

Assignment 3, COMP 540

Chen Zeng(cz39), Zhihui Xie(zx18)

February 16, 2018

1 MAP and MLE parameter estimation

Answer1 According to the properties of Bernoulli Distribution, we can get:

$$P(D|\theta) = \theta^{\sum x^{(i)}} (1 - \theta)^{m - \sum x^{(i)}}$$

According to the principle of maximum likelihood, we can get:

$$\theta_{MLE}^* = \operatorname{argmax}_{\theta} \log P(D|\theta) = \operatorname{argmax}_{\theta} \left(\sum x^{(i)} \log \theta + \left(m - \sum x^{(i)} \right) \log (1 - \theta) \right)$$

Make the derivation of $\log P(D|\theta)$ on θ and make it equal to 0, we can get:

$$\frac{\partial \log P(D|\theta)}{\partial \theta} = \frac{\partial \left(\sum x^{(i)} \log \theta + \left(m - \sum x^{(i)} \right) \log (1 - \theta) \right)}{\partial \theta} = \frac{\sum x^{(i)}}{\theta} - \frac{\left(m - \sum x^{(i)} \right)}{1 - \theta} = 0$$

We can solve the result:

$$\theta = \frac{\sum x^{(i)}}{m}$$

Therefore, here is the final result:

$$\theta_{MLE}^* = \frac{\sum x^{(i)}}{m}$$

This is the estimate of θ .

Answer2 Because of:

$$Beta(\theta|a, b) \propto \theta^{a-1} (1-\theta)^{b-1}$$

Therefore:

$$P(\theta|D) \propto P(D|\theta) P(\theta) \propto \left(\theta^{\sum x^{(i)}} (1-\theta)^{m-\sum x^{(i)}} \right) \left(\theta^{a-1} (1-\theta)^{b-1} \right) = \theta^{\sum x^{(i)} + a - 1} (1-\theta)^{m+b-\sum x^{(i)} - 1}$$

According to the MAP estimation of θ , we can get:

$$\theta_{MAP} = \operatorname{argmax}_{\theta} \log P(\theta|D) = \operatorname{argmax}_{\theta} \left(\left(\sum x^{(i)} + a - 1 \right) \log \theta + \left(m + b - \sum x^{(i)} - 1 \right) \log (1 - \theta) \right)$$

Make the derivation of $\log P(\theta|D)$ on θ and make it equal to 0, we can get:

$$\begin{aligned} \frac{\partial \log P(\theta|D)}{\partial \theta} &= \frac{\partial \left(\left(\sum x^{(i)} + a - 1 \right) \log \theta + \left(m + b - \sum x^{(i)} - 1 \right) \log (1 - \theta) \right)}{\partial \theta} \\ &= \frac{\sum x^{(i)} + a - 1}{\theta} - \frac{m + b - \sum x^{(i)} - 1}{1 - \theta} \\ &= 0 \end{aligned}$$

We can solve the result:

$$\theta = \frac{a + \sum x^{(i)} - 1}{a + b + m - 2}$$

Therefore, here is the final result:

$$\theta_{MAP}^* = \frac{a + \sum x^{(i)} - 1}{a + b + m - 2}$$

When it's under a uniform prior, $a = b = 1$, we can get:

$$\theta_{MAP}^* = \frac{\sum x^{(i)}}{m}$$

Combing the answer from problem1, we can get:

$$\theta_{MLE}^* = \theta_{MAP}^*$$

Therefore, the MAP and MLE estimates of θ are the same under a uniform prior.

2 Logistic regression and Gaussian Naive Bayes

Answer1 We can easily get the results from the logistic regression shown as below:

$$P(y = 1|x) = g(\theta^T x)$$

$$P(y = 0|x) = 1 - g(\theta^T x)$$

They are the expression in terms of the parameter θ and the sigmoid function.

