

1) Show that a Kernel function $K(x_1, x_2)$ satisfies the Cauchy-Schwarz generalization $K(x_1, x_2)^2 \leq K(x_1, x_1)K(x_2, x_2)$.

↪ Hint: $|u^T v|^2 \leq \|u\|^2 \|v\|^2 \quad \forall u, v \in \mathbb{R}^n$

Consider a feature mapping $\phi: \mathbb{R}^m \rightarrow \mathbb{R}^n$.

Then, the corresponding kernel is defined as $K(x_1, x_2) = \phi(x_1)^T \phi(x_2) \quad \forall x_1, x_2 \in \mathbb{R}^m$.

we have $K(x_1, x_2)^2 = (\phi(x_1)^T \phi(x_2))^2$

$$= |\phi(x_1)^T \phi(x_2)|^2 \quad \text{b/c } x^2 = \|x\|^2 \quad \forall x \in \mathbb{R}$$

$$\leq \|\phi(x_1)\|^2 \|\phi(x_2)\|^2 \quad \text{by Cauchy-Schwarz}$$

$$= \phi(x_1)^T \phi(x_1) \phi(x_2)^T \phi(x_2) \quad \text{b/c } u^T u = \|u\|^2$$

$$= K(x_1, x_1) K(x_2, x_2) \quad \text{by definition.}$$

$$\therefore K(x_1, x_2)^2 \leq K(x_1, x_1) K(x_2, x_2).$$

□

Zack
Beyer

[2] Given $K_1(x, x')$ and $K_2(x, x')$ are valid. Prove the following is valid.

a) $K(x, x') = K_1(x, x') + K_2(x, x')$.

$K_1(x, x')$ and $K_2(x, x')$ are valid kernels, and so are symmetric.

Thus, $K(x, x') = K_1(x, x') + K_2(x, x') = K_1(x', x) + K_2(x', x) = K(x', x)$.

$\therefore K(x, x')$ is symmetric.

The kernel matrix K for $K(x, x')$ is defined as $K = K_1 + K_2$.

Consider some arbitrary vector y . Then,

$$\begin{aligned} y^T K y &= y^T (K_1 + K_2) y \\ &= y^T K_1 y + y^T K_2 y \\ &\geq 0 + 0 \\ &= 0 \end{aligned}$$

$\therefore K$ is positive semidefinite.

Hence, $K(x, x')$ is a valid kernel.

Kernel matrices for $K_1(x, x')$ and $K_2(x, x')$.

Generally, to prove K is valid, must show it is symmetric and its matrix is positive Semidefinite.

□

Pointwise multiplication

b) $K(x, x') = K_1(x, x') K_2(x, x')$. The Kernel matrix K is defined as $K = K_1 \circ K_2$.

Positive semi-definite:

Consider some arbitrary vector y .

$$\begin{aligned} y^T K y &= y^T (K_1 \circ K_2) y \\ &= y^T K_1 y \circ y^T K_2 y \\ &\geq 0 \circ 0 \\ &= 0 \end{aligned}$$

$\therefore K$ is positive semidefinite.

Hence, $K(x, x')$ is a valid kernel.

Symmetry:

Also, since $K_1(x, x')$ and $K_2(x, x')$ are symmetric,

$$\begin{aligned} K(x, x') &= K_1(x, x') K_2(x, x') \\ &= K_1(x', x) K_2(x', x) \\ &= K(x', x) \end{aligned}$$

$\therefore K$ is symmetric.

□

$$c) K(x, x') = \exp(K_1(x, x')).$$

Suppose K_1, K_2 are valid kernels.

Claim ①: If $\alpha, \beta \geq 0$, then $K(x, x') := \alpha K_1(x, x') + \beta K_2(x, x')$ is a valid kernel.

Claim ②: If f is a polynomial with positive coefficients, then $K(x, x') := f(K_1(x, x'))$ is a valid kernel.

By definition, $\exp(x) = \sum_{i=0}^{\infty} \frac{x^i}{i!}$ is a polynomial with positive coefficients.

Then, $\exp(K_1(x, x'))$ is equivalent to writing an infinite polynomial with positive coefficients as a function of $K_1(x, x')$.

