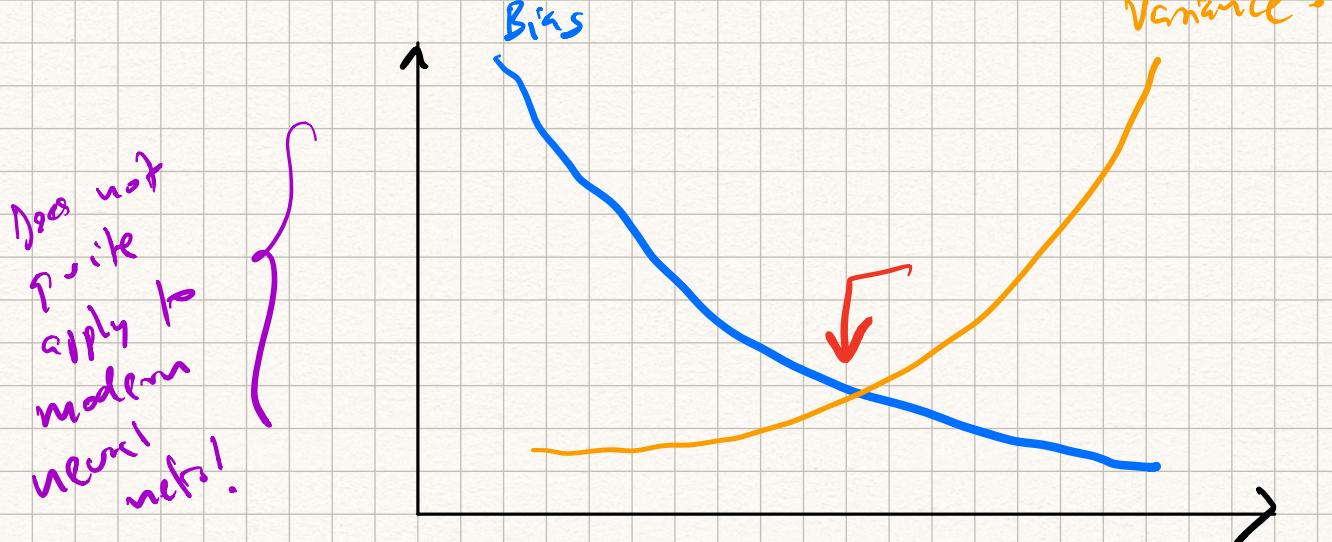
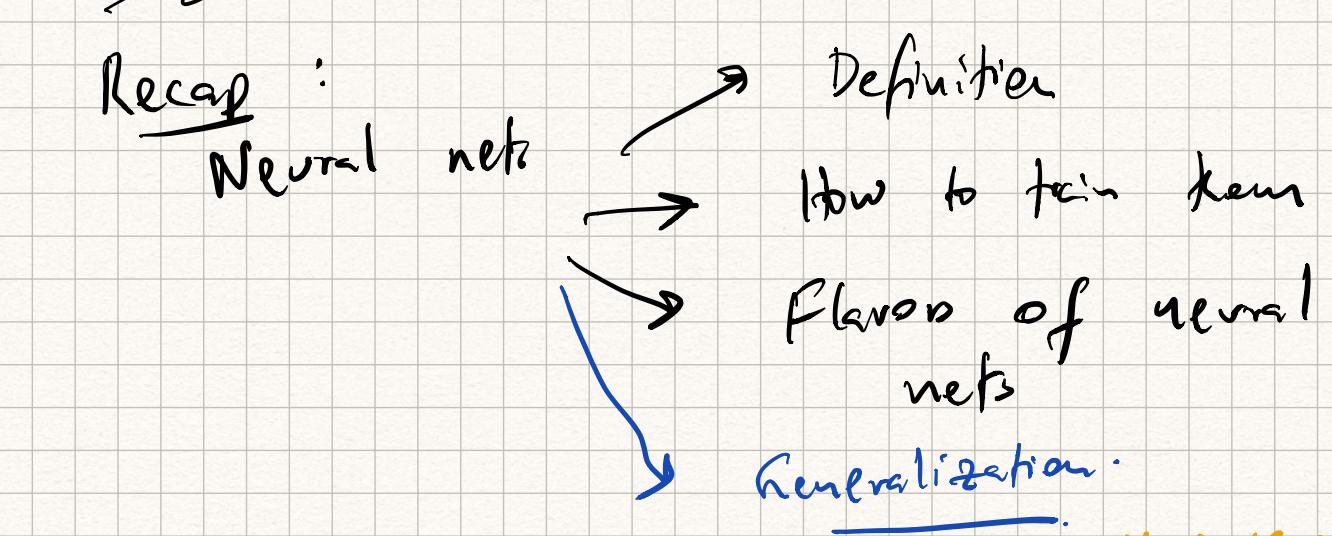


