

# CSE 517A

## Machine Learning

### THW2

#### Problems

**1. (30 points) Naïve Bayes**

- (a) We have  $2^2$  different  $\vec{x}$ s under 2 different categories  $y$ s, so in total we need  $2^2 * 2 = 8$  parameters for  $P(X|Y)$ .
- (b) This time we need  $2^d * 2 = 2^{d+1} = 2^{101}$  parameters.
- (c) If we make the naïve assumption, for each dimension  $X_i$  we need only to estimate  $P(X_i|Y = 0)$  and  $P(X_i|Y = 1)$ , therefore in total we need  $2d$  parameters, which is much smaller than that of not using naïve assumption ( $2^{d+1}$ ).  
Generally speaking, the naïve assumption can be seen as an approximation method which makes our problem easier to solve. So when the dimension is high, we'd better use it or the problem is almost intractable for having too much (exponential # of) parameters to estimate. However, in lower dimension where we are able to estimate all the parameters, we can get more precise result by not using it.

## 2. (30 points) Naïve Bayes, Part II

(a) Use  $\vec{x}$  as feature vector of document, we have:

$$\begin{aligned}
 h(d) &= \underset{y}{\operatorname{argmax}} P(Y = y) \prod_{i=1}^l P(W_i = w_i | Y = y) \\
 &= \underset{y}{\operatorname{argmax}} P(Y = y) \prod_{i=1}^d P(X_i = x_i | Y = y) \\
 &= \underset{y}{\operatorname{argmax}} P(Y = y) \prod_{i=1}^d P(a_i \in d | Y = y)^{x_i} \\
 &= \underset{y}{\operatorname{argmax}} \log \left( P(Y = y) \prod_{i=1}^d P(a_i \in d | Y = y)^{x_i} \right) \\
 &= \underset{y}{\operatorname{argmax}} \log P(Y = y) + \sum_{i=1}^d \log [P(a_i \in d | Y = y)] x_i \\
 &= \operatorname{sign} \left[ \log P(Y = 1) + \sum_{i=1}^d \log [P(a_i \in d | Y = 1)] x_i - \log P(Y = -1) - \sum_{i=1}^d \log [P(a_i \in d | Y = -1)] x_i \right] \\
 &= \operatorname{sign} \left[ \sum_{i=1}^d \log \frac{P(a_i \in d | Y = 1)}{P(a_i \in d | Y = -1)} x_i + \log \frac{P(Y = 1)}{P(Y = -1)} \right]
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 \vec{v} &= \left( \log \frac{P(a_1 \in d | Y = 1)}{P(a_1 \in d | Y = -1)}, \log \frac{P(a_2 \in d | Y = 1)}{P(a_2 \in d | Y = -1)}, \dots, \log \frac{P(a_d \in d | Y = 1)}{P(a_d \in d | Y = -1)} \right) \\
 b &= \log \frac{P(Y = 1)}{P(Y = -1)}
 \end{aligned}$$

(b) (i) Assuming we have in total  $n$  samples,  $X_i$  has  $K_i$  different categories and  $Y$  has  $K_y$  different classes,

$$\begin{aligned}
 Pr(X_i | Y) &= \prod_{j=1}^{K_i} Pr(X_i = j | y)^{I(X_i=j)} \\
 Pr(Y = c) &= \frac{\sum_{i=1}^K I(y_i = c)}{K}
 \end{aligned}$$

(ii) The MLE estimators for Gaussian distribution are:

$$\begin{aligned}
 \mu_i^* &= \bar{X}_i = \frac{1}{N} \sum_{j=1}^N X_{ij} \\
 \sigma_i^* &= \frac{1}{N}
 \end{aligned}$$

(iii) We have 2 classes, and for each we need to estimate  $d$  means,  $d$  variances and 1 prior  $p(Y)$ , so there should be  $2 * (d + d + 1) = 4d + 2$  parameters.

(c)

$$\begin{aligned}
 \because P(X|Y = 1) &= \prod_{i=1}^d P(X_i | Y = 1) \\
 &= \prod_{i=1}^d \frac{1}{\sqrt{2\pi}\sigma_i} \exp \left( -\frac{(X_i - \mu_{1,i})^2}{2\sigma_i^2} \right)
 \end{aligned}$$

