

Problem 1.

find the parameters that maximize

$$CL(\theta) = \frac{1}{N} \sum_{n=1}^N \log P(y^n | x^n, \theta)$$

where all (x^n, y^n) generated from $P(x, y | \theta^0) = P(y | x, \theta^0) P(x | \theta^0)$

Proof.

$$CL(\theta) = \frac{1}{N} \sum_{n=1}^N \log P(y^n | x^n, \theta)$$

$$= \langle \log P(y^n | x^n, \theta) \rangle_{p(\text{data})}, \text{ where } p(\text{data}) = P(x^n, y^n | \theta^0)$$

$$= \langle \log P(y^n | x^n, \theta) \rangle_{P(x^n, y^n | \theta^0)}.$$

When the amount of training data is large enough.

$$CL(\theta) = \int P(x^n, y^n | \theta^0) \log P(y^n | x^n, \theta) dn$$

$$= \int \underbrace{P(y^n | x^n, \theta^0)} P(x^n | \theta^0) \underbrace{\log P(y^n | x^n, \theta)} dn$$

$$= - \int P(x^n | \theta^0) P(y^n | x^n, \theta^0) \log \frac{P(y^n | x^n, \theta^0)}{P(y^n | x^n, \theta)} dn$$

$$+ \int P(x^n | \theta^0) P(y^n | x^n, \theta^0) \log P(y^n | x^n, \theta^0) dn$$

$$= - \int P(x^n | \theta^0) KL(P(y^n | x^n, \theta^0) \| P(y^n | x^n, \theta)) dn + f(\theta^0)$$

Thus, to maximize $CL(\theta)$, $KL(P(y^n | x^n, \theta^0) \| P(y^n | x^n, \theta)) = 0$.

which means $\theta = \theta^0$.

Therefore, $CL(\theta)$ has an optimum at θ^0 .



Problem 2.

Older than 60

C	F	S	B
1	0	0	0
1	0	0	1
1	1	1	1
0	0	0	1

younger than 60

C	F	S	B
0	1	1	0
1	1	1	0

New customer = 0110.

Proof. Naive Bayes assume all features to be independent

$$\text{Thus, for older people} = \left. \begin{aligned} P(C=1) &= \frac{3}{4}, P(F=1) = \frac{1}{4} \\ P(S=1) &= \frac{1}{4}, P(B=1) = \frac{3}{4} \end{aligned} \right\} P(\text{old}) = \frac{2}{3}$$

$$\text{for younger people:} \left. \begin{aligned} P(C=1) &= \frac{1}{2}, P(F=1) = 1 \\ P(S=1) &= 1, P(B=1) = 0 \end{aligned} \right\} P(\text{y}) = \frac{1}{3}$$

Thus, the maximum likelihood will give us

$$P(\text{new} = \text{old}) \propto \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{2}{3} = \frac{1}{16 \times 8 \times 3}$$

$$P(\text{new} = \text{young}) \propto \frac{1}{2} \times 1 \times 1 \times 1 \times \frac{1}{3} = \frac{1}{6}$$

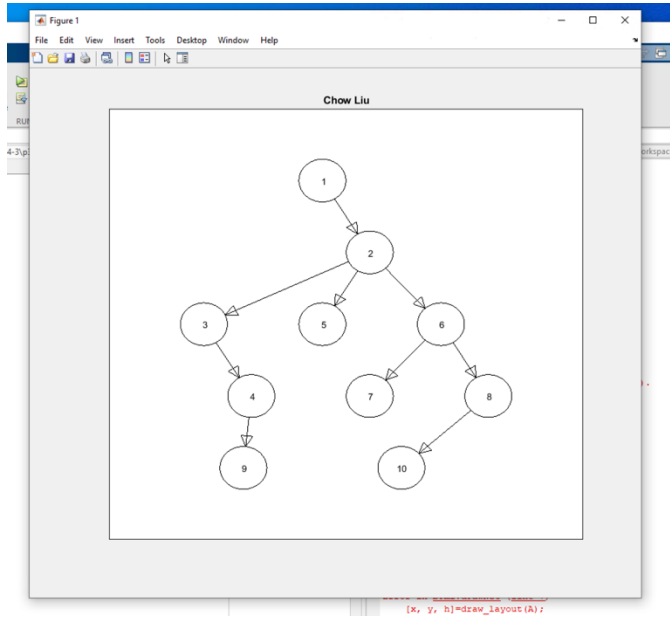
$$P(\text{new} = \text{old}) + P(\text{new} = \text{young}) = 1$$