By claim ②, $K = \exp(K_1(x, x'))$ is a valid kernel.

□

Proof of claim ①: Similar to part a, except $K = \alpha K_1 + \beta K_2$ and we note $y^T K y = \alpha y^T K_1 y + \beta y^T K_2 y \geq 0$.

□

Proof of claim ②: Each polynomial term is a product of kernels multiplied by some positive coefficient.

By ①, each term is a valid kernel.

By ①, the sum of the terms, is a valid kernel, multiplied by a coefficient

□

>Show that the parameter b can be determined using

$$b = \frac{1}{N_m} \sum_{n \in M} (y^{(n)} - \sum_{m \in S} \alpha_m y^{(m)} \langle x^{(n)}, x^{(m)} \rangle) \quad M = \text{index set for data st. } 0 < d_n < C \\ S = \text{index set for data st. } d_n \neq 0.$$

Consider Primal problem of Soft SVM.

$$\text{The Lagrangian is } \mathcal{L} = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i(w^T x_i + b)) - \sum_{i=1}^m \beta_i \xi_i$$

→ Here, α_i and β_i denote Lagrange multipliers. $\alpha_i, \beta_i \geq 0$.

$$\text{Taking derivatives, } \frac{\partial \mathcal{L}}{\partial w} = \frac{1}{2} \cdot 2w - \sum_{i=1}^m \alpha_i y_i x_i$$

$$\text{Taking } \frac{\partial \mathcal{L}}{\partial w} = 0, \text{ obtain } w - \sum_{i=1}^m \alpha_i y_i x_i = 0 \leftrightarrow w = \sum_{i=1}^m \alpha_i y_i x_i \text{ at the desired optimum.}$$

Let S denote the set of indices of data points having $d_i \neq 0$.

$$\text{Then, } w = \sum_{i=1}^m \alpha_i y_i x_i = \underbrace{\sum_{m \in S} \alpha_m y_m x_m}_{} = w$$

Now, let M denote the index set of points x_i for which $0 < d_i < C$.

↪ Recall the complementary slackness conditions at the optimum,

$$1. \alpha_i (y_i(w^T x_i + b) - 1 + \xi_i) = 0$$

$$2. \beta_i \xi_i = 0$$

Also note that $\frac{\partial \mathcal{L}}{\partial \xi_i} = (-\alpha_i - \beta_i) = 0$ at the optimum (KKT condition).

Because $0 < d_i < C$ and $C - d_i - \beta_i = 0$, $0 < \beta_i < C$.

Because $0 < \beta_i$ and $\beta_i \xi_i = 0$, $\underline{\xi_i = 0}$.

Because $\alpha_i > 0$ and $\alpha_i (y_i(w^T x_i + b) - 1 + \xi_i) = 0$, $y_i(w^T x_i + b) - 1 + \xi_i = 0$.

Because $\xi_i = 0$ and $y_i(w^T x_i + b) - 1 + \xi_i = 0$, $y_i(w^T x_i + b) - 1 = 0$.

$$\Rightarrow y_i(w^T x_i + b) = 1$$

By the above argument, for points x_i whose index is in M , $y_i(w^T x_i + b) = 1$.

(i.e. $0 < d_i < C$)



We can algebraically solve for b for any point whose index $\in \mathcal{M}$:

$$y_i(w^\top x_i + b) = 1$$

$$\downarrow \\ w^\top x_i + b = \frac{1}{y_i}$$

$$\downarrow \\ w^\top x_i + b = y_i \quad (b/c \quad y_i \in \{-1, 1\}, \quad \frac{1}{y_i} = y_i)$$

$$\downarrow \\ b = y_i - w^\top x_i$$

Using $w = \sum_{m \in S} \alpha_m y_m x_m$ as before, $b = y_i - (\sum_{m \in S} \alpha_m y_m x_m)^\top x_i$

$$\Rightarrow b = y_i - \sum_{m \in S} \alpha_m y_m \langle x_i, x_m \rangle$$

If we average b over all vectors whose index is in the set S , then we obtain the desired result, assuming we found optimal α :

$$b = \frac{1}{N_S} \sum_{n \in S} \left(y_n - \sum_{m \in S} \alpha_m y_m \langle x_n, x_m \rangle \right)$$

□

4) Consider RVs A, B, C with listed joint $P(A, B, C)$. Marginals computed by summing rows / columns...