Answer2 At first, we define the const A as below:

$$A = \frac{1}{P(x, y' = 1) + P(x, y' = 0)} = \frac{1}{P(x|y' = 1)\gamma + P(x|y' = 0)(1 - \gamma)}$$

Then, according to the definition and properties of Gaussian distribution and Bernoulli distribution, we have:

$$P(y = 1) = \gamma$$

$$P(y = 0) = 1 - \gamma$$

$$P(x|y = 1) = \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x^{(j)} - \mu_j^1)^2}{2\sigma_j^2}\right)$$

$$P(x|y = 0) = \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x^{(j)} - \mu_j^0)^2}{2\sigma_j^2}\right)$$

Therefore, according to the Bayes rule and Gaussian Naive Bayes model, we can get the following results:

$$P(y = 1|x) = \frac{P(x|y = 1)P(y = 1)}{\sum_{y' \in \{0,1\}} P(x|y')P(y')} = A\gamma \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x^{(j)} - \mu_j^1)^2}{2\sigma_j^2}\right)$$

$$P(y = 0|x) = \frac{P(x|y = 0)P(y = 0)}{\sum_{y' \in \{0,1\}} P(x|y')P(y')} = A(1 - \gamma) \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x^{(j)} - \mu_j^0)^2}{2\sigma_j^2}\right)$$

Here const A is defined at the beginning.

Answer3 Because class 1 and class 0 are equally likely, $\gamma = \frac{1}{2}$. Therefore, to simplify the expression for $P(y = 1|x)$, we can get:

$$P(y = 1|x) = \frac{A}{2} \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x^{(j)} - \mu_j^1)^2}{2\sigma_j^2}\right)$$

For the definition of const A :

$$A = \frac{1}{P(x, y' = 1) + P(x, y' = 0)} = \frac{2}{P(x|y' = 1) + P(x|y' = 0)}$$

Let's divide $P(y = 1|x)$ with $P(y = 0|x)$, using the formation in problem2, we can get:

$$\begin{aligned}
\frac{P(y = 1|x)}{P(y = 0|x)} &= \frac{A\gamma \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x^{(j)} - \mu_j^1)^2}{2\sigma_j^2}\right)}{A(1-\gamma) \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x^{(j)} - \mu_j^0)^2}{2\sigma_j^2}\right)} \\
&= \exp\left(\sum_{j=1}^d \left(-\frac{(x^{(j)} - \mu_j^1)^2}{2\sigma_j^2} + \frac{(x^{(j)} - \mu_j^0)^2}{2\sigma_j^2}\right)\right) \\
&= \exp\left(\sum_{j=1}^d \frac{(\mu_j^0)^2 - (\mu_j^1)^2}{2\sigma_j^2} + \sum_{j=1}^d \frac{x^{(j)} (\mu_j^1 - \mu_j^0)}{\sigma_j^2}\right)
\end{aligned}$$

According to the answer of problem1, we have:

$$\frac{P(y = 1|x)}{P(y = 0|x)} = \frac{g(\theta^T x)}{1 - g(\theta^T x)} = \exp(\theta^T x)$$

These two are equal and comparing them, we can get:

$$\exp\left(\sum_{j=1}^d \frac{(\mu_j^0)^2 - (\mu_j^1)^2}{2\sigma_j^2} + \sum_{j=1}^d \frac{x^{(j)} (\mu_j^1 - \mu_j^0)}{\sigma_j^2}\right) = \exp(\theta^T x)$$

It means we have the following equation:

$$\sum_{j=1}^d \frac{(\mu_j^0)^2 - (\mu_j^1)^2}{2\sigma_j^2} + \sum_{j=1}^d \frac{x^{(j)} (\mu_j^1 - \mu_j^0)}{\sigma_j^2} = \theta^T x$$

Comparing the parameters to x , we can get the appropriate parameterization:

$$\theta_0 = \sum_{j=1}^d \frac{(\mu_j^0)^2 - (\mu_j^1)^2}{2\sigma_j^2}$$

$$\theta_j = \frac{\mu_j^1 - \mu_j^0}{\sigma_j^2} \quad (j = 1, \dots, d)$$

Therefore, it shows that with appropriate parameterization, $P(y = 1|x)$ for Gaussian Naive Bayes with uniform priors is equivalent to $P(y = 1|x)$ for logistic regression.

