

①

a) y_i is distributed as $\mathcal{N}(\mu, \sigma^2)$, where $\mu = \log(w x_i)$ and $\sigma^2 = 1$

$$p(y_i; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y_i - \log(w x_i))^2}{2}\right)$$

The full data likelihood is written by (assuming i.i.d):

$$\begin{aligned} L &= \prod_{i=1}^N p(y_i; \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y_i - \log(w x_i))^2}{2}\right) \\ &= (2\pi)^{-N/2} \prod_{i=1}^N \exp\left(-\frac{(y_i - \log(w x_i))^2}{2}\right) \end{aligned}$$

Then, log-likelihood is:

$$\begin{aligned} LL &= \log\left[(2\pi)^{-N/2} \prod_{i=1}^N \exp\left(-\frac{(y_i - \log(w x_i))^2}{2}\right)\right] \\ &= -\frac{N}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^N (y_i - \log(w x_i))^2 \end{aligned}$$

To infer the unknown parameter w , we maximize the log-likelihood

which means we maximize the probability/likelihood that the data is generated by the distribution.

b) To maximize the likelihood, we maximize log-likelihood as it is computationally easier to work with:

For this, we take the derivative of LL w.r.t. w & set it to zero:

$$\frac{\partial LL}{\partial w} = 0 \Rightarrow \frac{\partial \sum_{i=1}^N (y_i - \log(w x_i))^2}{\partial w} = 0$$

$$\Rightarrow \sum_{i=1}^N \frac{x_i}{w x_i} (y_i - \log(w x_i)) = 0$$

$$\Rightarrow \frac{1}{w} \sum_{i=1}^N (y_i - \log(w x_i)) = 0$$

$$\Rightarrow \sum_{i=1}^N y_i = \sum_{i=1}^N \log(w x_i) = \sum_{i=1}^N \log(x_i) + N \log w$$

c) for scalar w we have: $\|w\|_2^2 = w^2$

The new log-likelihood is:

$$LL^{new} = -\frac{N}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^N (y_i - \log(w x_i))^2 - \alpha w^2$$

Thus the derivative is:

$$\frac{\partial LL^{new}}{\partial w} = -\frac{1}{2} \frac{\partial \sum_{i=1}^N (y_i - \log(w x_i))^2}{\partial w} - \frac{\alpha \partial w^2}{\partial w}$$

$$= -\frac{1}{2} \sum_{i=1}^N \frac{2}{w} (y_i - \log(w x_i)) - 2w\alpha$$

$$= \frac{+1}{w} \sum_{i=1}^N (y_i - \log(w x_i)) - 2w\alpha$$

$$\frac{\partial LL^{new}}{\partial w} = 0 \Rightarrow \frac{+1}{w} \sum_{i=1}^N (y_i - \log(w x_i)) - 2w\alpha = 0$$

$$\Rightarrow 2\alpha w^2 = \sum_{i=1}^N (y_i - \log(w x_i)) \Rightarrow$$

$$\Rightarrow \sum_{i=1}^N y_i = \sum_{i=1}^N \log(w x_i) + 2\alpha w^2$$

d) if $\alpha \rightarrow \infty \Rightarrow -\alpha \|w\|_2^2 \rightarrow -\infty$

Since we want to maximize the likelihood then $\|w\|_2^2$ should become small, i.e. $w \rightarrow 0$ as $\alpha \rightarrow \infty$ in order to maximize the likelihood.

Thus, this regularizer shifts w towards zero!