$$\text{Thus, } P(\text{young}) = \frac{\frac{1}{6}}{\frac{1}{6} + \frac{1}{16 \times 24}} = \frac{4 \times 16}{4 \times 16 + 1} = \boxed{\frac{64}{65}}$$



Problem 3

The function is inside p3ChowLiu.m, the main code is inside p3main.m, run the p3main.m to use the data and the function. From the main code we can get the following picture:



Problem 4.

NB classifier for $x_i \in \{0, 1\}$ $\leftarrow \begin{cases} \theta_i^1 = P(x_i=1 | \text{class}=1) \\ \theta_i^0 = P(x_i=1 | \text{class}=0) \end{cases}$

$$\text{Proof. } P(\text{class}=0 | \vec{x}) \propto P(\text{class}=0) P(\vec{x} | \text{class}=0) \\ = p_0 \prod_i P(x_i | \text{class}=0)$$

$$\text{similarly, } = p_0 \prod_i (\theta_i^0)^{x_i} (1 - \theta_i^0)^{1-x_i}$$

$$\hookrightarrow P(\text{class}=1 | \vec{x}) \propto p_1 \prod_i (\theta_i^1)^{x_i} (1 - \theta_i^1)^{1-x_i}$$

$$\text{decision bound } P(\text{class}=0 | \vec{x}) = P(\text{class}=1 | \vec{x})$$

$$p_0 \prod_i (\theta_i^0)^{x_i} (1 - \theta_i^0)^{1-x_i} = p_1 \prod_i (\theta_i^1)^{x_i} (1 - \theta_i^1)^{1-x_i}$$

$$\log p_0 + \sum_i [x_i \log \theta_i^0 + (1-x_i) \log (1 - \theta_i^0)] \\ = \log p_1 + \sum_i [x_i \log \theta_i^1 + (1-x_i) \log (1 - \theta_i^1)]$$

$$\sum_i x_i (\log \theta_i^0 - \log (1 - \theta_i^1)) - \log \theta_i^1 + \log (1 - \theta_i^0) \\ + \sum_i (\log (1 - \theta_i^0) - \log (1 - \theta_i^1)) + \log \frac{p_0}{p_1} = 0.$$

$$\underbrace{\sum_i x_i \left(\log \frac{\theta_i^0 (1 - \theta_i^1)}{\theta_i^1 (1 - \theta_i^0)} \right)}_{w_i} + \underbrace{\log \frac{p_0}{p_1} + \sum_i \log \frac{1 - \theta_i^0}{1 - \theta_i^1}}_b = 0.$$

Hence, the decision holds according to $Wx + b = 0$, while the value of w^T and b are stated as above. \square

Problem 5

1. derive expressions using maximum likelihood.

We know that training data = (x^n, c^n) , $n \in \{1, \dots, N\}$.

$$\text{Thus, } p(c=1) = \frac{\text{number of data with } c=1}{N} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{c^i=1\}$$

$$\text{where } \mathbb{I}\{c^i=1\} = \begin{cases} 0 & \text{when } c^i \neq 1 \\ 1 & \text{when } c^i = 1 \end{cases}$$

$$p(x_i=1|c=1) = \frac{p(x_i=1, c=1)}{p(c=1)}$$

$$= \frac{\text{number of data with } x_i=1 \text{ and } c=1}{\text{number of data with } c=1}$$

$$= \frac{\sum_{j=1}^N \mathbb{I}\{c^j=1 \text{ and } x_i^j=1\}}{\sum_{j=1}^N \mathbb{I}\{c^j=1\}}$$

$$\text{similarity, } p(x_i=1|c=0) = \frac{p(x_i=1, c=0)}{p(c=0)}$$

$$= \frac{\sum_{j=1}^N \mathbb{I}\{c^j=0 \text{ and } x_i^j=1\}}{\sum_{j=1}^N \mathbb{I}\{c^j=0\}}$$