$$\begin{aligned}
&= \left( \frac{1}{\sqrt{2\pi}} \right)^d \exp \left( \sum_{i=1}^d \frac{1}{2\sigma_i^2} (X_i - \mu_1)^2 \right) \prod_{i=1}^d \frac{1}{\sigma_i} \\
\therefore P(Y = +1|X) &= \frac{P(X|Y = +1)P(Y = +1)}{P(X)} \\
&= \frac{P(Y = +1) \left( \frac{1}{\sqrt{2\pi}} \right)^d \exp \left( \sum_{i=1}^d \frac{1}{2\sigma_i^2} (X_i - \mu_1)^2 \right) \prod_{i=1}^d \frac{1}{\sigma_i}}{P(Y = +1) \left( \frac{1}{\sqrt{2\pi}} \right)^d \exp \left( \sum_{i=1}^d \frac{1}{2\sigma_i^2} (X_i - \mu_1)^2 \right) \prod_{i=1}^d \frac{1}{\sigma_i} + P(Y = -1) \left( \frac{1}{\sqrt{2\pi}} \right)^d \exp \left( \sum_{i=1}^d \frac{1}{2\sigma_i^2} (X_i - \mu_{(-1)})^2 \right) \prod_{i=1}^d \frac{1}{\sigma_i}} \\
&= \frac{1}{1 + \frac{P(Y = -1)}{P(Y = +1)} \exp \left( \sum_{i=1}^d \left[ \frac{1}{\sigma_i^2} (\mu_1 - \mu_{(-1)}) X_i + \frac{1}{2\sigma_i^2} (\mu_{(-1)}^2 - \mu_1^2) \right] \right)} \\
&= \frac{1}{1 + \exp \left( \log P(Y = -1) - \log P(Y = +1) + \sum_{i=1}^d \frac{1}{2\sigma_i^2} (\mu_{(-1)}^2 - \mu_1^2) + \sum_{i=1}^d \frac{1}{\sigma_i^2} (\mu_{(-1)} - \mu_1) X_i \right)}
\end{aligned}$$

As a result,

$$w_i = \frac{1}{\sigma_i^2} (\mu_1 - \mu_{(-1)}), i = 1, 2, \dots, d$$

$$w_0 = \log P(Y = +1) - \log P(Y = -1) + \sum_{i=1}^d \frac{1}{2\sigma_i^2} (\mu_{(-1)}^2 - \mu_1^2) = \log \frac{P(Y = +1)}{P(Y = -1)} + \sum_{i=1}^d \frac{1}{2\sigma_i^2} (\mu_{(-1)}^2 - \mu_1^2)$$

**3. (40 points) Valid Kernels, Kernel Construction**

(a) Here  $D = 10$ . Define kernel function  $K(\cdot, \cdot)$ :

$$\begin{aligned}
 K(\mathbf{x}_i, \mathbf{x}_j) &= (1 + \mathbf{x}_i^T \mathbf{x}_j)^2 \\
 &= \left(1 + \sum_{d=1}^3 (x_i)_d (x_j)_d\right)^2 \\
 &= 1 + 2 \sum_{d=1}^3 (x_i)_d (x_j)_d + 2 \sum_{d=1}^3 (x_i)_d^2 (x_j)_d^2 + 2(x_i)_1 (x_i)_2 (x_j)_1 (x_j)_2 + 2(x_i)_2 (x_i)_3 (x_j)_2 (x_j)_3 \\
 &\quad + 2(x_i)_1 (x_i)_3 (x_j)_1 (x_j)_3 \\
 &= \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle
 \end{aligned}$$

We can see that transforming  $\mathbf{x}$  to  $\Phi(\mathbf{x})$  needs 12 productions, and compute  $\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$  need additional 10 productions and 9 summations, thus in total we need 22 productions and 9 summations.

However, compute  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$  need only 4 productions and 1 summation and can get the same result, which is highly efficient.

(b) We can judge from the eigenvalues of matrices.

(i) For  $A_1 = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ , its eigen vectors and eigenvalues are:

$$\mathbf{v}_1 = (1, 1)^T, \lambda_1 = 2 \geq 0$$

$$\mathbf{v}_2 = (-1, 1)^T, \lambda_2 = 0 \geq 0$$

Therefore,  $A_1$  is positive semidefinite.

(ii) For  $A_2 = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 0 \\ 1 & 0 & 2 \end{bmatrix}$ , its eigenvectors and eigenvalues are:

$$\mathbf{v}_1 = (\sqrt{2}, 1, 1)^T, \lambda_1 = 2 + \sqrt{2} > 0$$

$$\mathbf{v}_2 = (0, -1, -1)^T, \lambda_2 = 2 > 0$$

$$\mathbf{v}_3 = (-\sqrt{2}, -1, 1)^T, \lambda_3 = 2 - \sqrt{2} > 0$$

Therefore,  $A_2$  is positive definite.

