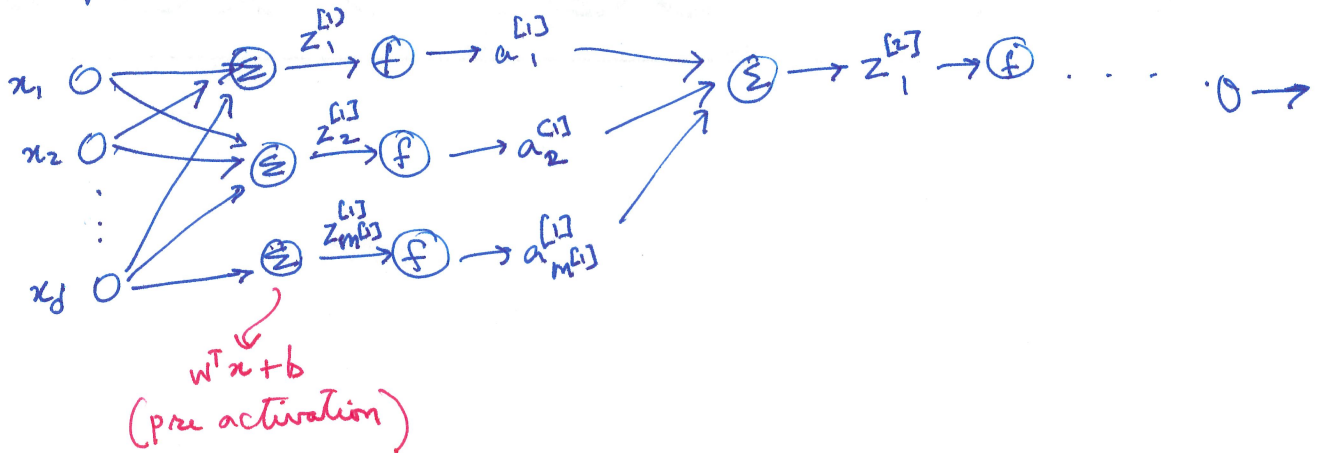


Neural Nets



Feed forward

Computing output by moving from left to right.

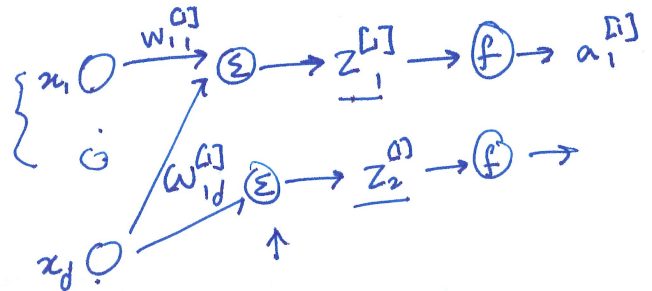
Gradients: Back propagation. Cost fn J .

$$\frac{\partial J}{\partial w_{jk}^{[l]}} = \frac{\partial J}{\partial o} \left(\frac{\partial o}{\partial w_{jk}^{[l]}} \right) \rightarrow \text{Chain rule.}$$

CNN

Filter + bias

Maxpooling



Unsupervised Learning:

$$S = \{ x^{(1)}, x^{(2)}, \dots, x^{(n)} \}$$

Reducing Dimensionality

Reducing Dimensionality:

Training Data : d -features.

↳ Reduce → Training : k -features.

Principal Component Analysis (PCA)

- Normalize data [mean data is 0, variance is 1]
- Find k -directions such that projecting the training data along those k dimensions maximizes the variance of the projected data.

$$S = \{x^{(1)}, \dots, x^{(n)}\}$$

$$X = \begin{bmatrix} -x^{(1)T} \\ \vdots \\ -x^{(d)T} \end{bmatrix}$$

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

Project X along vectors u_1, u_2, \dots, u_k where u_1, \dots, u_k are eigenvectors of norm 1 corresponding to the k -largest eigenvalues of A .

