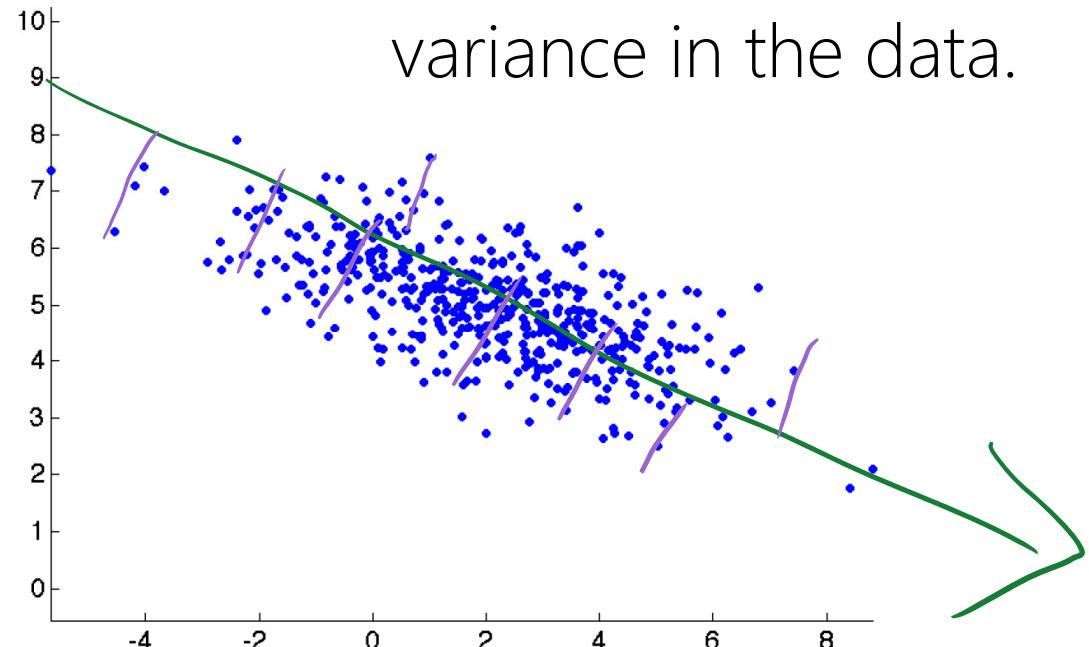
- “eigenfaces” are the PCA basis



- Look for the **direction** in the original space that “retains most of the “**information**” if you project your data down on to it.



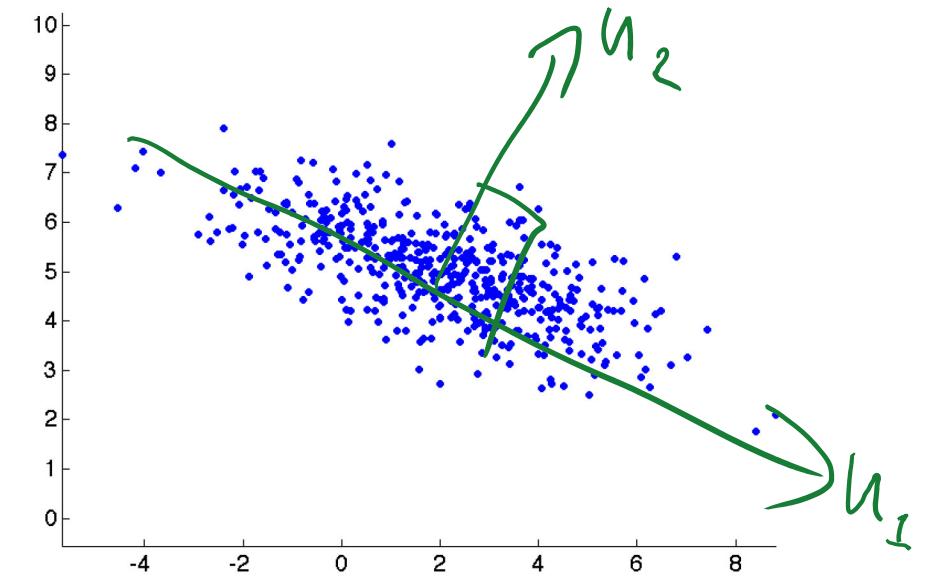
Direction with lowest *reconstruction loss*, is the direction with maximal variance in the data.



Principal Components Analysis (PCA)

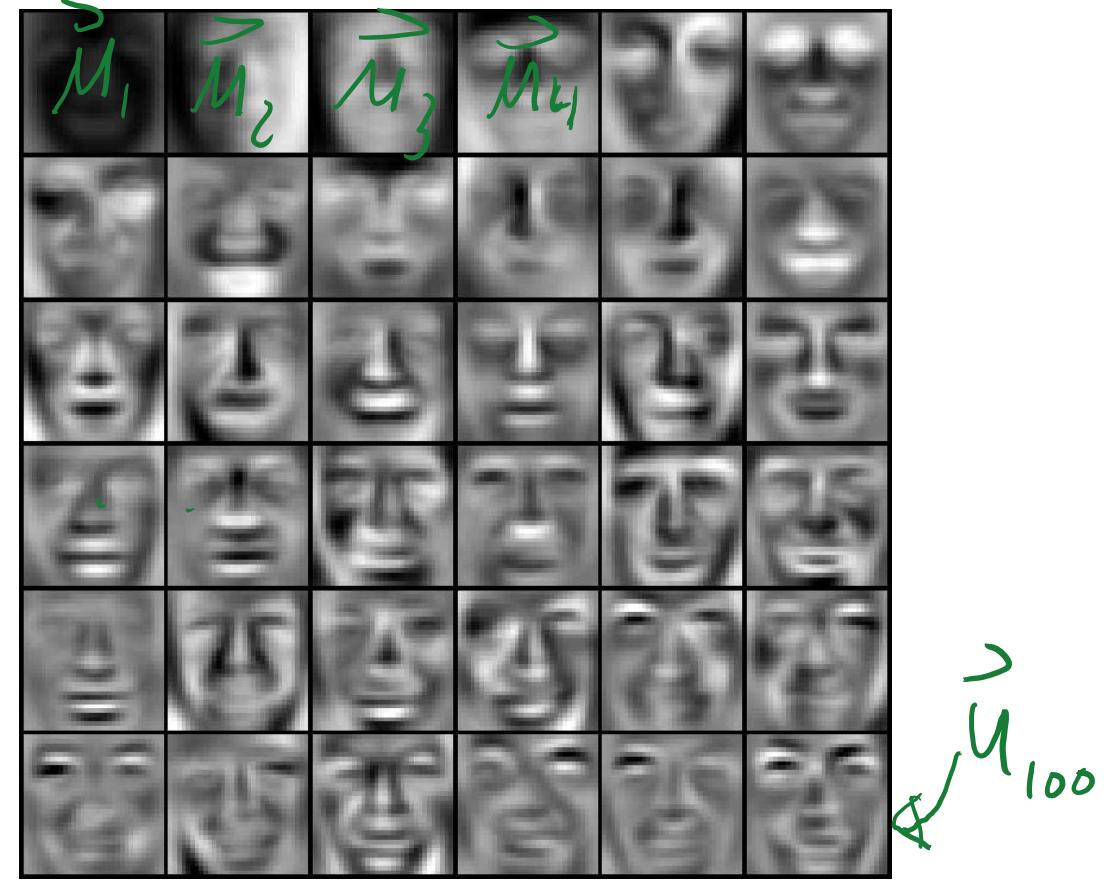
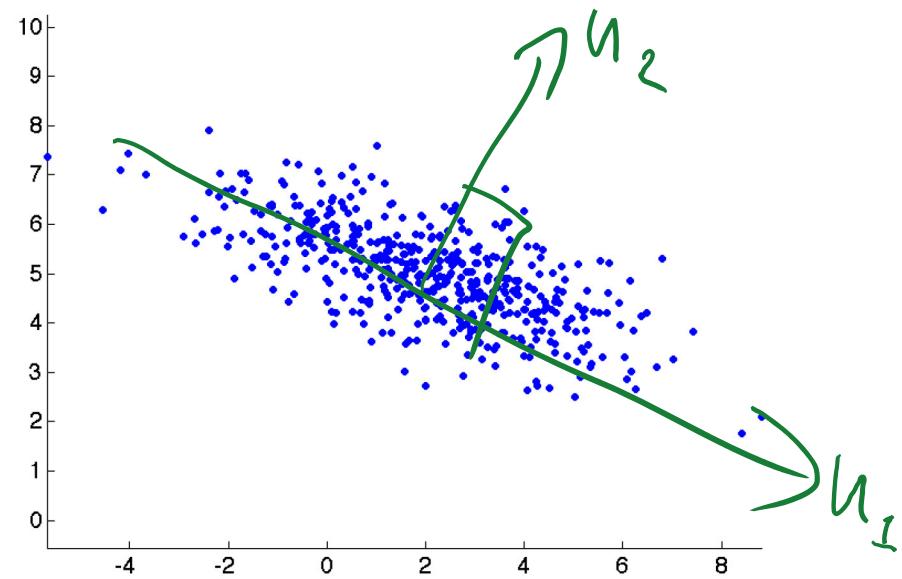
Recursively apply this idea to find 2nd best direction, then 3rd best:

- 2nd direction should be orthogonal to the first... and...
- ... be direction of most variance subject to that constraint.
- What's the maximum number of such directions we can find?



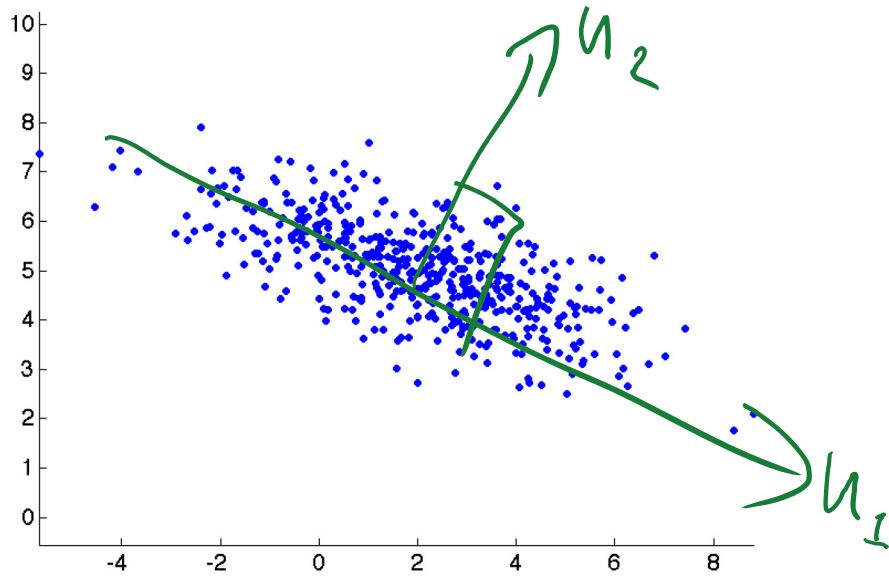
Principal Components Analysis (PCA)

- Each of the 100 eigenfaces was one of these special directions in the original 1024-dimensional image space.



