

①

a) $-1 < \cos(y-t) < +1 \Rightarrow \begin{cases} \hat{L}_{\max} = 2 \Rightarrow \cos(y-t) = -1 \Rightarrow y = \frac{7\pi}{4} \\ \hat{L}_{\min} = 0 \Rightarrow \cos(y-t) = 1 \Rightarrow y = \frac{3\pi}{4} \end{cases}$

b) $\frac{\partial L}{\partial w} = \frac{\partial (1 - \cos(w^T x + b - t))}{\partial w} = x \sin(w^T x + b - t)$

$\frac{\partial L}{\partial b} = \frac{\partial (1 - \cos(w^T x + b - t))}{\partial b} = \sin(w^T x + b - t)$

c) The batch cost function for N input data is written by:

$L(y, t) = \sum_{i=1}^N L(y_i, t_i) = \sum_{i=1}^N 1 - \cos(y_i - t_i)$

$\frac{\partial L}{\partial w} = \sum_{i=1}^N x_i \sin(w^T x_i + b - t_i)$

$\frac{\partial L}{\partial b} = \sum_{i=1}^N \sin(w^T x_i + b - t_i)$

Gradient Descent:

1) initialize \hat{w}, \hat{b}, γ

2) until convergence:

update $\hat{w} \leftarrow \hat{w} - \gamma \frac{\partial L}{\partial w}$
 update $\hat{b} \leftarrow \hat{b} - \gamma \frac{\partial L}{\partial b}$

②-1

* For simplicity, we assume two categories of age: $\begin{cases} \text{young if age} \leq 30 \\ \text{senior if age} > 30 \end{cases}$

a) For Bayes classifier, see Lecture 4 b.

* we assume a uniform prior distribution over class labels \Rightarrow

$$\begin{aligned} P(t_{\text{new}}=k | X, t, x_{\text{new}}) &= \frac{P(x_{\text{new}} | t_{\text{new}}=k, X, t) P(t_{\text{new}}=k)}{\sum_j P(x_{\text{new}} | t_{\text{new}}=j, X, t) P(t_{\text{new}}=j)} \\ &= \frac{P(x_{\text{new}} | t_{\text{new}}=k, X, t)}{\sum_j P(x_{\text{new}} | t_{\text{new}}=j, X, t)} \end{aligned}$$

$$\begin{aligned} P(t_{\text{new}}=L | X, t, x_{\text{new}}=(\text{young}, \text{sports})) &= \frac{P(x_{\text{new}} | t_{\text{new}}=L, X, t)}{\sum_{j \in \{L, H\}} P(x_{\text{new}} | t_{\text{new}}=j, X, t)} \\ &= \frac{2/2}{2/2 + 1/4} = \frac{4}{5} \end{aligned}$$

$$\begin{aligned} P(t_{\text{new}}=H | X, t, x_{\text{new}}=(\text{young}, \text{sports})) &= \frac{P(x_{\text{new}} | t_{\text{new}}=H, X, t)}{\sum_{j \in \{H, L\}} P(x_{\text{new}} | t_{\text{new}}=j, X, t)} \\ &= \frac{1/4}{2/2 + 1/4} = \frac{1}{5} \end{aligned}$$

b) Naive Bayes assumes the independence of components of x_{new} given the class label, i.e.,

$$P(x_{\text{new}} | X, t, t_{\text{new}}=k) = \prod_{d=1}^D P(x_d^{\text{new}} | X, t, t_{\text{new}}=k)$$

②-2

thus:

$$P(x_{\text{new}} | t_{\text{new}}=L, X, t) = P(x_{(1)}^{\text{new}} = \text{young} | t_{\text{new}}=L, X, t) P(x_{(2)}^{\text{new}} = \text{sports} | t_{\text{new}}=L, X, t)$$
$$= \frac{2}{2} \times \frac{2}{2} = 1$$

$$P(x_{\text{new}} | t_{\text{new}}=H, X, t) = P(x_{(1)}^{\text{new}} = \text{young} | t_{\text{new}}=H, X, t) P(x_{(2)}^{\text{new}} = \text{sports} | t_{\text{new}}=H, X, t)$$
$$= \frac{3}{4} \times \frac{1}{4} = \frac{3}{16}$$

$$\Rightarrow \begin{cases} P(t_{\text{new}}=L | X, t, x_{\text{new}}) = \frac{1}{1 + \frac{3}{16}} = \frac{16}{19} \\ P(t_{\text{new}}=H | X, t, x_{\text{new}}) = \frac{\frac{3}{16}}{1 + \frac{3}{16}} = \frac{3}{19} \end{cases}$$

c) The training data is not sufficient to include ~~data~~ ^{data} of type track.

In the naive-Bayes, we can write/assume:

$$P(x_{(2)}^{\text{new}} = \text{truck} | X, t, t_{\text{new}}=L) = P(x_{(2)}^{\text{new}} = \text{truck} | X, t, t_{\text{new}}=H)$$

$$\Rightarrow P(t_{\text{new}}=L | X, t, x_{\text{new}}) = \frac{2/2}{2/2 + 3/4} = \frac{4}{7}$$

$$P(t_{\text{new}}=H | X, t, x_{\text{new}}) = \frac{3/4}{2/2 + 3/4} = \frac{3}{7}$$

③

a) $O(nkd)$ \Rightarrow the method requires a full pass of data for each center to be sampled.

b, c, d): there are several solutions/possibilities, for instance:

in each iteration $j \in \{2, 3, \dots, k\}$ we construct a Markov chain of length m using Metropolis-Hasting method.

The proposal dist. is defined to be an independent and uniform distribution $q(x) = \frac{1}{N}$.

