

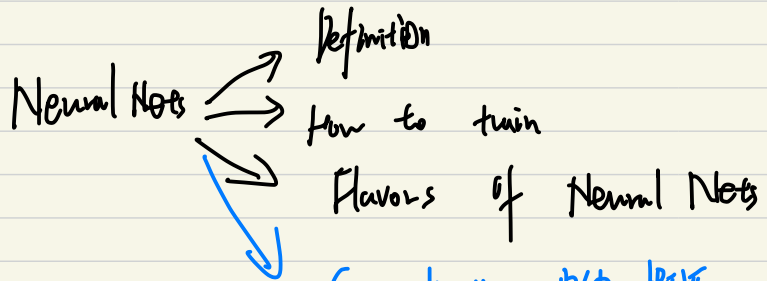
L11

★ Recap : Neural Net Architecture

★ Generalization in neural nets

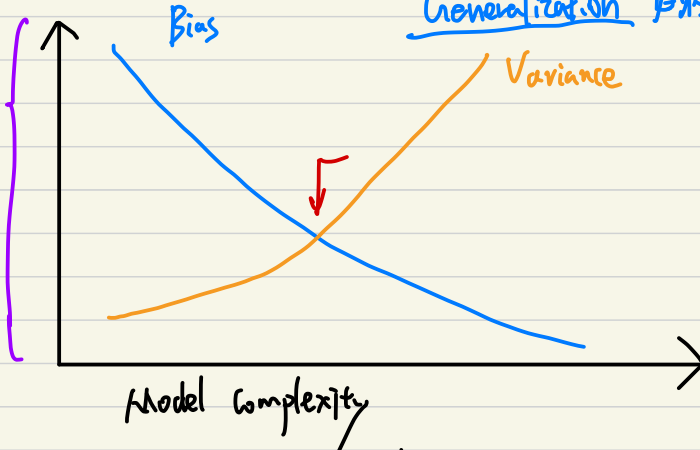
★ Intro to Unsupervised Learning { PCA
Autoencoders

Recap :



Generalization Bias-Variance

Does not
suit to
Neural Nets



Classic approach to control bias-variance — Regularization

io) 20

Issue with neural nets : Dataset size Lots of parameters

e.g. Object detection [ResNet - 1000]

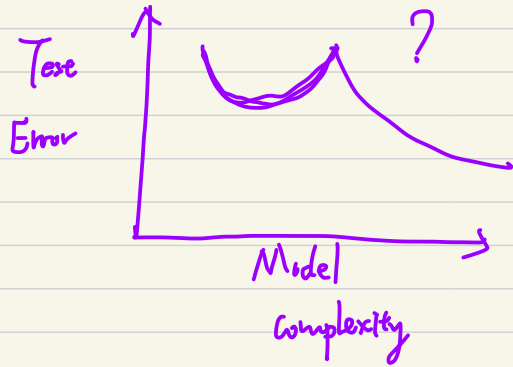
~ 928 million

(p) parameters

ImageNet ~ 14 million

(n)

$h \ll p$

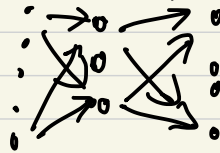


控制生成模型的手段:

Controlling generalization

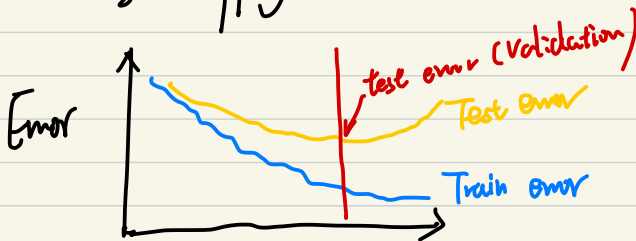
- Add regularization

- Add bottleneck layer

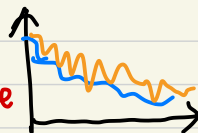


[Autoencoder / PCA]
bottleneck layer

- Early stopping



— Stop training as soon as test error increase



— May not always work (fluctuation in test error)

- Weight decay (form of regularization)

- Dataset augmentation

- apply transformation to input data samples

$$h \leftarrow p$$

eg. images shifting/rotation/flipping/cropping

- Drop out

$$h \leftarrow p$$

- zero out certain neurons

- Define a binary random variables

$$h_i = m_i \phi(z_i)$$

$m_i \sim 0$ with probability p

1 with $1-p$

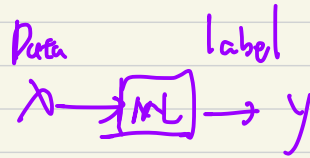
$$\frac{\partial L}{\partial z_i} = \frac{\partial L}{\partial h_i} \cdot m_i \phi'(z_i)$$