→ Use the SVD of X to discover the eigenvectors of A .

σ is a singular value of X iff σ^2 is eigenvalue of $X^T X$.

$$U = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_k \\ | & | & & | \end{bmatrix}$$

— left singular vectors of X
the columns of U are orthonormal.

$$\rightarrow U^T U = I.$$

$$U U^T = I \quad \leftarrow \text{True.}$$

} If B is ~~right~~ ^{left} inverse of A .

$$BA = I.$$

then B is a right inverse of A .

$$AB = I.$$

Clustering:

- Group data that is similar while keeping data that not similar in different clusters.

k-means: Find a k-partition $C_1 \dots C_k$ of S .

$C_i \cap C_j = \emptyset$ and $C_1 \cup C_2 \dots \cup C_k = S$.

s.t. ~~minimize~~

$$\min \sum_{j=1}^k \sum_{x^{(i)} \in C_j} \|x^{(i)} - \mu_j\|^2 \rightarrow \mu_j = \frac{\sum_{x^{(i)} \in C_j} x^{(i)}}{|C_j|}$$

k-means pick $\mu_1 \dots \mu_k$ randomly.

Repeat.

$$y^{(i)} = \operatorname{argmin}_j \|x^{(i)} - \mu_j\|^2$$

$$\mu_j = \frac{\sum_{i=1}^n x^{(i)} \mathbb{1}[y^{(i)} = j]}{\sum \mathbb{1}[y^{(i)} = j]}$$

$$\mathbb{1}[b] = \begin{cases} 1 & \text{if } b = \text{true} \\ 0 & \text{o.w.} \end{cases}$$

Linkage-based Clustering:

Start with every point being in its own cluster.

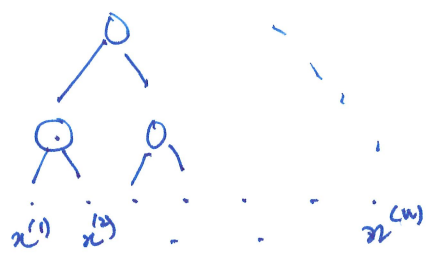
Repeat.

Merge the two clusters that are "closest".

$$\operatorname{dist}(C, D) = \min_{x \in C, y \in D} d(x, y)$$

$$= \operatorname{average} (d(x, y))$$

DENDROGRAM.



Training set are examples drawn from a distribution
Learning objective: Find this distribution.

- First pick $z \in \{1, \dots, k\}$. $p(z) \sim$ Multinomial
- Pick x from a distribution that depends on z .
 $p(x|z) \sim$ Gaussian distribution.
Bernoulli.

Given training set S , find parameters that define this distribution
parameter of $p(z)$
parameter of $p(x|z)$

Find parameters that maximize log-likelihood.

$$\begin{aligned}
 L(\theta) &= \log \prod_{i=1}^n p(x^{(i)}) \\
 &= \sum_{i=1}^n \log p(x^{(i)}) \\
 &= \sum_{i=1}^n \log \sum_{z=1}^k p(x^{(i)}, z) \\
 &\geq \sum_{i=1}^n \sum_{z=1}^k \underbrace{Q_i(z)} \log \frac{p(x^{(i)}, z)}{Q_i(z)} = F(Q, \theta)
 \end{aligned}$$

x, y ↗

↳ prob distribution on Z .
 $\forall z Q_i(z) \geq 0$ & $\sum_{z=1}^k Q_i(z) = 1$.

Goal: Find θ that maximize $L(\theta)$.

Repeat:

$Q_i(z) = p(z|x^{(i)})$
↳ Depends on θ^t

$Q^{(t+1)} = \operatorname{argmax}_Q F(Q, \theta^t)$

$\theta^{(t+1)} = \operatorname{argmax}_\theta F(Q^{(t+1)}, \theta)$

Assumption: Can solved computationally.