3 Reject option in classifiers

Answer1 To calculate the minimum risk, we discuss it in the following two cases.

When $i = 1, 2, \dots, C$, we have:

$$\begin{aligned} R(\alpha_i|x) &= \sum_{j=1}^C L(\alpha_i|y=j) P(y_j|x) \\ &= \lambda_s \sum_{j=1, j \neq i}^C P(y=j|x) \\ &= \lambda_s (1 - P(y=i|x)) \end{aligned}$$

When $i = C+1$, we have:

$$R(\alpha_{C+1}|x) = \lambda_r$$

Therefore, if we want to decide $y = j$ to obtain the minimum risk, the following two requirements should be fulfilled:

$$R(\alpha_j|x) \leq R(\alpha_k|x), (j, k = 1, 2, \dots, C \text{ \& } j \neq k)$$

$$R(\alpha_j|x) \leq R(\alpha_{C+1}|x), (j = 1, 2, \dots, C)$$

Combining with the results we get before, we need to solve the following two inequalities:

$$\lambda_s (1 - P(y=j|x)) \leq \lambda_s (1 - P(y=k|x))$$

$$\lambda_s (1 - P(y=j|x)) \leq \lambda_r$$

Finally we get the results:

$$P(y=j|x) \geq P(y=k|x)$$

$$P(y=j|x) \geq 1 - \frac{\lambda_r}{\lambda_s}$$

Therefore, if it fulfilled the above requirements, we decide $y = j$ to obtain the minimum risk.

Answer2 When $\lambda_r = 0$, which makes $\lambda_r/\lambda_s = 0$, we do always reject. As it increases but not equal to 1, the relative cost of rejection increases, we are more likely not to reject, but to assign x to class j . When $\lambda_r/\lambda_s = 1$, we should assign x to class j . Here class j meets with the requirements stated in the last problem.

4 Kernelizing k-nearest neighbors

Answer According to the definition of the dot product of vector x_1 and vector x_2 from wikipedia, we have:

$$\langle x_1, x_2 \rangle = x_1 \cdot x_2 = \sum_{i=1}^d x_1^{(i)} x_2^{(i)}$$

To re-express this classification rule in terms of dot products, we have:

$$d = \sqrt{\|x_1 - x_2\|^2} = \sqrt{\sum_{i=1}^d (x_1^{(i)} - x_2^{(i)})^2} = \sqrt{\sum_{i=1}^d x_1^{(i)} x_1^{(i)} + \sum_{i=1}^d x_2^{(i)} x_2^{(i)} - 2 \sum_{i=1}^d x_1^{(i)} x_2^{(i)}} = \sqrt{\langle x_1, x_1 \rangle + \langle x_2, x_2 \rangle - 2 \langle x_1, x_2 \rangle}$$

To make use of the kernel trick to formulate the k-nearest neighbor algorithm, we can just replace the dot products in the above equation with kernel function K and then we can get:

$$d = \sqrt{K(x_1, x_1) + K(x_2, x_2) - 2K(x_1, x_2)}$$

Instead of using the Euclidean distance to calculate the relation between two vectors, we use this function to calculate it. Here is the concrete algorithm:

Algorithm 1 Kernel Trick Based K Nearest Neighbour

```

1: Inputs:
    $S = \{s_i\} = \{(x_i, y_i, d_i)\}$ 
2: Initialize:
    $d_i = 0, i = 1, \dots, n$ 
3: for  $i = 1$  to  $n$  do
4:    $d_i = \sqrt{K(x_i, x_i) + K(x, x) - 2K(x, x_i)}$ 
5: end for
6:  $S' = \{s'_j\} = \{(x'_j, y'_j, d'_j)\} = (\text{sort}(s_i.d_i)) [1 : k], i = 1, \dots, n, j = 1, \dots, k$ 
7:  $cnt1 = 0$ 
8:  $cnt2 = 0$ 
9: for  $j = 1$  to  $k$  do
10:  if  $s'_j.y'_j = 1$  then
11:     $cnt1 = cnt1 + 1$ 
12:  else
13:     $cnt2 = cnt2 + 1$ 
14:  end if
15: end for
16:  $label = cnt1 > cnt2 ? 1 : -1$ 