		$C=0$	$C=1$
		$B=0$	$B=1$
$A=0$	$B=0$	0.096	0.024
	$B=1$	0.24	0.03
$A=1$	$B=0$	0.224	0.056
	$B=1$	0.27	0.03

$$P(A=0) = 0.42$$

$$P(A=1) = 0.58$$

$$P(C=0) = 0.4 \quad P(C=1) = 0.6$$

$$P(B=0) = 0.86 \quad P(B=1) = 0.14$$

$$a) P(A|C=0) = \frac{P(A, C=0)}{P(C=0)} = \begin{cases} \rightarrow A=0: \frac{(0.096 + 0.024)}{0.40} = 0.3 \\ \rightarrow A=1: \frac{(0.224 + 0.056)}{0.40} = 0.7 \end{cases}$$

$$P(B|C=0) = \frac{P(B, C=0)}{P(C=0)} = \begin{cases} \rightarrow B=0: \frac{(0.096 + 0.224)}{0.40} = 0.8 \\ \rightarrow B=1: \frac{(0.024 + 0.056)}{0.40} = 0.2 \end{cases}$$

$$P(A, B|C=0) = \begin{cases} \rightarrow A=0, B=0: \frac{0.096}{0.40} = 0.24 & \rightarrow A=1, B=0: \frac{0.224}{0.40} = 0.56 \\ \rightarrow A=0, B=1: \frac{0.024}{0.40} = 0.06 & \rightarrow A=1, B=1: \frac{0.056}{0.40} = 0.14 \end{cases}$$

$$P(A|C=0) = 0.3 \text{ if } A=0, 0.7 \text{ if } A=1$$

$$P(B|C=0) = 0.8 \text{ if } B=0, 0.2 \text{ if } B=1$$

$$P(A, B|C=0) = \begin{cases} 0.24 \text{ if } A=B=0, 0.14 \text{ if } A=B=1, \\ 0.06 \text{ if } A=0, B=1, \\ 0.56 \text{ if } A=1, B=0 \end{cases}$$

$$b) P(A|C=1) = \begin{cases} A=0: & \frac{(0.27+0.03)}{0.6} = \frac{1}{2} \\ A=1: & \frac{(0.27+0.03)}{0.6} = \frac{1}{2} \end{cases}$$

$$P(B|C=1) = \begin{cases} B=0: & \frac{(0.27+0.27)}{0.6} = 0.9 \\ B=1: & \frac{(0.03+0.03)}{0.6} = 0.1 \end{cases}$$

$$P(A, B|C=1) = \begin{cases} A=0, B=0: & \frac{0.27}{0.60} = 0.45 \\ A=0, B=1: & \frac{0.03}{0.60} = 0.05 \\ A=1, B=0: & \frac{0.27}{0.60} = 0.45 \\ A=1, B=1: & \frac{0.03}{0.60} = 0.05 \end{cases}$$

$P(A C=1) = 0.5$ if $A=0$, 0.5 if $A=1$
$P(B C=1) = 0.9$ if $B=0$, 0.1 if $B=1$
$P(A, B C=1) = 0.45$ if $A=B=0$, 0.05 if $A=B=1$, 0.05 if $A=C, B=1$ 0.45 if $A=1, B=0$

c) 8 cases to check to see if A conditionally independent of B given C.

- $P(A=0, B=0|C=0) = 0.24 = 0.24 = P(A=0|C=0)P(B=0|C=0)$
- $P(A=0, B=1|C=0) = 0.06 = 0.06 = P(A=0|C=0)P(B=1|C=0)$
- $P(A=1, B=0|C=0) = 0.56 = 0.56 = P(A=1|C=0)P(B=0|C=0)$
- $P(A=1, B=1|C=0) = 0.14 = 0.14 = P(A=1|C=0)P(B=1|C=0)$
- $P(A=0, B=0|C=1) = 0.45 = 0.45 = P(A=0|C=1)P(B=0|C=1)$
- $P(A=0, B=1|C=1) = 0.05 = 0.05 = P(A=0|C=1)P(B=1|C=1)$
- $P(A=1, B=0|C=1) = 0.45 = 0.45 = P(A=1|C=1)P(B=0|C=1)$
- $P(A=1, B=1|C=1) = 0.05 = 0.05 = P(A=1|C=1)P(B=1|C=1)$