- Transfer Learning

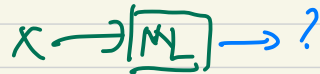
- Fine Tune using an already existing networks

Unsupervised learning

Supervised learning



Unsupervised learning

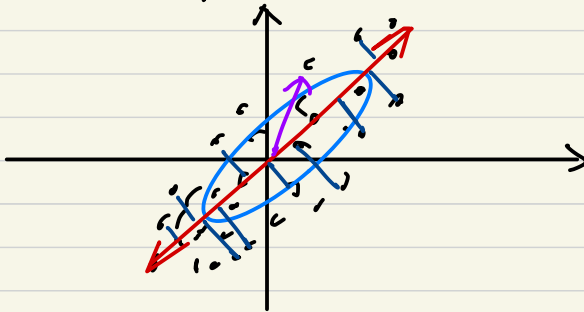


— Visualization

— Dimensional reduction

— Clustering

Principal Components Analysis (PCA)



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



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

Zero-mean

$$S(w) = \frac{1}{n} \sum_{i=1}^n \langle x_i, w \rangle^2 \quad \text{s.t. } \|w\|_2 = 1 \quad \text{单位方向型 (normalized)}$$

投影方差

$$= \frac{1}{n} \|XW\|_2^2 = \frac{1}{n} (XW)^T XW$$

$$= \frac{1}{n} W^T X^T X W$$

$$\max_{\|w\|_2=1} W^T \left[\frac{X^T X}{n} \right] W$$

$$\|w\|_2=1$$

A

$$\begin{aligned} Av &= \lambda v \\ v^T A v &= \lambda v^T v \\ &= \lambda \end{aligned}$$

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

$$A \rightarrow \begin{bmatrix} (\lambda_1, v_1), \dots, (\lambda_d, v_d) \\ (\lambda_2, v_2), \dots \end{bmatrix}$$

$v_1 \rightarrow$ "First principal component direction"

$$X_{v_1} = [x_1^T v_1, x_2^T v_1, \dots, x_n^T v_1]$$

\rightarrow First principal component scores

Iteratively, repeat for v_2, v_3, \dots

\rightarrow Normalized, orthogonal

$$X_i = \sum_{j=1}^d \underbrace{\langle x_i, v_j \rangle}_{\text{scores}} \underbrace{v_j}_{\text{directions}}$$



Points to use

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



PCA vs. neural networks

PCA: "linear auto encoder"

$$[xv_1, xv_2, \dots, xv_k]_{n \times k}$$

$$\approx XW$$

dimension reduce

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

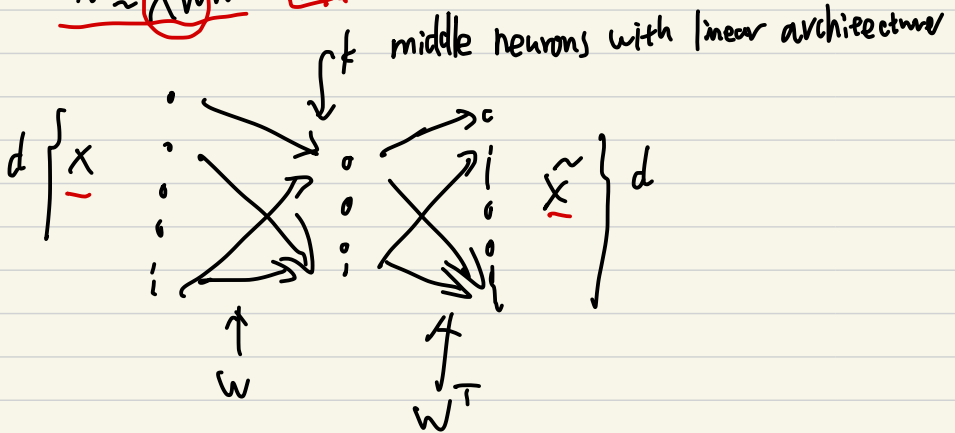
PCA

$$\tilde{x} \approx XWWT$$

还原

PCA

神经网络



"Auto encoders" = Encoder + Decoder

→ Denoising auto encoders

→ Sparse auto encoders