(iii) For  $A_3 = \begin{bmatrix} 2 & 1 & -1 \\ 1 & 1 & 1 \\ -1 & 1 & 2 \end{bmatrix}$ , its eigenvectors and eigenvalues are:

$$\mathbf{v}_1 = (-1, 0, 1)^T, \lambda_1 = 3 > 0$$

$$\mathbf{v}_2 = \left(1, \frac{-2 + \sqrt{2}}{-1 + \sqrt{2}}, 1\right)^T, \lambda_2 = 1 + \sqrt{2} > 0$$

$$\mathbf{v}_3 = \left(1, -\frac{2 + \sqrt{2}}{1 + \sqrt{2}}, 1\right)^T, \lambda_3 = 1 - \sqrt{2} < 0$$

Therefore,  $A_1$  is neither positive definite or semidefinite.

(c) (i) Assuming A and B are both  $n \times n$  matrices, since A and B are both positive semi-definite matrices, we have:

$$\forall \lambda \in \mathbb{R}^n, \lambda^T A \lambda \geq 0, \lambda^T B \lambda \geq 0.$$

Therefore,

$$\forall \lambda \in R^n, \lambda^T(A+B)\lambda = \lambda^T A \lambda + \lambda^T B \lambda \geq 0$$

Thus (A+B) is positive semi-definite.

(ii) Using property of Kronecker product <sup>[1]</sup>:

Suppose that A and B are square matrices of size n and m respectively. Let  $\lambda_1, \dots, \lambda_n$  be the eigenvalues of A and  $\mu_1, \dots, \mu_m$  be those of B (listed according to multiplicity). Then the eigenvalues of  $A \otimes B$  are

$$\lambda_i \mu_j, \quad i = 1, \dots, n, j = 1, \dots, m.$$

Since A and B are both positive semi-definite, we know:

$$\lambda_1, \dots, \lambda_n \geq 0, \quad \mu_1, \dots, \mu_m \geq 0$$

Therefore,  $\lambda_i \mu_j \geq 0, i = 1, \dots, n, j = 1, \dots, m$ , i.e. all eigenvalues of  $A \otimes B$  nonnegative. As a result  $A \otimes B$  is positive semi-definite.

(d) By Taylor expansion:

$$\begin{aligned} K(x_i, x_j) &= \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) \\ &= \sum_{n=0}^{\infty} \frac{(x_i^T x_j)^n}{n! (2\sigma)^n} \exp\left(-\frac{1}{2\sigma^2} \|x_i\|^2\right) \exp\left(-\frac{1}{2\sigma^2} \|x_j\|^2\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} \|x_i\|^2\right) \exp\left(-\frac{1}{2\sigma^2} \|x_j\|^2\right) \sum_{n=0}^{\infty} \frac{(x_i^T x_j)^n}{n! (2\sigma)^n} \\ &= \Phi(x_i)^T \Phi(x_j) \end{aligned}$$

Where

$$\Phi(x) = \exp\left(-\frac{1}{2\sigma^2} \|x\|^2\right) \left[1, \frac{x}{1! (2\sigma)^1}, \dots, \frac{x^n}{n! (2\sigma)^n}, \dots\right]$$

Which is a vector in infinite dimensional space.

References:

[1] [Kronecker product - Wikipedia](#)

**4. (30 points) Kernelize the Perceptron Algorithm**

(a) Note where  $\alpha_i$  is the number of times  $\mathbf{x}_i$  was misclassified, then:

$$\mathbf{w} = \sum_i^n \alpha_i y_i \mathbf{x}_i$$

(b) From (a):

$$\begin{aligned} h(\mathbf{x}) &= \text{sign}(\mathbf{w}^T \mathbf{x}) \\ &= \text{sign} \left[ \left( \sum_i^n \alpha_i y_i \mathbf{x}_i \right)^T \mathbf{x} \right] \\ &= \text{sign} \left[ \sum_{i=1}^n \alpha_i y_i (\mathbf{x}_i^T \mathbf{x}) \right] \end{aligned}$$

If we use feature transformation  $\Phi(\cdot)$ , and kernel function  $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$ , then:

$$h(\mathbf{x}) = \text{sign} \left[ \sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) \right]$$

(c) The kernelized version of perceptron is:

- (0) Initialize  $\mathbf{w} = \mathbf{0}$   
REPEAT until convergence:
- (1) Pick  $(\mathbf{x}_j, y_j)$  randomly from  $D$
- (2) If  $h(\mathbf{x}_j) = \text{sign}[\sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_j)] \neq y_j$ , update  $\alpha_j$ :  $\alpha_j \leftarrow \alpha_j + 1$