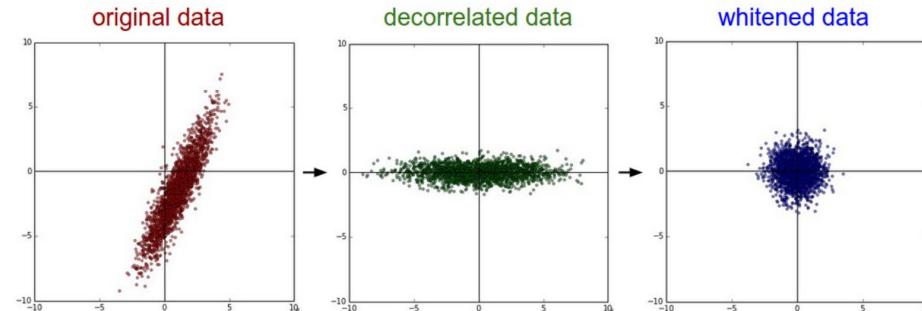
Principal Components Analysis (PCA)

Is this starting to remind you of anything?



Diagonalizing a MVG ("spherizing")

- To sphere a MVG is to make all its contour lines be spheres (also called "whitening").
- Thus we need to make the ellipses look like spheres.
- To do this, we need to understand how to diagonalize a matrix.



Can be derived (16.2.1-16.2.3) with MLE for params μ, W , assuming $p(x \in \mathbb{R}^d) = N(\mu + xW; I\sigma^2)$, for $\sigma^2 \rightarrow 0$ and $W \in \mathbb{R}^{d \times k}$.

Principal Components Analysis (PCA)

Recall: to diagonalize a MVG distribution, we made use of a special factorization of its covariance matrix, $A = QDQ^T$, an "eigen" or "spectral"-decomposition:

Linear Algebra: Diagonalizing a matrix

Spectral theorem:

When A is symmetric $A = A^T$

$A = QDQ^T$ with real eigenvalues in D

and orthonormal vectors in $S = Q$

Q is an orthonormal matrix

$[q_1 \ q_2 \dots \ q_n]$

each is
of length 1

any two
are orthogonal

true for col and row

$$\downarrow$$

$$Q^{-1} = Q^T$$

(rotations & reflections)

$$A = QDQ^{-1}$$

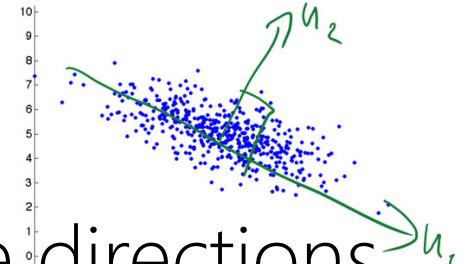
Can use this to do PCA.

(Even if data are not Gaussian).

PCA overview

Compute covariance matrix of the data and pick off the directions with the top k eigenvalues (all $\lambda_i \geq 0$ because covariance is PSD):

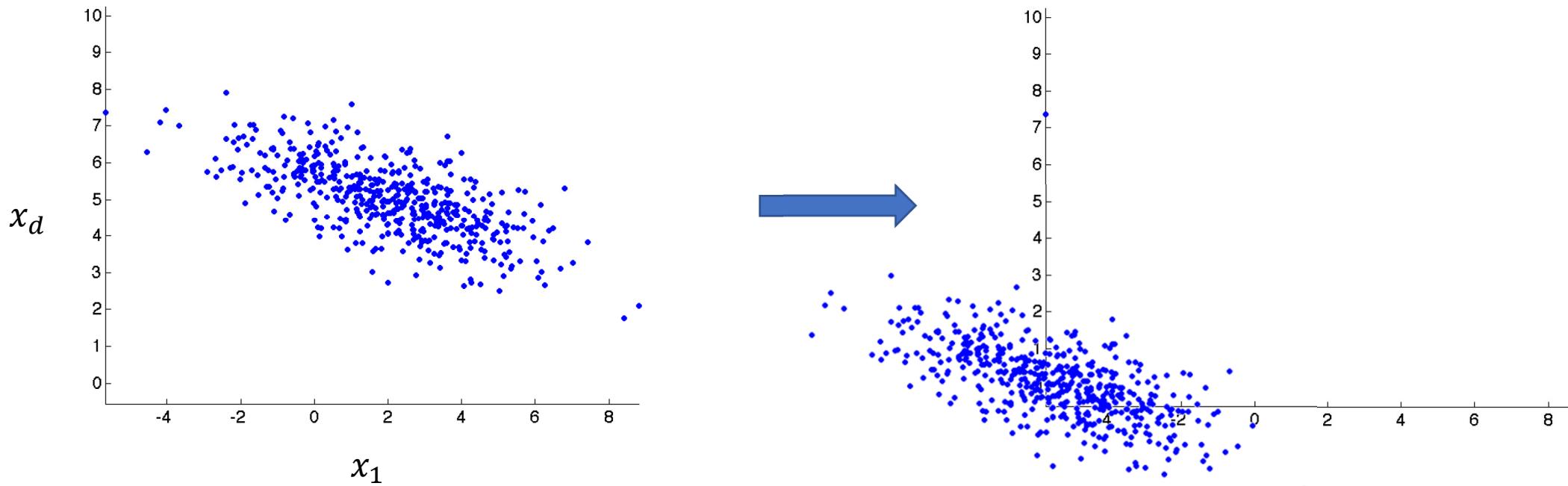
- Given data matrix, $X \in \mathbb{R}^{n \times d}$.
- Construct, $X^T X \in \mathbb{R}^{d \times d}$ (after mean-centering each feature).
- Apply spectral theorem, $X^T X = Q D Q^T$ to pick off k directions.
- Now approximate this covariance matrix with the “best” low rank approximation to it (rank k). (*Eckart-Young theorem*)
- Yields lowest “reconstruction loss”, and highest variance directions.



PCA step-by-step

Given n data points of dimension d , $X \in \mathbb{R}^{n \times d}$, to perform PCA, we:

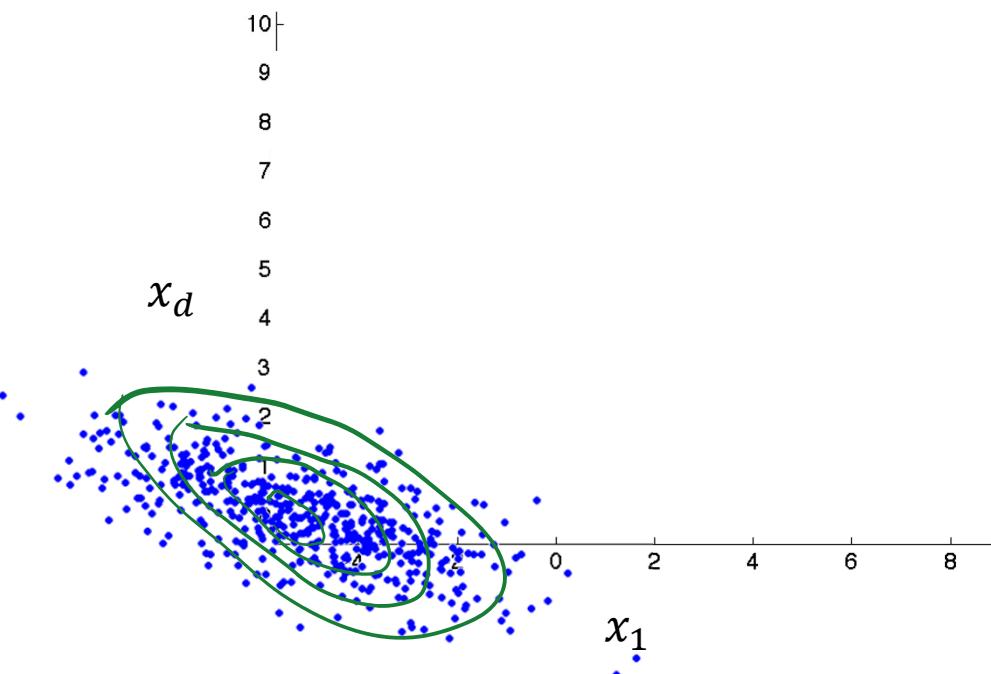
1. Center the data (make each feature zero mean), $\bar{X} = X - [\bar{x}_1, \dots, \bar{x}_d]$. (We will continue on only with the matrix \bar{X}).



