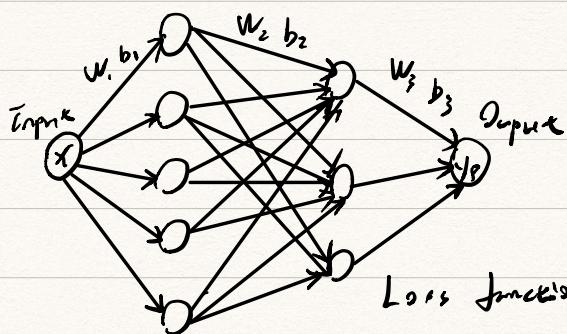


# 1. Backpropagation Algorithm for Neural Networks

(a)



$$\begin{aligned} \text{Loss function: } & \sum_{i=1}^n L(y_i, f(x_i | w, b)) \\ & = \sum_{i=1}^n (y_i - f(x_i | w, b))^2 \end{aligned}$$

$$(b) \text{MSE}(\hat{y}) = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\begin{aligned} \frac{\partial \text{MSE}(\hat{y})}{\partial \hat{y}} &= \frac{1}{2n} \times 2(y_i - \hat{y}_i) \times (-1) \\ &= \frac{(\hat{y}_i - y_i)}{n} \end{aligned}$$

The code is included in the appendix.

(c) Linear function

$$\delta_{\text{linear}}(z) = z$$

$$\frac{\partial \delta_{\text{linear}}(z)}{\partial z} = 1$$

ReLU function

$$\delta_{\text{ReLU}}(z) = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}$$

$$\begin{cases} z & \text{otherwise} \\ \end{cases}$$

$$\frac{\partial \sigma_{\text{ReLU}}(z)}{\partial z} = \begin{cases} 0 & z \leq 0 \\ 1 & \text{otherwise} \end{cases}$$

tanh function

$$\sigma_{\tanh}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\begin{aligned} \frac{\partial \sigma_{\tanh}(z)}{\partial z} &= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2} \\ &= 1 - \sigma_{\tanh}^2(z) \end{aligned}$$

The code is included in the appendix.

- (d) The code is included in the appendix

$$\left( \frac{\partial \text{MSE}}{\partial y_p} \times \frac{\partial y_p}{\partial z_i} \right) \times \frac{\partial z_i}{\partial a_i} \times \frac{\partial a_i}{\partial z_{i-1}}$$

- (e) With the same batch size, the training error decreases as the increase of epoch. This is obvious because the number of iterations increases. However, with the same epoch, the training error increases not with the increase of batch size. This maybe due to as the increase of batch size, large batch tend to converge sharp minimizers.

The performance of three activation function is: tanh > ReLU > linear.  
The code is included in appendix.

(f) With 2 different activation function : tanh, ReLU. The MSE both decrease

2. Regularized and Kernel k-means

(a) 0

(b)

$$\lambda \times 2 \mu_i + \sum_{x_j \in C_i} 2(\mu_i - x_j) = 0.$$

first and then increase as the increase of network width. The best width is around 8. This matches my expectation. This is the trade off between width and MSE.

$$\lambda \mu_i + |C_i| \mu_i - \sum_{x_j \in C_i} x_j = 0$$

$$\mu_i = \frac{1}{|C_i| + \lambda} \sum_{x_j \in C_i} x_j$$

(c)

$$\min_{C_1, C_2, \dots, C_K} \sum_{i=1}^k \left( \|\mu_i\|_2 + \sum_{x_j \in C_i} \|x_j - \mu_i\|_2 \right)$$

(d)

$$\text{Class}(j) = \arg \min_k \|x_j - \mu_k\|^2$$

$$= \arg \min_k \|x_j - \frac{\sum_{l \in S_k} x_l}{|S_k|}\|^2$$

$$= \arg \min_k k(x_j, x_j) - \frac{2 \sum_{l \in S_k} k(x_j, x_l)}{|S_k|}$$

$$+ \frac{\sum_{l \in S_k} \sum_{i \neq l} k(x_l, x_i)}{|S_k|^2}$$

(e) I see there are multiple  $k(x_1, x_2)$  through out the algorithm. In order to save

time, all kernels shouldn't be calculated twice or more. So the idea is to precompute  $k(x_i, x_j)$  for all possible combinations of  $i$  and  $j$  and store them for use in the algorithm. Where  $i, j \in \{1, 2, \dots, n\}$  and  $i$  can be equal to  $j$ .