```

5 Constructing kernels

Answer1 According to the Mercer's theorem, if k_1 is a valid kernel, then the Gram matrix K_1 whose elements are $k_1(x, x'), 1 \leq i, j \leq m$ is positive definite. According to the definition

of positive definite, the scalar $z^T K_1 z$ should always be positive for every non zero column vector z . Because we have $k(x, x'), 1 \leq i, j \leq m = ck_1(x, x'), 1 \leq i, j \leq m$, we can know the Gram matrix for k is $K = cK_1$. So we have $z^T K z = cz^T K_1 z$. It should also be positive and K should be positive definite if $c > 0$. Therefore, according to the Mercer's theorem, under the case of $c > 0$, k should also be valid kernels.

Answer2 Here $k(x, x')$ may not be valid kernel and it depends on $f(x)$. For example, if $f(x) = 0$, then K will be a zero matrix and it's not positive definite. Therefore, according to the Mercer's theorem, under this case, k is not a valid kernel.

However, if $f(x)$ is not always equal to 0, k is a valid kernel. Here is the proof:

According to the Mercer's theorem, because k_1 is a valid kernel, we have $k_1(x, x') = \langle \phi(x), \phi(x') \rangle$.

Therefore, we can make the following transformation:

$$\begin{aligned} k(x, x') &= f(x) k_1(x, x') f(x') \\ &= f(x) f(x') \langle \phi(x), \phi(x') \rangle \\ &= \langle f(x) \phi(x), f(x') \phi(x') \rangle \\ &= \langle \phi'(x), \phi'(x') \rangle \end{aligned}$$

Here we define:

$$\phi'(x) = f(x) \phi(x)$$

Therefore, we have:

$$k(x, x') = \langle \phi'(x), \phi'(x') \rangle$$

According to the Mercer's theorem, k is a valid kernel

Answer3 This is similar to answer1. Here $K = K_1 + K_2$ and $z^T K z = z^T K_1 z + z^T K_2 z > 0$ for every non zero column vector z . According to the Mercer's theorem, K is positive definite and therefore k should also be valid kernels.

6 One-vs-all logistic regression

Answer

7 Softmax regression

Answer 7.1 CIFAR-10 database has 10 labels and θ matrix is randomly initialized, so $\frac{\exp(\theta^{(k)T} x^{(i)})}{\sum_{j=1}^K \exp(\theta^{(j)T} x^{(i)})}$

is approximately 0.1, which means $-\log(\frac{\exp(\theta^{(k)T} x^{(i)})}{\sum_{j=1}^K \exp(\theta^{(j)T} x^{(i)})})$ is $-\log(0.1)$. Because $\lambda = 0$, $J(\theta)$ is about $-\log(0.1)$.

Classes	plane	car	bird	cat	deer	dog	frog	horse	ship	truck
OVA	0.465	0.463	0.193	0.161	0.234	0.272	0.457	0.406	0.546	0.42
Softmax Regression	0.465	0.486	0.242	0.266	0.299	0.347	0.502	0.425	0.621	0.459

Answer 7.9

Softmax regression classifier performs better than OVA binary logistic regression. I would recommend Softmax regression classifier for the CIFAR-10 classification problem. For OVA method, we need to change multiple classes into two classes and we will build K logistic classifiers if it has K classes. In this problem, these 10 classes are mutually exclusive, so softmax regression classifier would be appropriate and it is faster. What's more, one sample may be classified into several classes or non-class.