PCA step-by-step

Given n data points of dimension d , $X \in \mathbb{R}^{n \times d}$, to perform PCA, we:

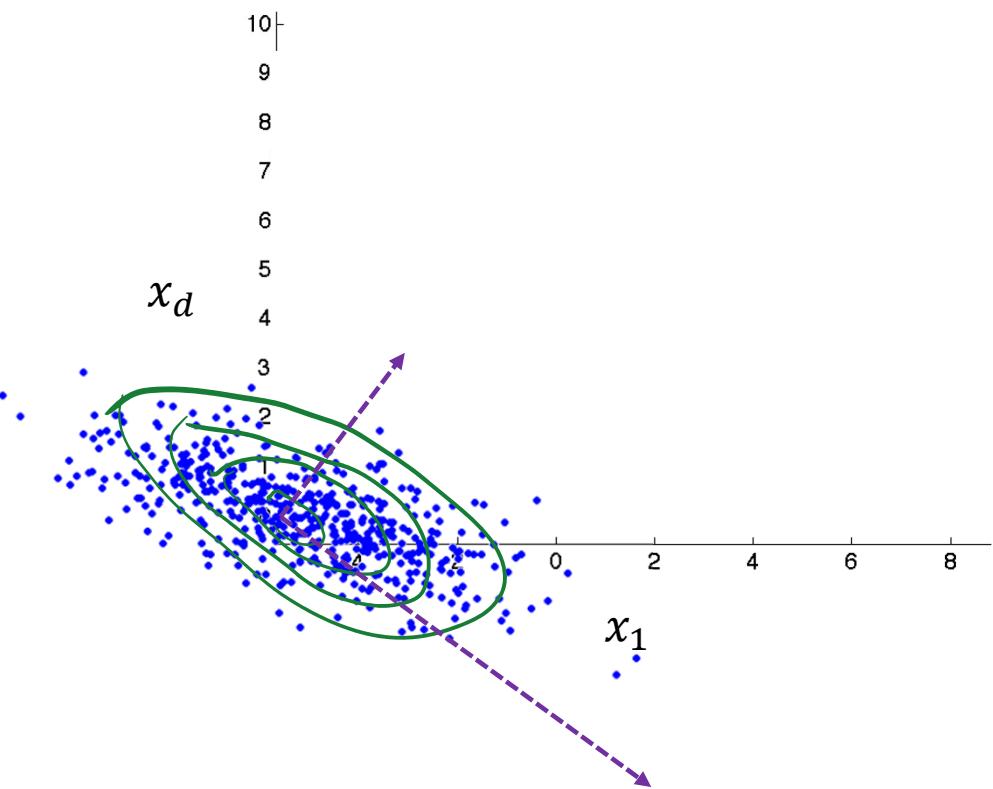
2. Compute the covariance matrix, $\Sigma = \bar{X}^T \bar{X}$.



PCA step-by-step

Given n data points of dimension d , $X \in \mathbb{R}^{n \times d}$, to perform PCA, we:

2. Compute the covariance matrix, $\Sigma = \bar{X}^T \bar{X}$.
3. Compute $\bar{X}^T \bar{X} = Q D Q^T$ to get eigenvectors (aka *principal directions*).

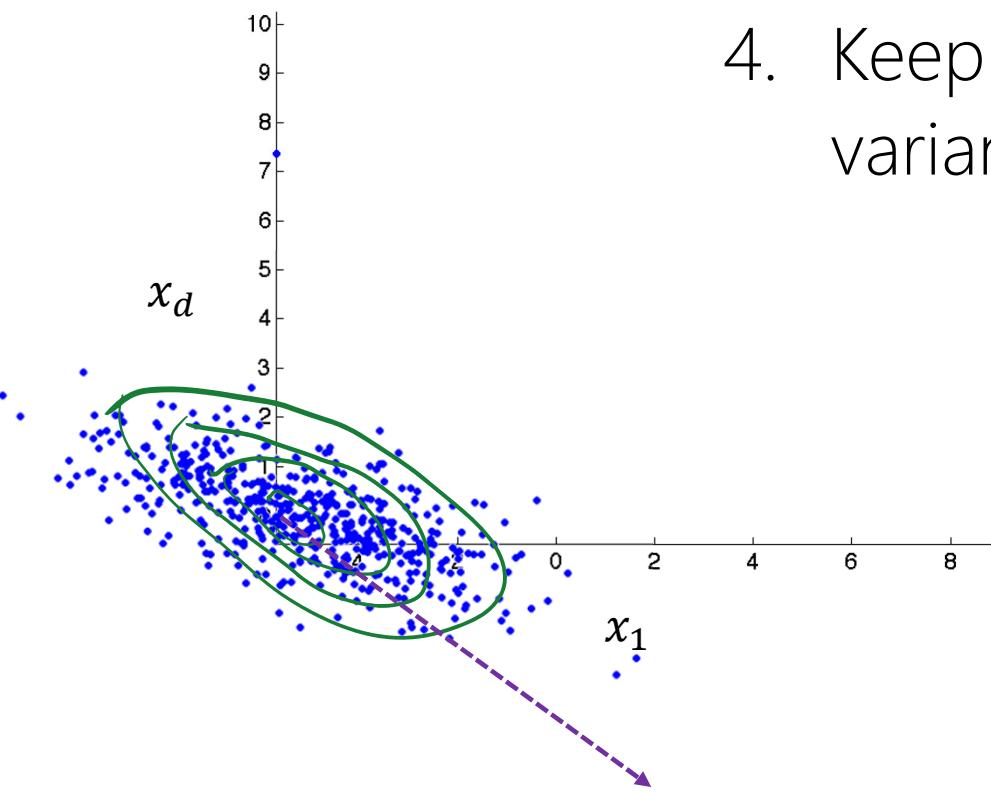


- This decomposition is typically implemented with a call to "eig" function (linalg package in python), but can also be obtained from "svd".
- Are the same for PSD matrices, but svd may be more stable.

PCA step-by-step

Given n data points of dimension d , $X \in \mathbb{R}^{n \times d}$, to perform PCA, we:

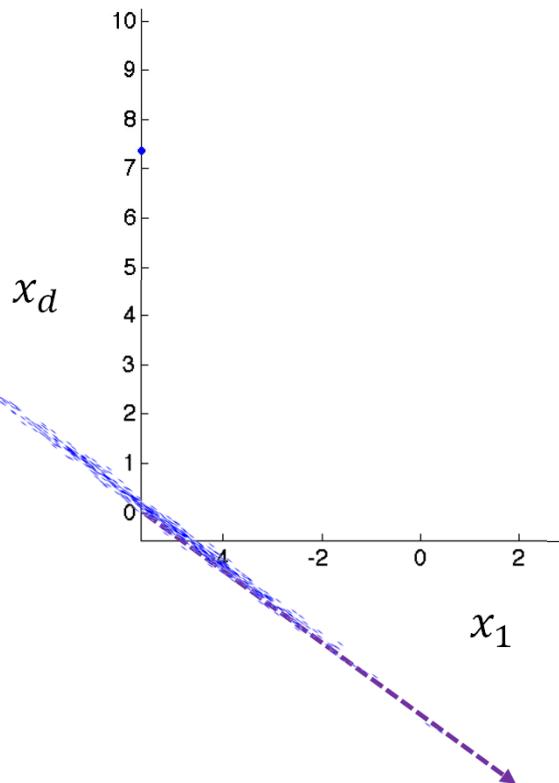
2. Compute the covariance matrix, $\Sigma = \bar{X}^T \bar{X}$.
3. Compute $\bar{X}^T \bar{X} = Q D Q^T$ to get eigenvectors (aka *principal component axes*).
4. Keep the k eigenvectors, $Q_k \equiv Q_{:,1:k}$, with the most variance (highest eigenvalues in D).



PCA step-by-step

Given n data points of dimension d , $X \in \mathbb{R}^{n \times d}$, to perform PCA, we:

2. Compute the covariance matrix, $\Sigma = \bar{X}^T \bar{X}$.
3. Compute $\bar{X}^T \bar{X} = Q D Q^T$ to get eigenvectors (aka *principal component axes*).
4.
 4. Keep the k eigenvectors, $Q_k \equiv Q_{:,1:k} \in \mathbb{R}^{n \times k}$, with the most variance (highest eigenvalues in D).
 5. Project your points (original, or new ones) down to this subspace, $\bar{X}_k = \bar{X} Q_k \in \mathbb{R}^{n \times k}$, these are your *principal component scores*.

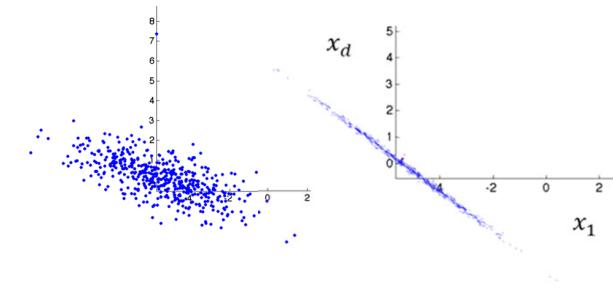


Final, dimensionality-reduced data is.

PCA step-by-step

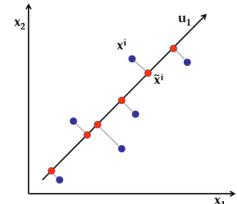
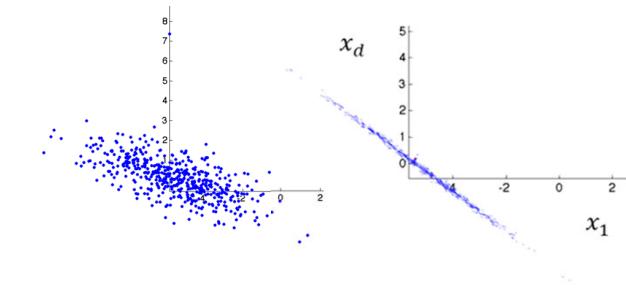
- Final, dimensionality-reduced data is $\bar{X}_k = \bar{X}Q_k \in \mathbb{R}^{n \times k}$.
- What if we wanted our data in the original dimension (d), but with only the information retained from our PCA- k analysis? (e.g. reconstructed faces)
- We need to “reconstruct” the original points. We can do this by [expanding back from PCA basis to original basis](#), but with only the first k PCA basis vectors, $\bar{X}_{recon-k} = \bar{X}_k Q_k^T = \bar{X} Q_k Q_k^T \in \mathbb{R}^{n \times d}$

$$\begin{matrix} n \times k & k \times d \\ \text{---} & \text{---} \\ d \times d & \end{matrix} = \begin{matrix} d \times k \\ \text{---} \\ \text{rank } k \times d \end{matrix}$$



PCA step-by-step

- Final, dimensionality-reduced data is $\bar{X}_k = \bar{X}Q_k \in \mathbb{R}^{n \times k}$.
- What if we wanted our data in the original dimension (d), but with only the information retained from our PCA- k analysis? (e.g. reconstructed faces)
- We need to “reconstruct” the original points. We can do this by [expanding back from PCA basis to original basis](#), but with only the first k PCA basis vectors, $\bar{X}_{recon-k} = \bar{X}_k Q_k^T = \bar{X}Q_k Q_k^T \in \mathbb{R}^{n \times d}$
- Now we talk about the “reconstruction loss”, as the difference between original and reconstructed data, $\|\bar{X}_{recon-k} - \bar{X}\|_F$.
- The larger the reconstruction loss, the more information we have lost.