### 3. Expectation Maximization (EM) Algorithm: in action!

(a) The initial guesses for 3 different algorithms are very different because the distance between different clusters are small.

After looking at the results of 3 algorithms. The EM performs the best, next is Kmeans, the worse is kDA

(b) After the increase of factor from 1 to 10, the distance between different clusters are much larger (The clusters are much obvious). Now the initial guesses are much closer to the truth. And the results show that the performances of the 3 algorithms are very close. They all perform very well.

### 4. One dimensional Mixture of Two Gaussians

$$(a) P_Z(z) = \begin{cases} \rho & z=1 \\ 1-\rho & z=2 \end{cases}$$

$$(X | Z=1) \sim N(\mu_1, \sigma_1^2)$$

$$\begin{aligned}\therefore P_\theta(X=x_i, Z_i=1) &= P_\theta(X=x_i | Z_i=1) P_Z(Z_i=1) \\ &= \frac{1}{\sqrt{2\pi} \sigma_1} e^{-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}} \times \beta = \frac{\beta}{\sqrt{2\pi} \sigma_1} e^{-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}}\end{aligned}$$

$$\begin{aligned}P_\theta(X=x_i, Z_i=2) &= P_\theta(X=x_i | Z_i=2) P_Z(Z_i=2) \\ &= \frac{1}{\sqrt{2\pi} \sigma_2} e^{-\frac{(x_i - \mu_2)^2}{2\sigma_2^2}} \times (1-\beta) = \frac{1-\beta}{\sqrt{2\pi} \sigma_2} e^{-\frac{(x_i - \mu_2)^2}{2\sigma_2^2}}\end{aligned}$$

$$\therefore P_\theta(X=x_i) = \frac{\beta}{\sqrt{2\pi} \sigma_1} e^{-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}} + \frac{1-\beta}{\sqrt{2\pi} \sigma_2} e^{-\frac{(x_i - \mu_2)^2}{2\sigma_2^2}}$$

$$(b) \quad l_\theta(x) = \sum_i \log \sum_{k=1}^2 N(x_i | \mu_k, \sigma_k^2) P(Z=k)$$

$$= \sum_i \log \left( \frac{\beta}{\sqrt{2\pi} \sigma_1} e^{-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}} + \frac{1-\beta}{\sqrt{2\pi} \sigma_2} e^{-\frac{(x_i - \mu_2)^2}{2\sigma_2^2}} \right)$$

It is obvious that this log likelihood function is very hard to take derivative. So it is very hard to set its derivative to 0 and get the closed form solution of MLE.

$$(c) \quad l_\theta(x_i) \approx \log \left( \frac{\beta}{\sqrt{2\pi} \sigma_1} e^{-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}} + \frac{1-\beta}{\sqrt{2\pi} \sigma_2} e^{-\frac{(x_i - \mu_2)^2}{2\sigma_2^2}} \right)$$