So with start with a random initial state x_1 and in each iteration $i \in \{2, \dots, m\}$ we sample a candidate y_i using $q(x)$.

We accept this candidate (i.e., $x_i = y_i$) with probability

$$\min \left(\frac{p(y_i)}{p(x_{i-1})} \frac{q(x_{i-1})}{q(y_i)}, 1 \right) = \min \left(\frac{d^2(y_i, c)}{d^2(x_{i-1}, c)}, 1 \right)$$

Then, in each of the $k-1$ iterations, we only need to compute the distances between m objects and $k-1$ centers.

Thus the complexity would be $O(mk^2d)$.

* for a very detailed & accurate solution see:

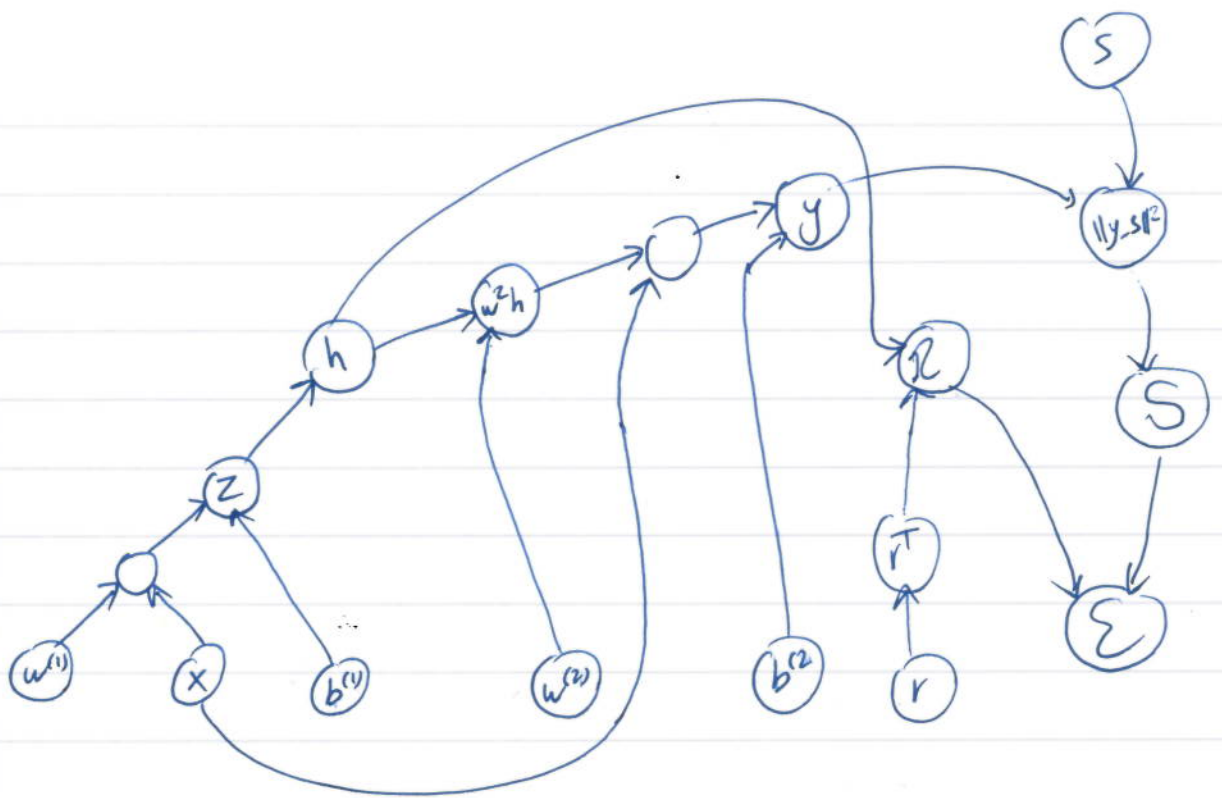
[Out of the scope of the course and the exam]

las.inf.ethz.ch/files/bachem16fast.pdf

or search: Fast and Provably Good seeding for k-means.

4

a)



b) $N \times k + k + N \times k + N$

c) we use chain rule in backpropagation, for example:

$$\frac{\partial \Sigma}{\partial w^{(1)}} = \frac{\partial R}{\partial w^{(1)}} + \frac{\partial S}{\partial w^{(1)}}$$

$$\frac{\partial R}{\partial w^{(1)}} = \frac{\partial R}{\partial h} \frac{\partial h}{\partial w^{(1)}} = \frac{\partial R}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial w^{(1)}}$$

$$\frac{\partial S}{\partial w^{(1)}} = \frac{\partial S}{\partial y} \frac{\partial y}{\partial w^{(1)}} = \frac{\partial S}{\partial y} \frac{\partial y}{\partial h} \frac{\partial h}{\partial w^{(1)}} = \frac{\partial S}{\partial y} \frac{\partial y}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial w^{(1)}}$$

Similarly:

$$\frac{\partial \Sigma}{\partial w^{(2)}} = \frac{\partial R}{\partial w^{(2)}} + \frac{\partial S}{\partial w^{(2)}}$$

$$\frac{\partial R}{\partial w^{(2)}} = \frac{\partial R}{\partial h} \frac{\partial h}{\partial w^{(2)}}$$

$$\frac{\partial S}{\partial w^{(2)}} = \frac{\partial S}{\partial y} \frac{\partial y}{\partial w^{(2)}}$$

similarly we can obtain $\frac{\partial \Sigma}{\partial b^{(1)}}$ and $\frac{\partial \Sigma}{\partial b^{(2)}}$.