and Frobenius norm as

$$\|X\|_F = \sqrt{\sum_{ij} X_{ij}^2} = \sqrt{\text{tr}(X^\top X)} = \sqrt{\sum s_i^2},$$

where s_i are singular values of X , i.e. diagonal

Original, and reconstructed faces

original faces, $x \in \mathbb{R}^{1024}$

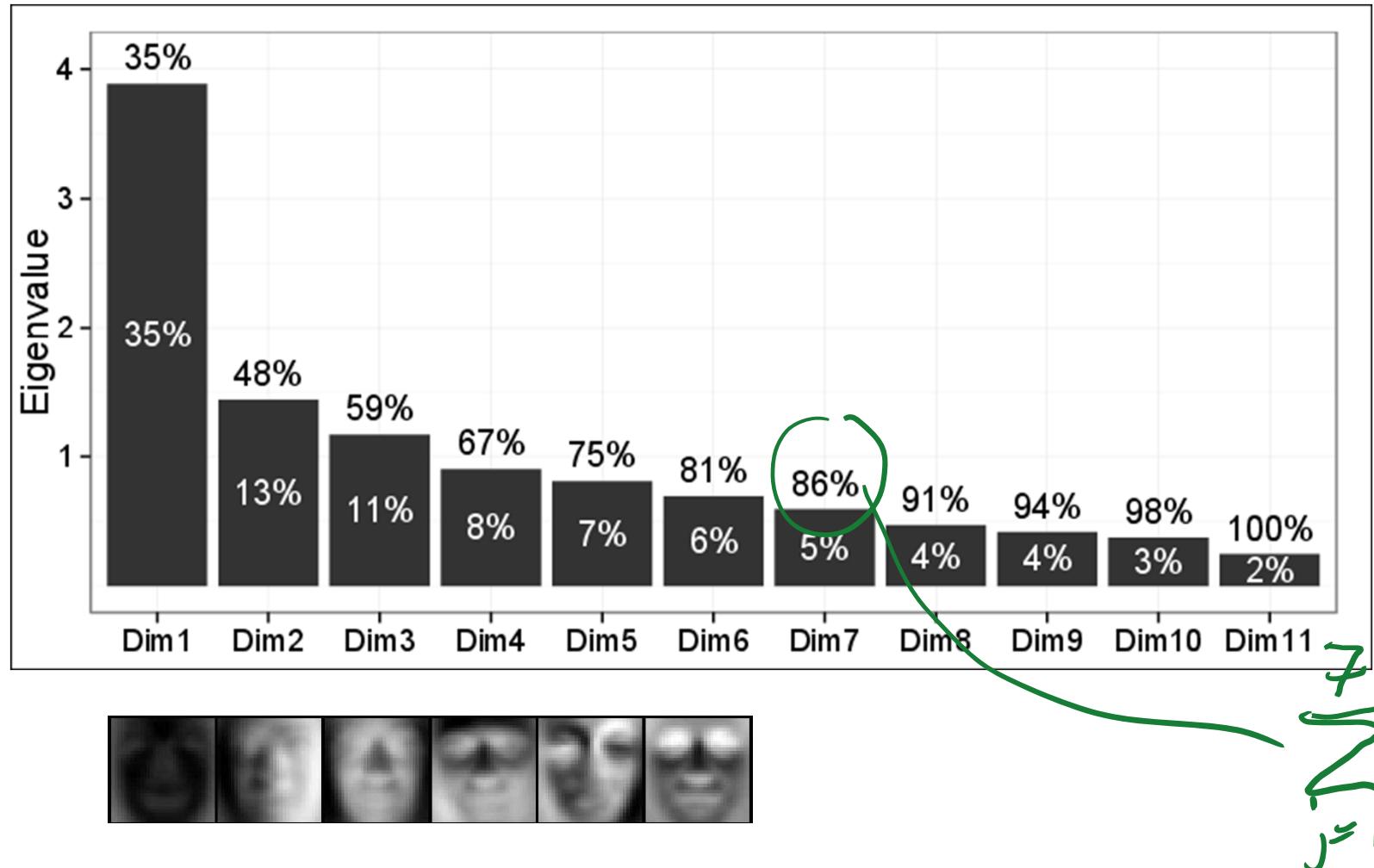


reconstructed faces

$$x' \in \mathbb{R}^{100} \rightarrow x_k \in \mathbb{R}^{1024}$$



PCA: % variance explained to help pick hyper-parameter, k



$$\mathbf{X}^T \mathbf{X} = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$$

If you normalize the eigenvalues in \mathbf{D} by their sum, then they correspond to % variance explained.

$$V_i \equiv \frac{D_{i,i}}{\sum D_{j,j}}$$

$$V_j \quad j=1$$

PCA: reconstruction loss and % variance

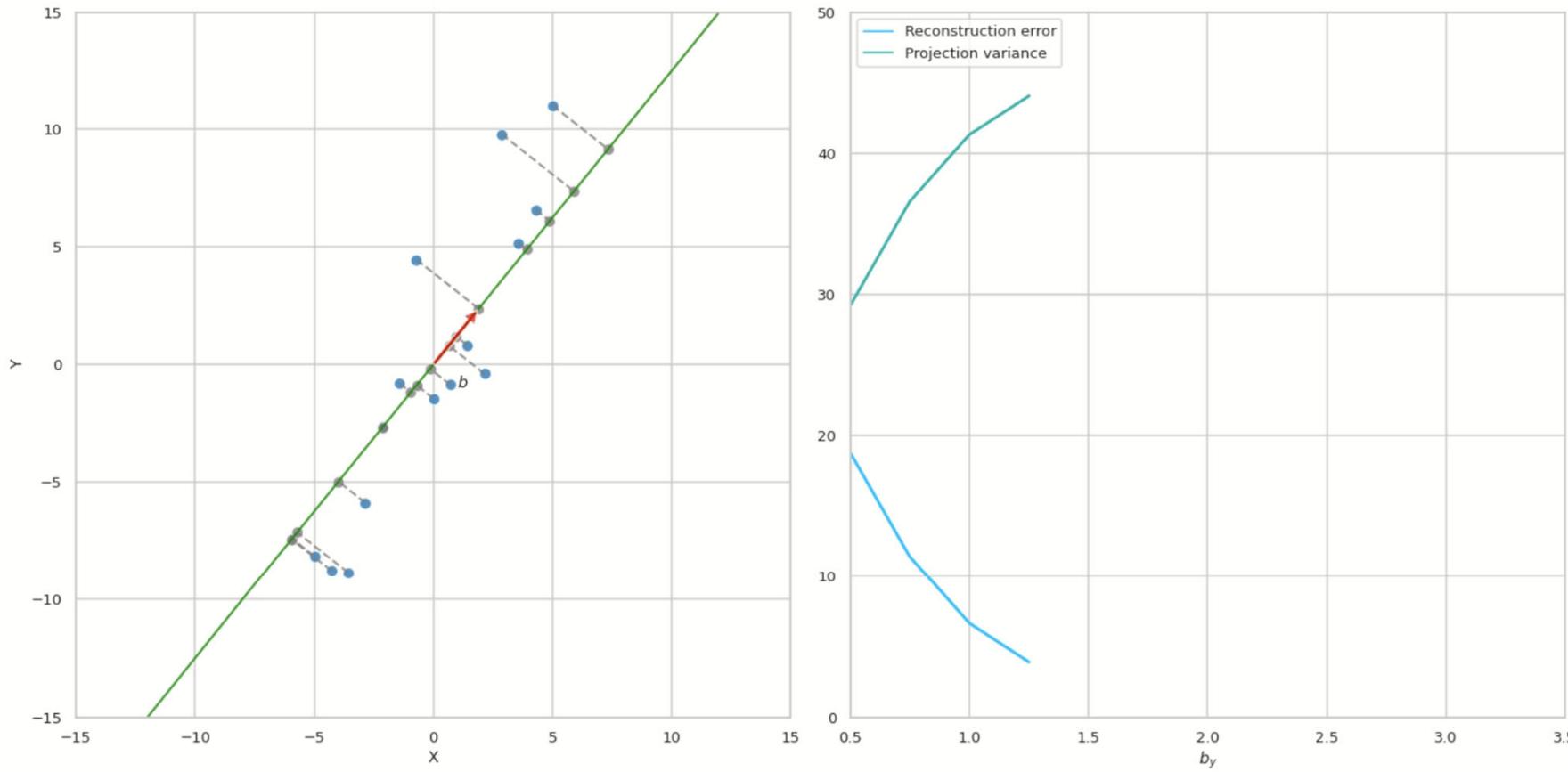
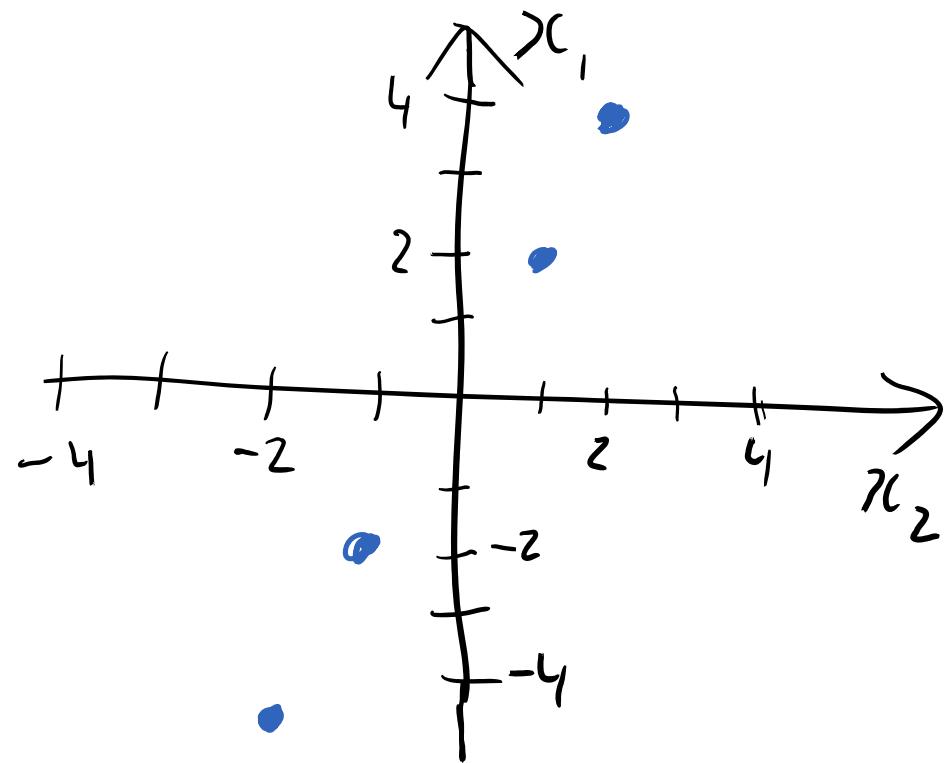


Fig.5. We rotate basis vector b and project data on corresponding subspace, then calculate reconstruction error L and projection variance $\text{Var}(\lambda)$.