According to Jensen's Inequality

$$\log(a_1 p_1 + a_2 p_2) \leq a_1 p_1 + a_2 p_2 \quad (p_1, p_2 > 0)$$

$$l_{\theta}(x_i) = \log p(x_i | \theta) = \log \sum_{z_i} p(x_i, z_i | \theta)$$

$$= \log \sum_{z_i} \frac{q(z_i | x_i, \theta) p(x_i, z_i | \theta)}{q(z_i | x_i, \theta)}$$

$$= \log \mathbb{E}_q \left( \frac{p(x_i, z_i | \theta)}{q(z_i | x_i, \theta)} \right)$$

there is concave function

$$\geq \mathbb{E}_q \left[ \log \left( \frac{p(x_i, z_i | \theta)}{q(z_i | x_i, \theta)} \right) \right]$$

which  $\log(x) \geq \mathbb{E}_q \left[ \log \frac{p(x=x_i, z_i=k)}{q(x=x_i, z_i=k | x_i)} \right]$

$$(d) \text{ lower bound : } \mathbb{E}_q \left[ \log \frac{p_0(x=x_i, z_i=k)}{q_0(z_i=k | x=x_i)} \right]$$

$$= - \sum_z q_0(z_i=k | x=x_i) \log [q_0(z_i=k | x=x_i)]$$

$$+ \sum_z q_0(z_i=k | x=x_i) \log [p_0(x=x_i, z_i=k)]$$

$$= H(q_0(z_i=k | x=x_i)) + \mathbb{E}_q [L_c(x_i, z_i | \theta)]$$

$$\equiv F_i(q, \theta)$$

$$\therefore \text{In E-step : lower bound : } F_i(q, \theta^t)$$

$$\text{In M-step : lower bound : } F_i(q^{t+1}, \theta)$$

$$\therefore \text{In I-step : } q^{t+1} = \arg \max_q F(q, \theta^t)$$

$$\text{In M-step : } \theta^{t+1} = \arg \max_{\theta} F(q^{t+1}, \theta)$$

$\therefore$  Alternating between E-step and M-step until

both try to maximize  $F(g, \theta)$

$$(e) P(z_i=k | x_i, \theta^t) = \frac{P(z_i, x_i | \theta^t)}{P(x_i | \theta^t)}$$

$$= \frac{P(z_i, x_i | \theta^t)}{\sum_{z_i} P(x_i, z_i | \theta^t)} = \frac{P(x_i | z_i, \theta^t) P(z_i | \theta^t)}{\sum_{z_i} P(x_i | z_i, \theta^t) P(z_i | \theta^t)}$$

$$\therefore g_{i,1}^{t+1} = \frac{\alpha_k^t N(x_i | \mu_k^t, \sigma_k^t)}{\sum_k \alpha_k^t N(x_i | \mu_k^t, \sigma_k^t)}$$

$$\therefore g_{i,1}^{t+1} = \frac{\frac{\beta^t}{\sqrt{2\pi}\sigma_1^t} e^{-\frac{(x_i - \mu_1^t)^2}{2\sigma_1^{t^2}}}}{\frac{\beta^t}{\sqrt{2\pi}\sigma_1^t} e^{-\frac{(x_i - \mu_1^t)^2}{2\sigma_1^{t^2}}} + \frac{1-\beta^t}{\sqrt{2\pi}\sigma_2^t} e^{-\frac{(x_i - \mu_2^t)^2}{2\sigma_2^{t^2}}}}$$

$$g_{i,2}^{t+1} = \frac{\frac{1-\beta^t}{\sqrt{2\pi}\sigma_2^t} e^{-\frac{(x_i - \mu_2^t)^2}{2\sigma_2^{t^2}}}}{\frac{\beta^t}{\sqrt{2\pi}\sigma_1^t} e^{-\frac{(x_i - \mu_1^t)^2}{2\sigma_1^{t^2}}} + \frac{1-\beta^t}{\sqrt{2\pi}\sigma_2^t} e^{-\frac{(x_i - \mu_2^t)^2}{2\sigma_2^{t^2}}}}$$

$$(f) \mu_1^{t+1} = \frac{\sum_i g_{i,1}^{t+1} x_i}{\sum_i g_{i,1}^{t+1}}$$

(g) MOG fits very well with EM. But k-means  
 and soft k-means  
 does not fit well with EM.

↪ k-means gives more flexibility than k-means.

MoG allows non-spherical clusters And different scales  
of covariance.

MoG makes explicit assumptions in the form of  
statistical distributions thus easier to generalize.