2) We first introduce some notations:

\hat{y} : output of neural network

a) z_0 : input to the activation of the output node

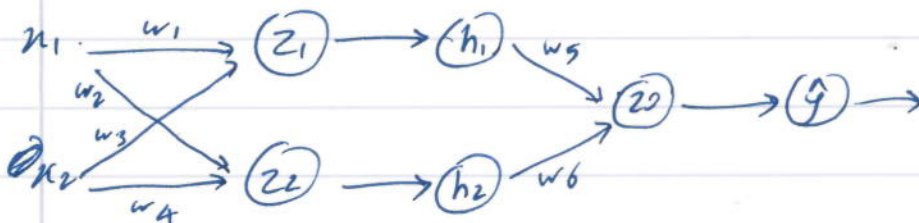
h_1 : the output of the first hidden node

h_2 : the output of the second hidden node

z_1 : input to the activation of first hidden node

z_2 : input to the activation of second hidden node

The network then can be seen as:



Thus:

$$\frac{\partial \mathcal{L}}{\partial w_5} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0} \frac{\partial z_0}{\partial w_5}$$

$$\frac{\partial \mathcal{L}}{\partial w_6} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0} \frac{\partial z_0}{\partial w_6}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0} \frac{\partial z_0}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0} \frac{\partial z_0}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} \\ &= \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0} \frac{\partial z_0}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} \end{aligned}$$

$$\frac{\partial \Sigma}{\partial w_2} = \frac{\partial \Sigma}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_2} = \frac{\partial w}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0} \frac{\partial z_0}{\partial w_2} = \frac{\partial \Sigma}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0} \frac{\partial z_0}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

$$= \frac{\partial \Sigma}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_0} \frac{\partial z_0}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$

similarly, you should obtain $\frac{\partial \Sigma}{w_3}, \frac{\partial \Sigma}{w_4}$

b)

1. initialize the parameters and set a learning rate γ .

2. repeat:

$$w_j^{(t+1)} = w_j^{(t)} - \gamma \frac{\partial \Sigma}{\partial w_j^{(t)}}, \quad j = 1 \dots 6$$

until convergence.

↳ the error does not change or the parameters do not change anymore.

⇒ yields a local optima.

c)

the activation function can be written as:

$$f(z) = \frac{1}{2} (2z)^q$$

if $q = 1 \Rightarrow f(z) = 2z \Rightarrow$ it is just a simple ^{linear} transformation.

thus:

$$\left. \begin{aligned} \hat{y} &= f(z_0) = 2z_0 = 2(w_5 h_1 + w_6 h_2) \\ h_1 &= 2z_1 = 2(w_1 x_1 + w_3 x_2) \\ h_2 &= 2z_2 = 2(w_2 x_1 + w_4 x_2) \end{aligned} \right\}$$

$$\Rightarrow \hat{y} = 2(w_5(2(w_1 x_1 + w_3 x_2)) + w_6(2(w_2 x_1 + w_4 x_2)))$$

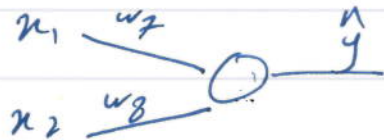
$$= 4w_1 w_5 x_1 + 4w_3 w_5 x_2 + 4w_2 w_6 x_1 + 4w_4 w_6 x_2$$

$$= (4w_1 w_5 + 4w_2 w_6) x_1 + (4w_3 w_5 + 4w_4 w_6) x_2$$

$$\Rightarrow \hat{y} = w_7 x_1 + w_8 x_2$$

$$\text{where } \begin{cases} w_7 = 4w_1w_5 + 4w_2w_6 \\ w_8 = 4w_3w_5 + 4w_4w_6 \end{cases}$$

Thus, we can replace the network with a simpler network that just computes the weighted sum of the input attributes.



and we do not need any activation function $f(z)$

d) yes: as we saw in (c), for $q=1$ we have $\hat{y} = w_7 x_1 + w_8 x_2$

i.e.: $\hat{y} = w^T x$, where $w^T = [w_7, w_8]$

This is just the formulation of ~~simple~~ ^{a basic} linear regression model.

3

a) The training error would be related to misclassification error, i.e., number of training items that are wrongly classified.