Example of PCA

$$X = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ -1 & -2 \\ -2 & -4 \end{bmatrix} \quad n=4 \quad d=2$$

$x_1 \quad x_2$



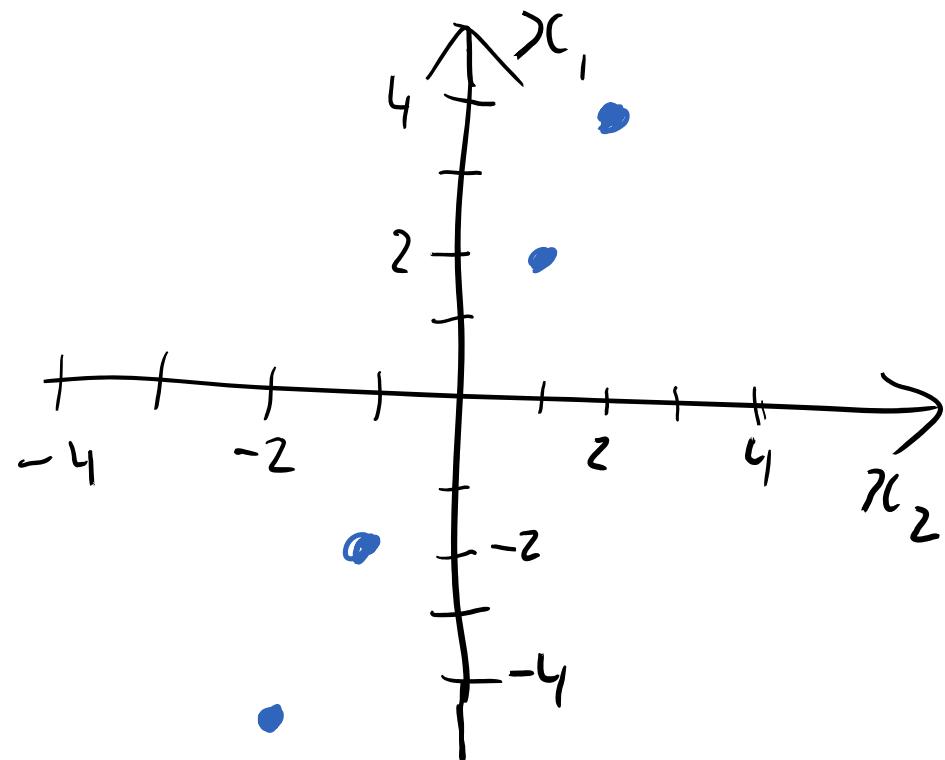
Example of PCA

$$X = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ -1 & -2 \\ -2 & -4 \end{bmatrix} \quad n=4 \quad d=2$$

$x_1 \quad x_2$

mean
of
each
dimension

$$\begin{bmatrix} 0 & 0 \end{bmatrix} \Rightarrow \bar{X} = X$$

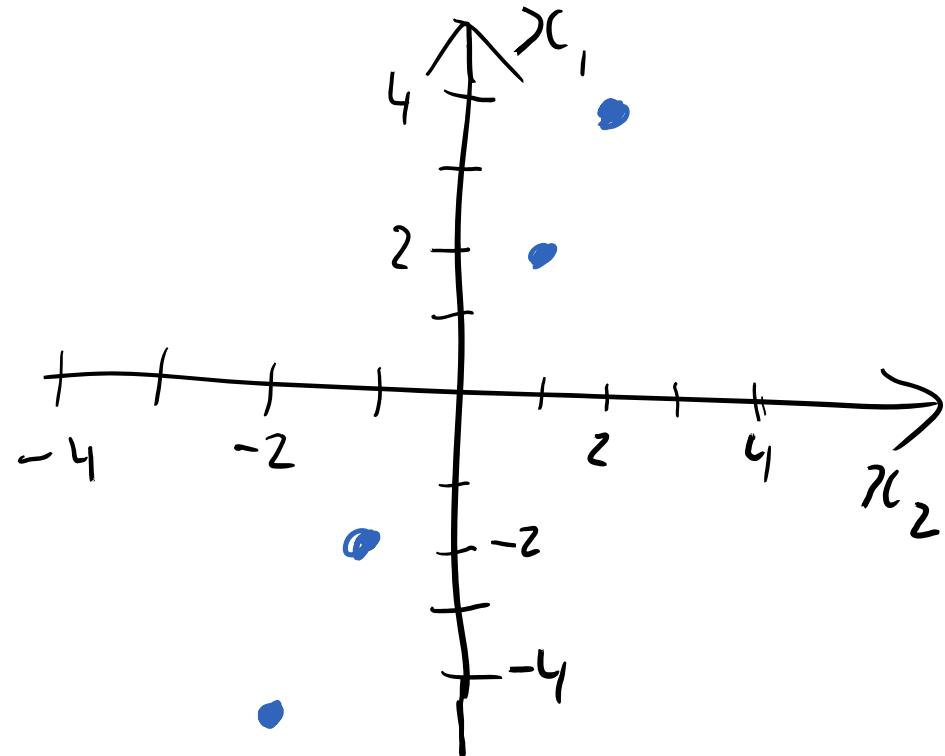


Example of PCA

$$X = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ -1 & -2 \\ -2 & -4 \end{bmatrix} \quad d=2 \quad n=4$$

$x_1 \quad x_2$

$$X^T X = \begin{bmatrix} 1 & 2 & -1 & -2 \\ 2 & 4 & -2 & -4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ -1 & -2 \\ -2 & -4 \end{bmatrix} = \begin{bmatrix} 10 & 20 \\ 20 & 40 \end{bmatrix}$$



Example of PCA

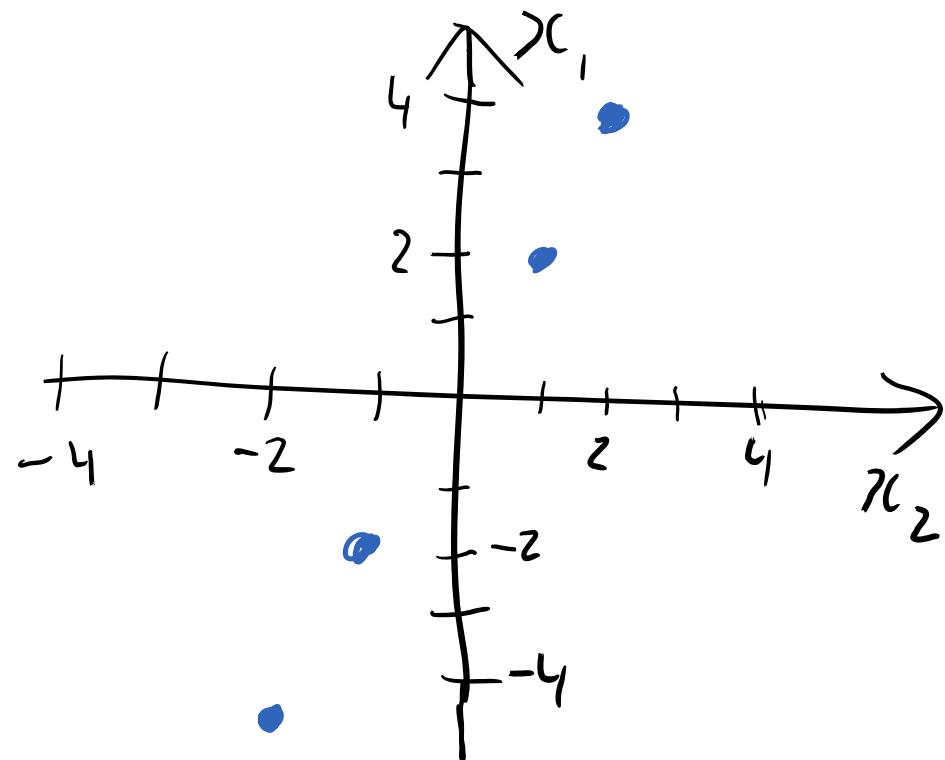
$$X = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ -1 & -2 \\ -2 & -4 \end{bmatrix}$$

$n \times d$

$x_1 \quad x_2$

$$X^T X = \begin{bmatrix} 10 & 20 \\ 20 & 40 \end{bmatrix}$$

$d \times d$



Example of PCA

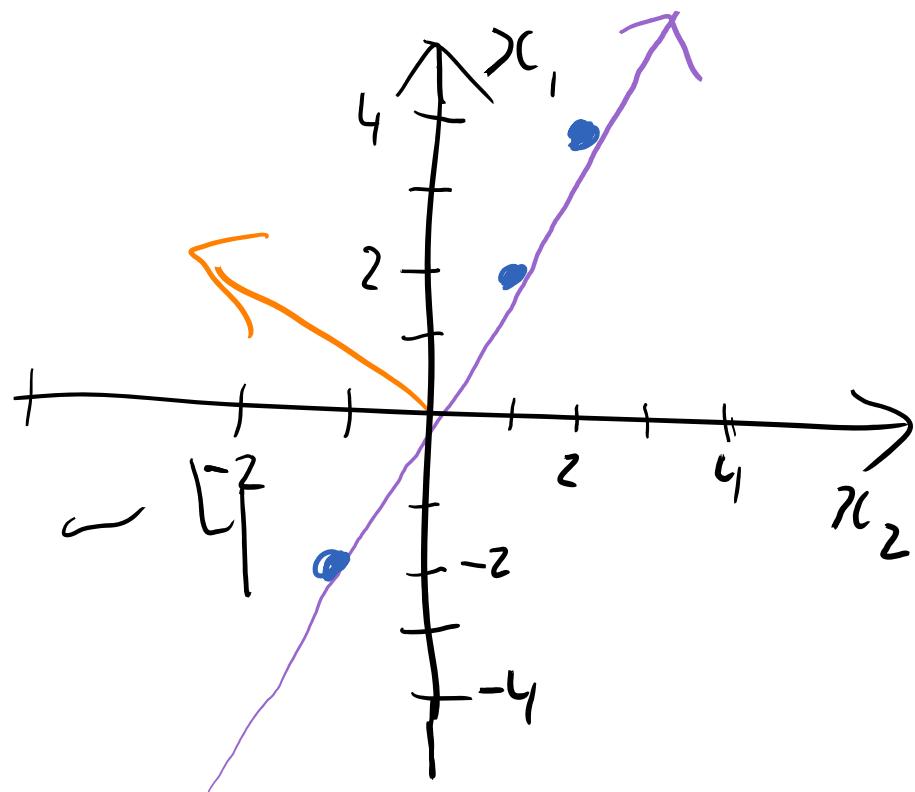
$$X = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ -1 & -2 \\ -2 & -4 \end{bmatrix} \quad \text{linalg.eig}(X^T X)$$