# ECE-GY 6143

## Intro to ML

- Recap: Neural net architectures
- Generalization in neural nets
- Intro to Unsupervised Learning
  - PCA
  - Autoencoders.



# Model Complexity

Classical approach to control bias-variance

## - Regularization.

Issue with neural nets : Dataset size.

e.g. Object detection

$\sim 928$  millions.

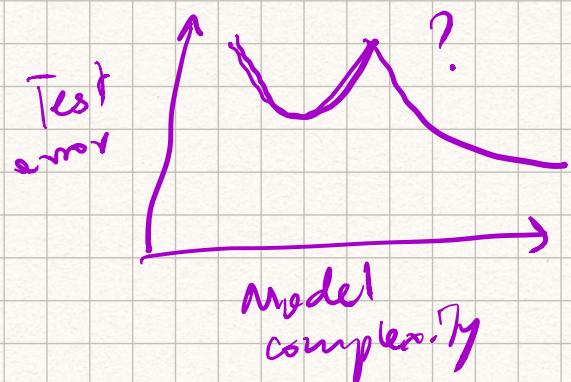
ImageNet  $\sim 14$  million

$$n \ll p$$

[ResNet - 1000].

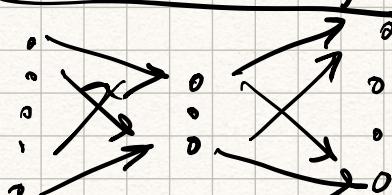
(P)

(n)



## Controlling Generalization

- Add regularization
- Adding bottleneck layers.



[Autoencoders/  
PCA].

## • Early stopping.

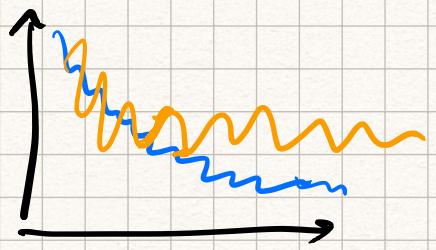


— Test error  
— Train error

Iteration count

- Stop training as soon as test error increases.

- May not always work  
(fluctuations in test error).



- Weight decay. [form of regularization].

- Dataset augmentation.  $n \uparrow \ll p$

- Apply transformations to input data samples

- Using images :  
Shifting      Shearing  
Rotation      Flipping.  
Cropping

- Dropout.

$$n \ll p$$

- Zero out certain neurons

- Define a binary random variable

$$h_i = m_i \phi(z_i) \quad m_i \sim 0 \text{ with prob } p$$

$$\frac{\partial}{\partial z_i} L = \frac{\partial}{\partial h_i} L \cdot m_i \cdot \phi'(z_i) \quad | \text{ with } 1-p.$$

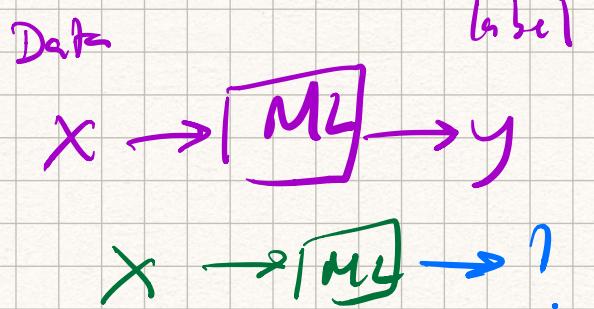
- Transfer learning.

- Fine tune using an already existing network.