2. We know that $p(x, c) = p(x|c)p(c) = p(c|x)p(x)$

$$\text{Thus, } p(c|\vec{x}) = \frac{p(\vec{x}, c)}{p(\vec{x})} = \frac{p(\vec{x}, c)}{p(\vec{x}, c=0) + p(\vec{x}, c=1)}$$

$$= \frac{p(c) \prod_{i=1}^N p(x_i|c)}{\sum_{c \in \{0,1\}} p(c) \prod_{i=1}^N p(x_i|c)}$$



3.

① what effect: since this word never appear in training dataset

$$\Rightarrow p(c | \text{"viagra"} = 1 | c) = 0.$$

when doing classification, $p(c | x, \text{"viagra"} = 1)$

$$= \frac{p(c) \prod_{i=1}^N p(x_i | c)}{\sum_{c \in \{0,1\}} p(c) \prod_{i=1}^N p(x_i | c)} = \frac{0}{0+0}$$

which is impossible to deduce.

Also, in this case, $p(c=1 | x, \text{"viagra"} = 1) = p(c=0 | x, \text{"viagra"} = 1)$

★ Therefore, the appearance of a new word will mess up the performance of the classifier, and make the classifier ignore all other words and give an arbitrary prediction.

② how to counter this effect: Laplace smoothing

change the original $p(x_i | c) = \frac{\sum_{i=1}^N \mathbb{I}\{x_i \text{ and } c\}}{\sum_{i=1}^N \mathbb{I}\{c\}}$

into new $p(x_i | c) = \frac{\sum_{i=1}^N \mathbb{I}\{x_i \text{ and } c\} + 1}{\sum_{i=1}^N \mathbb{I}\{c\} + 2}$

Thus, when doing the training, the $p(x_i | c)$ for unseen words would be 50%, which will not effect the performance.

Also, another trick would be: replace all infrequent words into a mark "RARE", and calculate all new unseen words according to the possibility of "RARE".

② How spammer fool NB filter:

add a lot of non-spam-like words in the end of the spam email will reduce its chance of being recognized as spam.

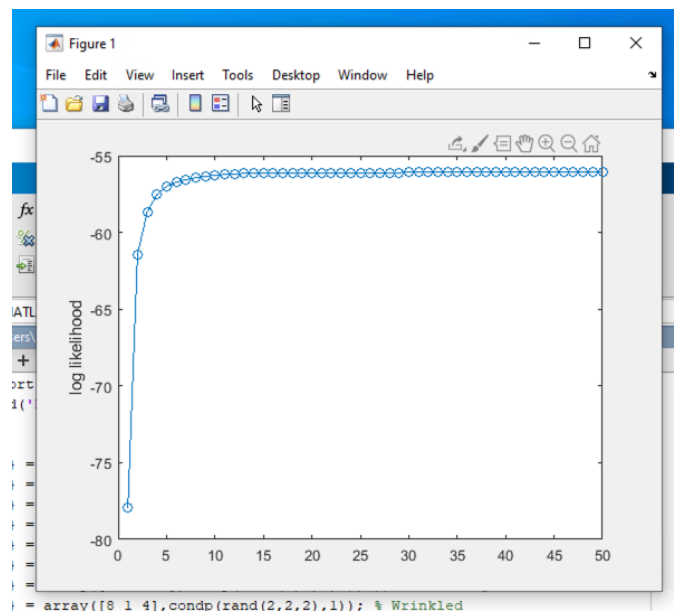
Problem 6

The code is inside the file p6.m, we can get the result of

3.810283e-01

6.189717e-01

Is the probability of drum unit problem, where 0.619 is the possibility of there is a drum unit problem, and 0.381 is the possibility of there is no drum unit problem. From the plotprogress variable of the function EMbeliefnet, we can have the following plot which monitors when the function converges.



It converges well.