$n \times d$

$x_1 \quad x_2$

\Rightarrow eigenvalues [50] [0]

$$X^T X = \begin{bmatrix} 10 & 20 \\ 20 & 40 \end{bmatrix} \quad d \times d$$



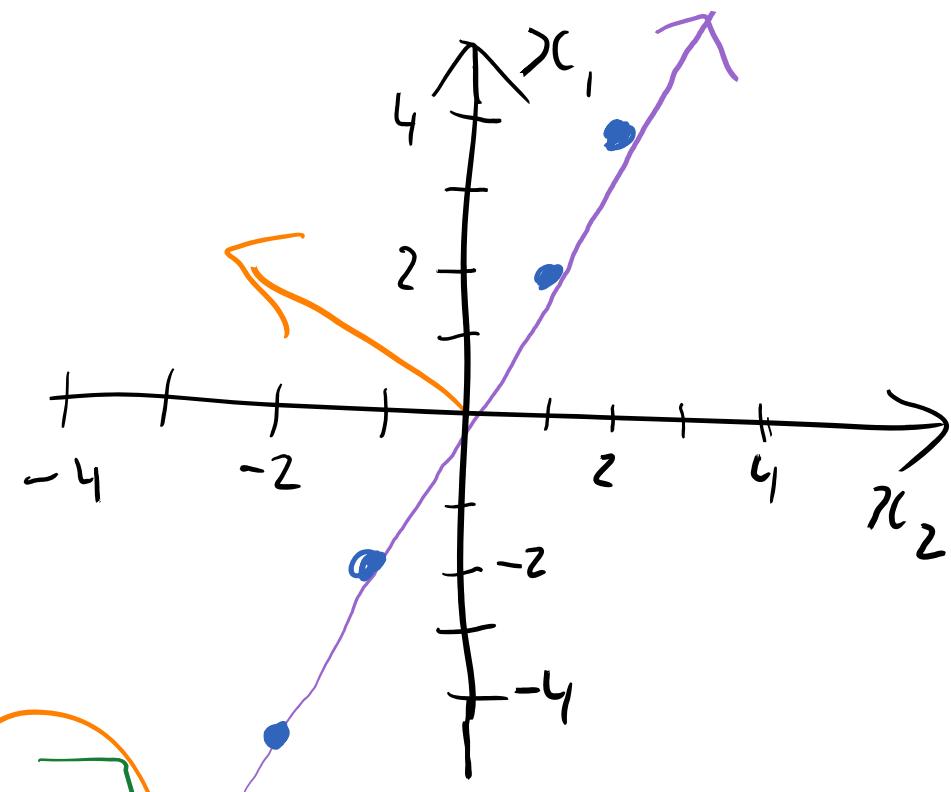
Example of PCA

$X = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ -1 & -2 \\ -2 & -4 \end{bmatrix}$ $\text{linalg.eig}(X^T X)$
 $n \times d$ \Rightarrow eigenvalues [50] (purple) [0] (orange)
\$x_1\$ \$x_2\$
& eigen vectors

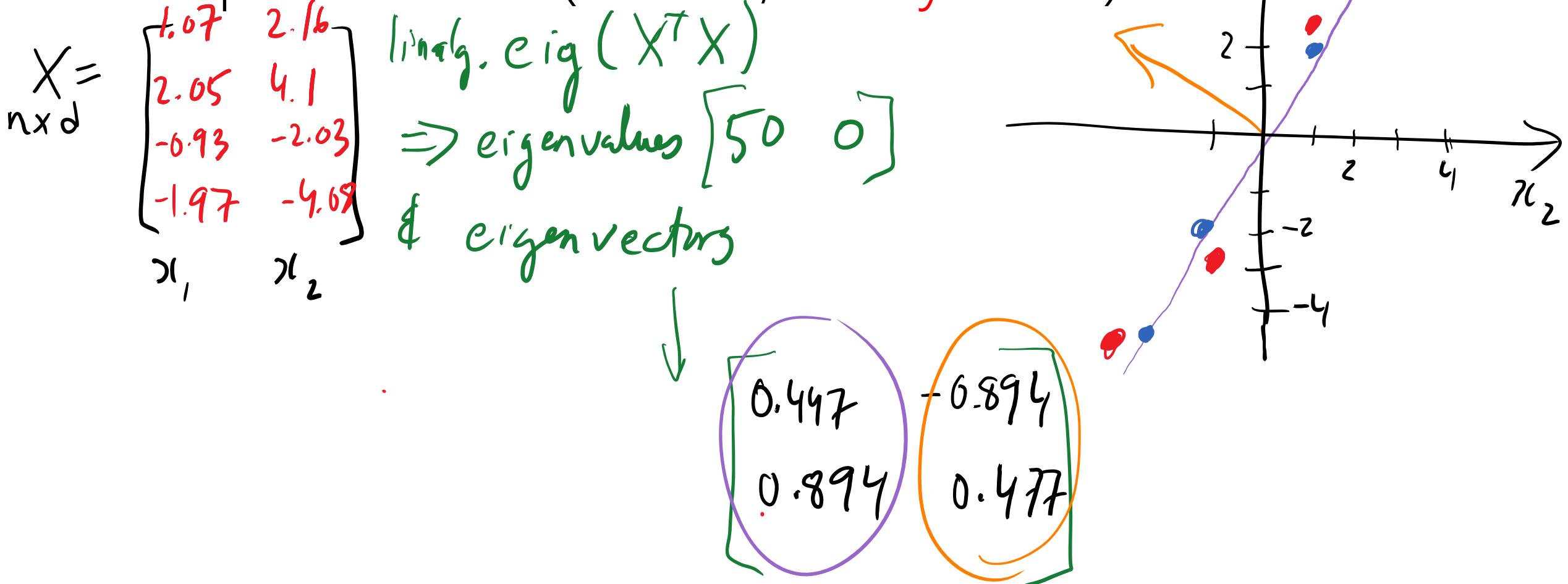
$X^T X = \begin{bmatrix} 10 & 20 \\ 20 & 40 \end{bmatrix}$

↓

0.447	-0.894
0.894	0.477



Example of PCA (take 2, noisy data)



Example of PCA (take 2, noisy data)

$X =$
 $n \times d$

$$\begin{bmatrix} 1.07 & 2.16 \\ 2.05 & 4.1 \\ -0.93 & -2.03 \\ -1.97 & -4.69 \end{bmatrix}$$

$x_1 \quad x_2$

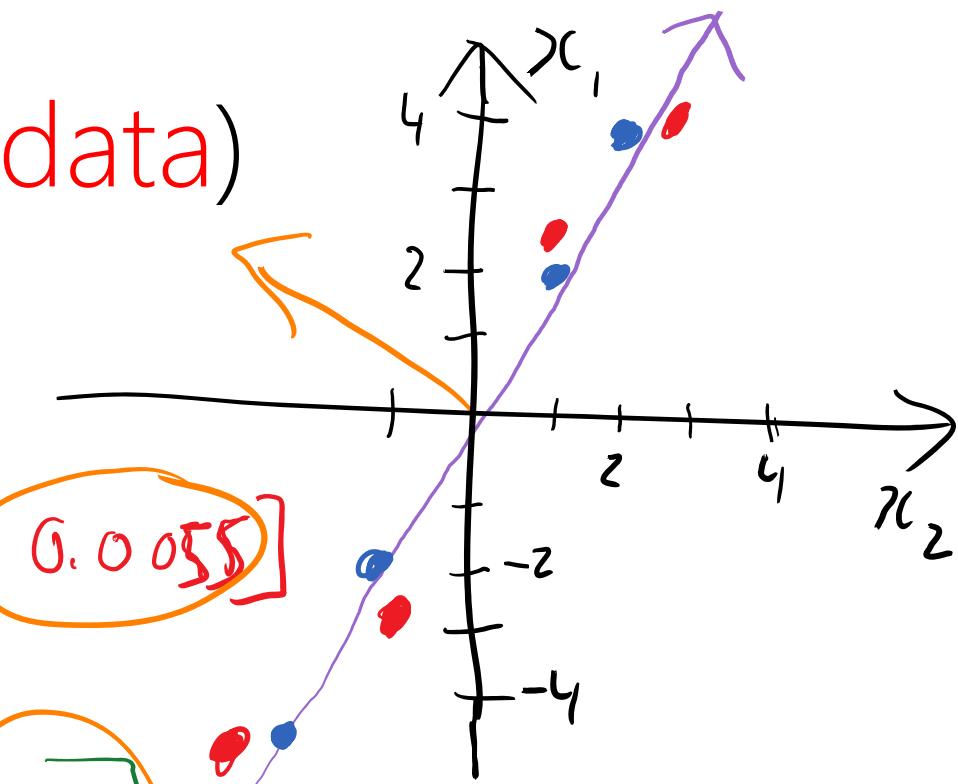
linalg.eig($X^T X$)

\Rightarrow eigenvalues $[50 \quad 0]$
\$ \& \$ eigen vectors $\begin{bmatrix} 52.34 \\ 0.0055 \end{bmatrix}$

\downarrow

$$\begin{bmatrix} 0.439 \\ 0.898 \end{bmatrix} \quad \begin{bmatrix} -0.898 \\ 0.439 \end{bmatrix}$$

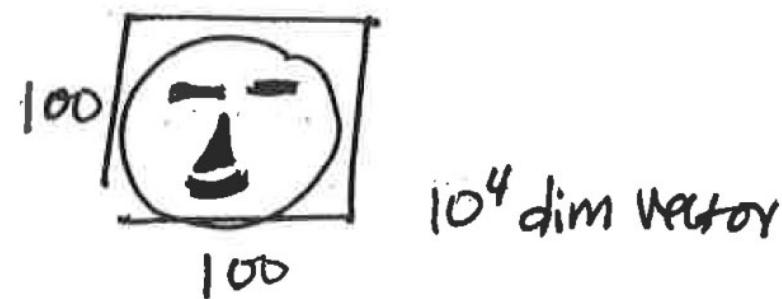
$$\begin{bmatrix} 0.447 \\ 0.894 \\ -6.894 \\ 0.477 \end{bmatrix}$$