## Unsupervised Learning

Supervised learning

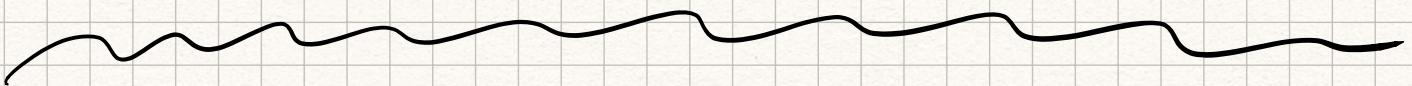


Unsupervised learning

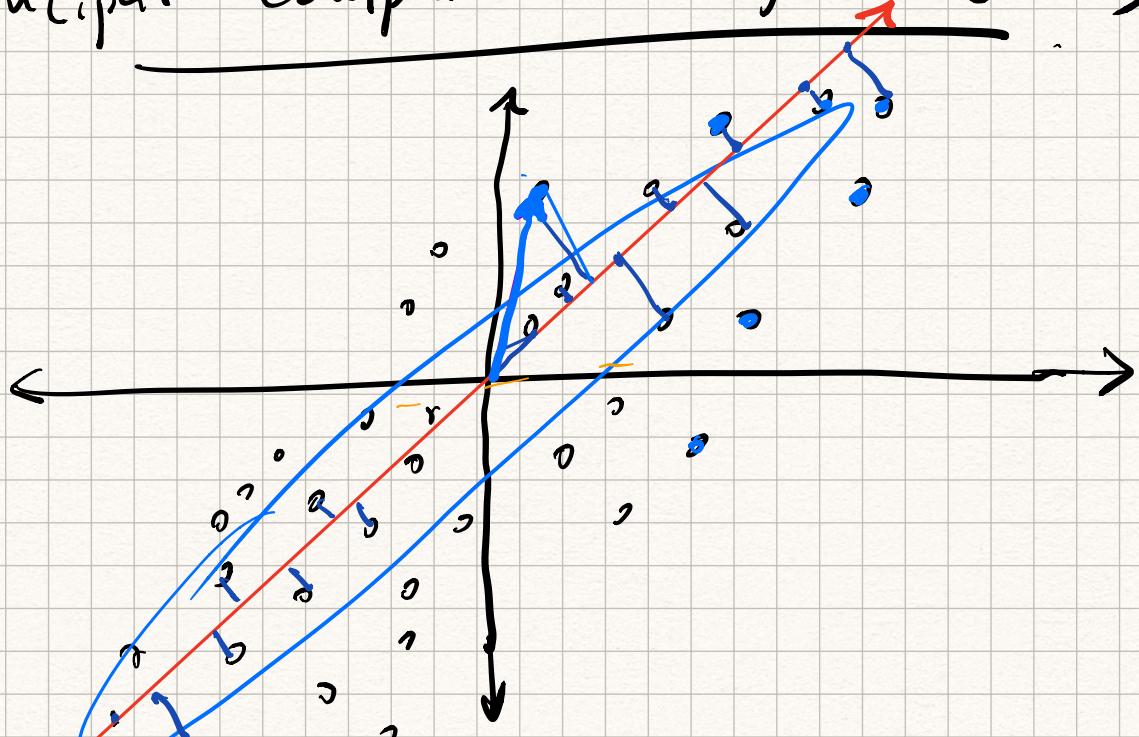
— Visualization

— Dimensionality reduction

— Clustering



## Principal Components Analysis (PCA)



$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n \times d}$$

Goal: Compute direction vector  $w$  such that  $w$  captures maximal variance.

$$\max S(w) = \frac{1}{n} \sum_{i=1}^n \langle x_i, w \rangle^2$$

s.t.  $\|w\|_2 = 1$

$$= \frac{1}{n} \|Xw\|_2^2$$

$\|v\|_2^2 \hat{=} v^T v$

$$= \frac{1}{n} (Xw)^T Xw$$

$$= \frac{1}{n} \overbrace{w^T X^T X w}^{w^T A^T A w}$$

$$\max_{\|w\|_2=1} w^T \underbrace{\frac{X^T X}{n}}_{A^T A} w$$

$$A_{d \times d}$$

$$Av = \lambda v$$

$$\sqrt{\lambda} Av \hat{=} \sqrt{\lambda} v$$

$$= \lambda$$

Achieved when  $w$  is the top eigenvector of  $A$  [ $=$  sample covariance matrix of your data].

$$A \xrightarrow{\text{---}} \left[ \begin{array}{c} (\lambda_1, v_1) \\ (\lambda_2, v_2) \\ \vdots \end{array} \right], \dots, (\lambda_d, v_d).$$

$v_1 \rightarrow$  "First principal component direction".

$$\underline{Xv_1} = \left[ x_1^T v_1 ; x_2^T v_1 ; \dots ; x_n^T v_1 \right]$$

$\rightarrow$  "First principal component scores"

Iteratively repeat for  $v_2, v_3, \dots$

$\overbrace{v_1, v_2, v_3, v_4, \dots, v_d}$  } Basis for the data space.