In this data set, for optimal solution, no item would be misclassified, thus ~~training~~ training error is zero.

b) The support vectors corresponding to γ 's:
 see 5.3.2.4 in the text book page (189-190) for details:
 ... at the optimum, all of the d_n that do not correspond to support vectors will be zero, ...

in optimum: $t_{new} = \text{sign}\left(\sum_{n=1}^N d_n t_n x_n^T x_{new} + b\right)$, $b = t_n - \sum_{m=1}^N d_m t_m x_m^T x_{new}$

Thus zero α_n 's will not have an impact on t_{new} , b , and they can be discarded. Therefore we can train the model using only the support vectors.

c) $O(d)$

see 5.3.2.3 : Making predictions.

we have:

$$t_{new} = \text{sign} \left(\underbrace{\sum_{n=1}^N t_n t_n X_n^T X_{new}}_{w^*} + b \right), \quad b = t_n - \sum_{m=1}^N \alpha_m t_m X_m X_n^T$$

computing w^* , b in general takes $O(dN)$, but they can be computed ~~are~~ independent of test (prediction). In other words, they can be obtained immediately after training and do NOT depend on particular test data point \Rightarrow we do ~~not~~ take them into account for test.

Note that if number of support vectors $\ll N$

\Rightarrow complexity of w^* , b would be ~~$O(dN)$~~ $O(d)$ and $O(d)$ (instead of $O(Nd)$)

\Rightarrow for test : we need compute $w^* X_{new}$ and $\text{sign}(w^* X_{new} + b)$ $\rightarrow O(d)$

\Rightarrow test complexity = $O(d)$

4

a) K-means cost function is written by:

$$R(\mu, Z; X) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|x_n - \mu_k\|_2^2, \quad z_{nk} \in \{0, 1\}, \sum_k z_{nk} = 1$$

if $k=1 \Rightarrow z_{nk} \dots$ is always 1:

$$\begin{aligned} R(\mu, Z; X) &= \sum_{n=1}^N z_{nk} \|x_n - \mu_k\|_2^2 = \sum_{n=1}^N z_{nk} (x_n - \mu_k)^T (x_n - \mu_k) \\ &= \sum_{n=1}^N z_{nk} (x_n - \mu_k)^T (x_n - \mu_k) \end{aligned}$$

There is only one mean: \Rightarrow

$$R(\mu, Z; X) = \sum_{n=1}^N (x_n - \mu)^T (x_n - \mu)$$

$$\frac{\partial R}{\partial \mu} = \frac{\partial \sum_{n=1}^N (x_n - \mu)^T (x_n - \mu)}{\partial \mu} = 2 \sum_{n=1}^N x_n - \mu$$

$$\frac{\partial R}{\partial \mu} = 0 \Rightarrow \sum_{n=1}^N (x_n - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{n=1}^N x_n$$

b) We know $R(\mu, Z; X) \geq 0$

Then if we show that for a solution $R(\mu, Z; X) = 0$ then this solution will be the optimal solution, because it yields the minimum possible cost.

if $\mu_k = x_k \quad \forall 1 \leq k \leq K$

$$\Rightarrow R(\mu, Z; X) = \sum_{i=1}^N z_{ni} \underbrace{\|x_n - x_{n_i}\|_2^2}_{=0} = 0$$

Thus the solution in which each ~~mean~~ mean is exactly one of the datapoints yields $R=0$ which is an optimal solution.

c) for $k^{\min} \leq k \leq k^{\max}$

compute AEC_k

compute BEC_k

choose k_{AIC}^* and k_{BIC}^* for which AIC and BIC numbers are minimal, i.e.:

$$k_{AIC}^* = \arg \min_k AEC_k$$

$$k_{BIC}^* = \arg \min_k BEC_k$$

AEC_k is obtained by: $\underbrace{-ll}_{\text{negative log-likelihood}} + \underbrace{c(U, Z)}_{\text{complexity}}$

BEC_k is obtained by: $\underbrace{-ll}_{\text{negative log-likelihood}} + \frac{1}{2} \underbrace{c(U, Z)}_{\text{complexity}} \ln N$

$$c(U, Z) = k * d + \text{related to } \pi_i \text{'s} \quad \left(\frac{N}{k} \right)$$

related to means

*all clusters have same size

if we know the size (of one cluster (π_i)) the size of other cluster is the same.