Linear algebra for PCA

$$\bar{X} \in \mathbb{R}^{n \times d}$$

- Need to compute the principal axes, Q , of $\bar{X}^T \bar{X} = Q D Q^T \in \mathbb{R}^{d \times d}$.
- This is an *eigendecomposition* of the matrix $\bar{X}^T \bar{X}$.
- Computing its eigendecomposition has time complexity $O(d^3)$.
- What if $d \gg n$? e.g., images of $d = 100 \times 100 = 10^4$ pixels, and $n = 1,000$ images.
- Could we do something cheaper?
- Yes. Need to understand the SVD.



$$\bar{X} \in \mathbb{R}^{n \times d}$$

Linear algebra for PCA

Recall the spectral theorem (principal axis theorem) from MVG lecture, which gives a *spectral (eigen) decomposition*:

When A is symmetric $A = A^T$

$A = Q D Q^T$ with real eigenvalues in D and orthonormal vectors in $S = Q$

Q is an orthonormal matrix $[q_1 \ q_2 \dots q_n]$

\Downarrow

$Q^{-1} = Q^T$
(rotations & reflections)

each is of length 1
any two are orthogonal
true for col and row

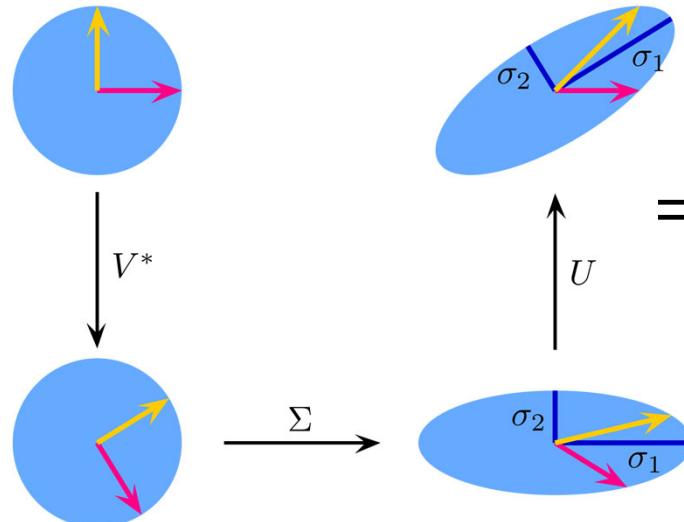
$A = Q D Q^{-1}$

- The covariance matrix for PCA, $\bar{X}^T \bar{X}$, is symmetric (and PSD).
- It turns out, there is a generalization of the spectral decomposition, for non-symmetric and non-square matrices, the SVD that will be helpful.

Singular Value Decomposition (SVD)

Can think of M as linear transformation broken down into three steps, by looking at its effect on the unit disc and the two canonical unit vectors e_1 and e_2 :

1. **Left:** V^T rotates the disc and unit vectors.
2. **Bottom:** Σ stretches scales axes by $\sigma_i = \Sigma_{i,i}$ (singular values).
3. **Right:** U performs another rotation.



$$M = U \cdot \Sigma \cdot V^*$$

Can be applied to any matrix M .

$$\begin{matrix} \text{Matrix } M \\ m \times n \end{matrix} = \begin{matrix} \text{Orthogonal Matrix } U \\ m \times m \end{matrix} \begin{matrix} \text{Diagonal Matrix } \Sigma \\ m \times n \end{matrix} \begin{matrix} \text{Orthogonal Matrix } V^* \\ n \times n \end{matrix}$$

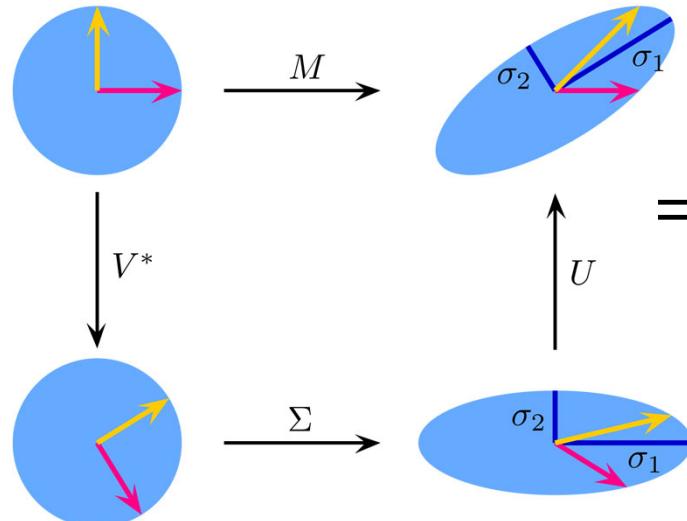
$$\begin{aligned} Mv_1 &= \Sigma_{1,1}u_1 \\ Mv_2 &= \Sigma_{2,2}u_2 \\ &\dots \\ Mv_r &= \Sigma_{r,r}u_r \end{aligned}$$

$$\begin{aligned} r &= \text{rank}(M) \\ &\leq \min(m, n) \end{aligned}$$

Singular Value Decomposition (SVD)

Can think of M as linear transformation broken down into three steps, by looking at its effect on the unit disc and the two canonical unit vectors e_1 and e_2 :

1. **Left:** V^T rotates the disc and unit vectors.
2. **Bottom:** Σ stretches scales axes by $\sigma_i = \Sigma_{i,i}$ (singular values).
3. **Right:** U performs another rotation.



$$M = U \cdot \Sigma \cdot V^*$$

Can be applied to any matrix M .

$$\begin{array}{cccc} \text{M} & = & \mathbf{U} & \Sigma \\ m \times n & & m \times m & m \times n \\ \mathbf{U} & & \mathbf{U}^* & = \mathbf{I}_m \\ & & \mathbf{V} & \mathbf{V}^* = \mathbf{I}_n \end{array}$$

- Has time complexity $O(\min(m^2n, mn^2))$.
- Σ is unique (if in descending order), but V and U are generally not: e.g. sign flips.
- (Eigendecomposition is unique if all eigenvalues are unique)
- If M is square+symmetric, yields the spectral decomposition.

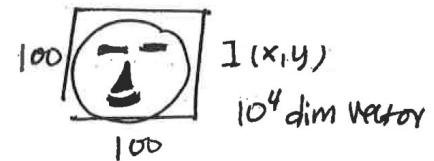
Singular Value Decomposition (SVD)

- Columns in \mathbf{U} are the eigenvectors of $\mathbf{M}\mathbf{M}^T$, called the *left singular vectors* of \mathbf{M} ($\mathbf{M}\mathbf{M}^T = \mathbf{U}\Sigma\mathbf{V}^T\mathbf{V}\Sigma^T\mathbf{U}^T = \mathbf{U}\Sigma^2\mathbf{U}^T$).
 - Columns in \mathbf{V} are the eigenvectors of $\mathbf{M}^T\mathbf{M}$, called the *right singular vectors* of \mathbf{M} ($\mathbf{M}^T\mathbf{M} = \mathbf{V}\Sigma^T\mathbf{U}^T\mathbf{U}\Sigma\mathbf{V}^T = \mathbf{V}\Sigma^2\mathbf{V}^T$).
 - Both spectral decompositions at once!
 - Eigenvalues are the same, given by $\lambda_i = \Sigma_{i,i}^2$ ($\Sigma_{i,i}$ are the *singular values* of \mathbf{M}).
- Can be applied to any matrix \mathbf{M} .

$$\mathbf{M}_{m \times n} = \mathbf{U}_{m \times m} \Sigma_{m \times n} \mathbf{V}^*_{n \times n}$$
$$\mathbf{U} \quad \mathbf{U}^* = \mathbf{I}_m$$
$$\mathbf{V} \quad \mathbf{V}^* = \mathbf{I}_n$$

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

Singular Value Decomposition (SVD)



- Recall this example with $d \gg n$, e.g. $d = 10^4$ pixels, $n = 1000$ images.
- How can we make use of what we just learned to do PCA faster than the eigendecomposition $O(d^3)$?
- Instead of spectral decomposition of $X^T X$...
- ...directly use SVD of the data matrix: $SVD(X) = U\Sigma V^T$
- This SVD has time complexity $O(dn^2)$.
- $V \in \mathbb{R}^{d \times d}$ are the needed eigenvectors for $X^T X \in \mathbb{R}^{d \times d}$.
- $\lambda_i = \Sigma_{i,i}^2$ are needed eigenvalues.

$$\begin{matrix} M \\ m \times n \end{matrix} = \begin{matrix} U \\ m \times m \end{matrix} \begin{matrix} \Sigma \\ m \times n \end{matrix} \begin{matrix} V^* \\ n \times n \end{matrix}$$

Singular Value Decomposition (SVD)

$$\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$$

$m \times n$ $m \times m$ $m \times n$ $n \times n$

For PCA we want projections onto top k PCs.

- When we used a spectral decomposition, $\mathbf{X}^T \mathbf{X} = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$, we compute: $\mathbf{X}_k = \mathbf{X} \mathbf{Q}_k \in \mathbb{R}^{n \times k}$ (\mathbf{Q} are eigvecs of $\mathbf{X}^T \mathbf{X}$).
- When using the SVD of \mathbf{X} , we can instead get this from:
- $\mathbf{X}_k = \mathbf{X} \mathbf{V}_{:,1:k}$

columns \mathbf{V} are
eigenvectors of $\mathbf{X}^T \mathbf{X}$

Singular Value Decomposition (SVD)

$$\mathbf{X} = \mathbf{U}_{m \times m} \mathbf{\Sigma}_{m \times n} \mathbf{V}^*_{n \times n}$$

For PCA we want projections onto top k PCs.