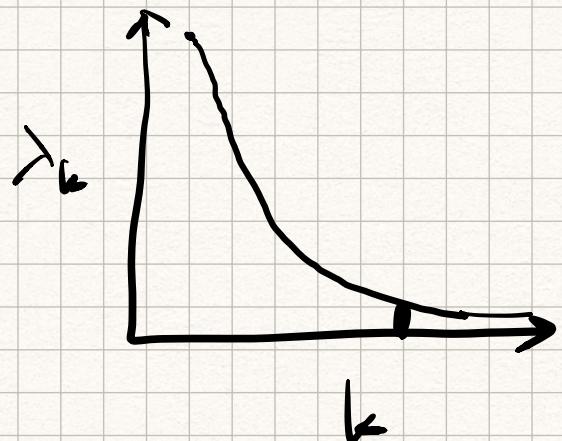
$\rightarrow$  Normalized, orthogonal.

$$x_i = \sum_{j=1}^d \frac{\langle x_i, v_j \rangle}{\|v_j\|} v_j$$

$\downarrow$  Scores       $\uparrow$  directions.

## Points to note

- Always center data.
- Make sure there are no outliers.
- Choose # principal component scores ( $k$ ) wisely.



## PCA vs. networks

PCA : "linear autoencoder".

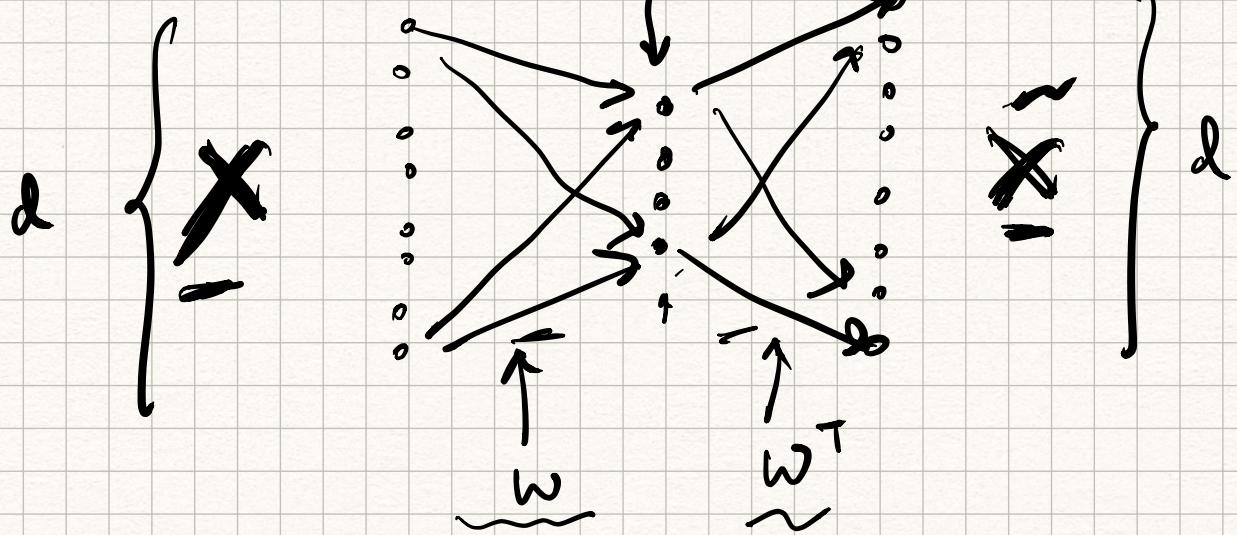
$$[X_{v_1}, X_{v_2}, \dots, X_{v_k}]_{n \times k}$$

$$\therefore = X W$$

$$W = [v_1, v_2, \dots, v_k]$$

$$\tilde{X} \approx \underline{X} \underline{W} \underline{W}^T$$

$\hookrightarrow$  hidden neurons - with linear activations



"Autoencoders" = Encoders  
+ Decoders

→ Denoising autoencoders

→ Sparse autoencoders

→ Stacked autoencoders