- When we used a spectral decomposition, $\mathbf{X}^T \mathbf{X} = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$, we compute: $\mathbf{X}_k = \mathbf{X} \mathbf{Q}_k \in \mathbb{R}^{n \times k}$ (\mathbf{Q} are eigvecs of $\mathbf{X}^T \mathbf{X}$).
- When using the SVD of \mathbf{X} , we can instead get this from:
- $\mathbf{X}_k = \mathbf{X} \mathbf{V}_{:,1:k} = \mathbf{U}_{:,1:k} \mathbf{\Sigma}_{1:k,1:k} \in \mathbb{R}^{n \times k}$ ("scores" in PCA basis). $\mathbf{X} \mathbf{v}_i = \mathbf{\Sigma}_{i,i} \mathbf{u}_i$
- $\mathbf{X}_{k\text{-recon}} = \mathbf{X} \mathbf{V}_{:,1:k} (\mathbf{V}_{:,1:k})^T = \mathbf{U}_{:,1:k} \mathbf{\Sigma}_{1:k,1:k} (\mathbf{V}_{:,1:k})^T$
- We don't need to compute covariance matrix, or do the projections, we just need $SVD(\mathbf{X})$!

“Eckart Young theorem” 1936

- The SVD “k-reconstruction” produces the best k -rank approximation by the matrix norm, $\|X - X_{recon-k}\|_F$.
- First proven by Schmidt (of Gram-Schmidt fame) in 1907 for Froebenius norm.
- Later rediscovered by Eckart & Young 1936, also generalized to other norms..
- Thus, PCA provides the best low rank approximation to the data matrix.

$$\|\mathbf{A}\|_F \equiv \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

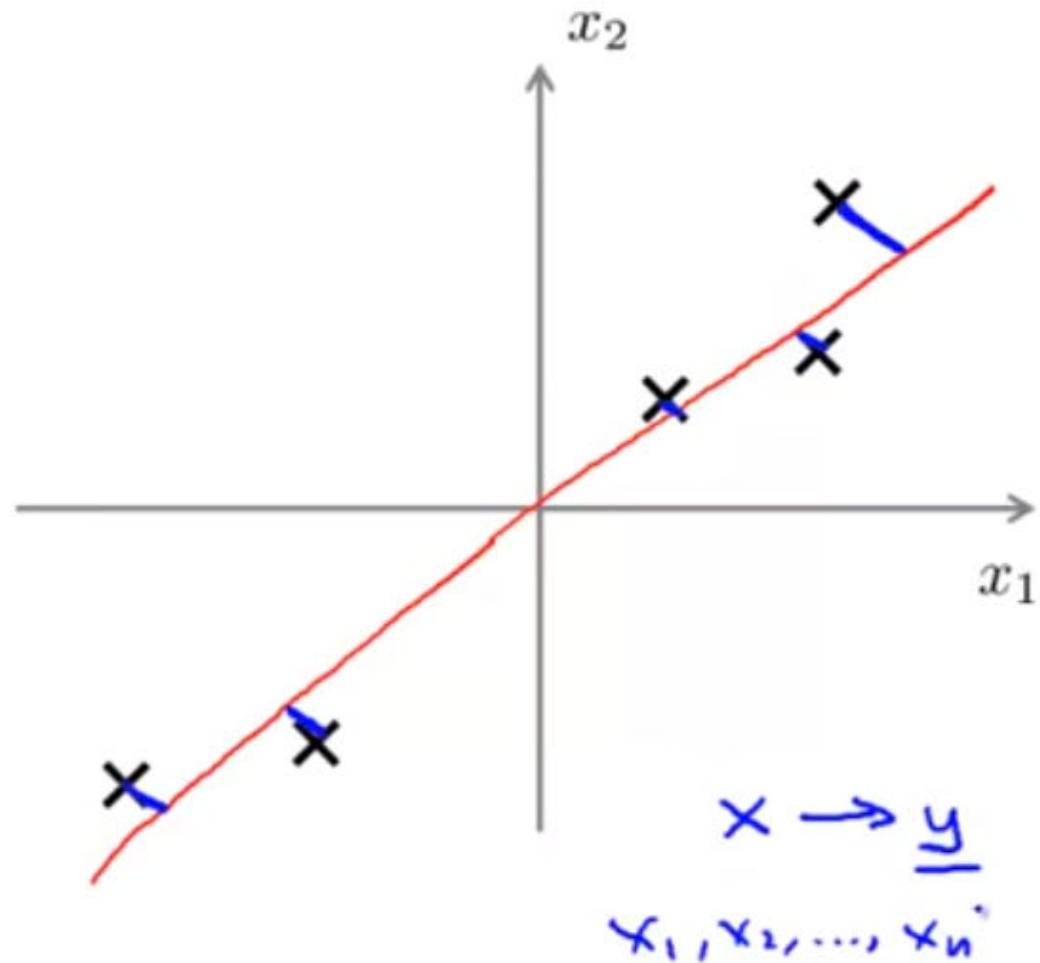
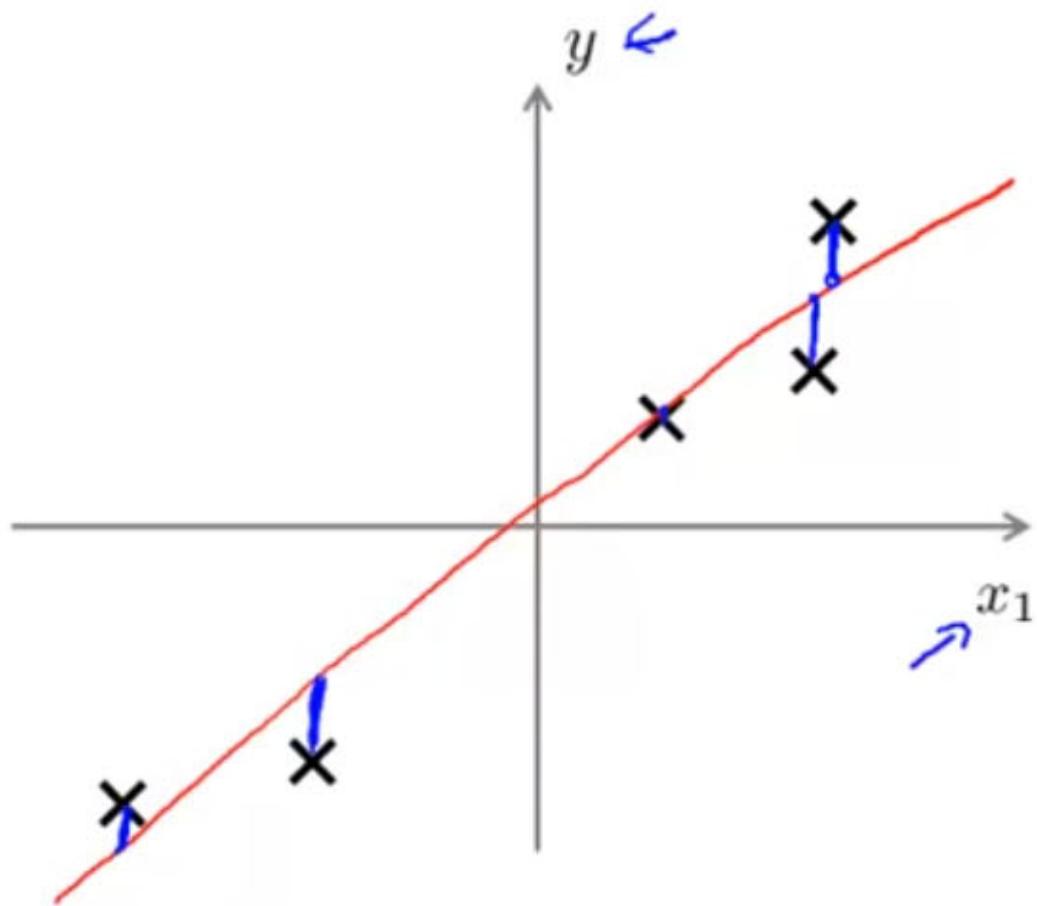
Practicalities: Reduced SVDs

For PCA and other applications, don't need the entire SVD, and can make do with "trimmed down" versions:

1. Full SVD
2. Thin SVD (remove columns of U not corresponding to rows of V^*)
3. Compact SVD (remove vanishing singular values and corresponding columns/rows in U and V^*),
4. Truncated SVD (keep only largest t singular values and corresponding columns/rows in U and V^*)

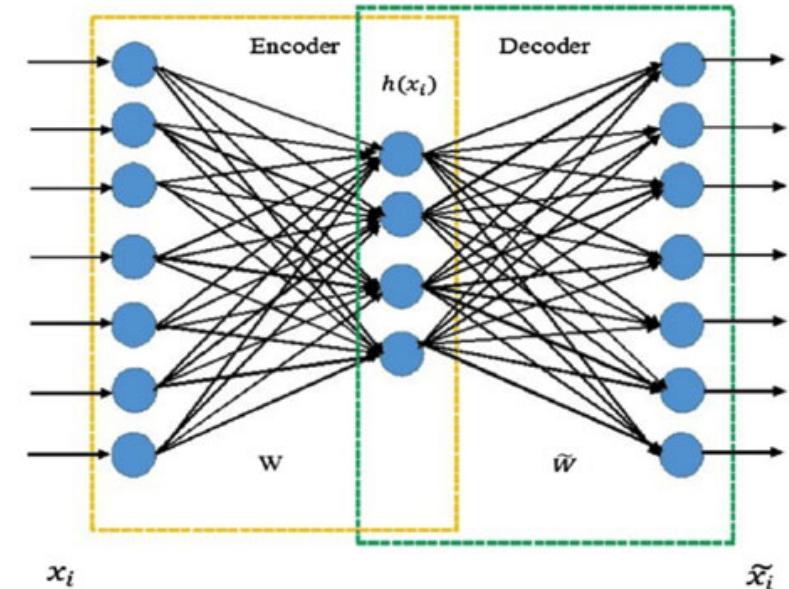
M $m \times n$	U $m \times m$	Σ $m \times n$	V^* $n \times n$
M $m \times n$	U_n $m \times n$	Σ_n $n \times n$	V^* $n \times n$
M $m \times n$	U_r $m \times r$	Σ_r $r \times r$	V_r^* $r \times n$
\bar{M} $m \times n$	U_t $m \times t$	Σ_t $t \times t$	V_t^* $t \times n$

PCA is not linear regression



PCA from neural networks!

- Special kind of neural network, called an *autoencoder* recovers the same subspace as PC-k.
- Autoencoder tries to compress a data set by trying to predict itself back after going through a bottleneck.



- For a linear autoencoder with one hidden layer containing k nodes, using squared loss, the hidden layer representation is equivalent to that found from PCA-k (although W may correspond to different eigenvectors that span the same space).
- Can generalize by making non-linear transfers, and more layers, etc.