$\Rightarrow A$ is conditionally independent of B given C Yes

$$d) P(A) = \begin{cases} \rightarrow A=0: 0.42 \\ \rightarrow A=1: 0.58 \end{cases} \quad P(B) = \begin{cases} \rightarrow B=0: 0.86 \\ \rightarrow B=1: 0.14 \end{cases}$$

$$P(A, B) = \begin{cases} \rightarrow A=0, B=0: 0.27 + 0.096 = 0.366 \\ \rightarrow A=1, B=1: 0.056 + 0.03 = 0.086 \\ \rightarrow A=0, B=1: 0.024 + 0.03 = 0.054 \\ \rightarrow A=1, B=0: 0.224 + 0.27 = 0.494 \end{cases}$$

$P(A) = 0.42 \text{ if } A=0, 0.58 \text{ if } A=1$
$P(B) = 0.86 \text{ if } B=0, 0.14 \text{ if } B=1$
$P(A, B) = 0.366 \text{ if } A, B=0, 0.086 \text{ if } A, B=1$
$0.054 \text{ if } A=0, B=1$
$0.494 \text{ if } A=1, B=0$

e) 4 cases to consider to see if A is independent of B.

- $P(A=0, B=0) = 0.366 \neq 0.3612 = P(A=0)P(B=0)$

\Rightarrow Since $P(A=0, B=0) \neq P(A=0)P(B=0)$, A is NOT independent of B. NO

5] Use Naive Bayes to decide if 9, 10 good or not.

a) calculate MLE of $P(G)$, $P(x|G)$ for $x \in \{0, 1, C, A\}$.

$$\hookrightarrow P(G=1) = \frac{6}{8}$$

$$P(O=1|G=1) = \frac{P(O=1, G=1)}{P(G=1)} = \frac{\frac{3}{8}}{\frac{6}{8}} = \frac{3}{6} \quad P(O=1|G=0) = \frac{2}{2}$$

$$P(B=1|G=1) = \frac{2}{6} \quad P(B=1|G=0) = \frac{2}{2}$$

$$P(C=1|G=1) = \frac{3}{6} \quad P(C=1|G=0) = \frac{1}{2}$$

$$P(A=1|G=1) = \frac{5}{6} \quad P(A=1|G=0) = \frac{0}{2}$$

b) Decide on #9, #10 with $G_i = \operatorname{argmax}_{G \in \{0, 1\}} P(O_i, B_i, C_i, A_i | G_i)$.

\hookrightarrow use conditional independence to write

$$P(O_i, B_i, C_i, A_i | G_i) = P(O_i | G_i) P(B_i | G_i) P(C_i | G_i) P(A_i | G_i)$$

• For sample #9, $(O=0, B=1, C=0, A=1)$

$$\begin{aligned} \rightarrow G=0: & P(G=0) P(O=0|G=0) P(B=1|G=0) P(C=0|G=0) P(A=1|G=0) \\ & = \left(\frac{2}{8}\right) \left(\frac{0}{2}\right) \left(\frac{2}{2}\right) \left(\frac{1}{2}\right) \left(\frac{0}{2}\right) = 0 \end{aligned}$$

$$\begin{aligned} \rightarrow G=1: & P(G=1) P(O=0|G=1) P(B=1|G=1) P(C=0|G=1) P(A=1|G=1) \\ & = \left(\frac{6}{8}\right) \left(\frac{3}{6}\right) \left(\frac{2}{6}\right) \left(\frac{3}{6}\right) \left(\frac{5}{6}\right) = \frac{5}{96} \end{aligned}$$

Since $\frac{5}{96} > 0$, $G=1$

• For sample #10, $(O=1, B=1, C=1, A=1)$

$$\begin{aligned} \rightarrow G=0: & P(G=0) P(O=1|G=0) P(B=1|G=0) P(C=1|G=0) P(A=1|G=0) \\ & = \left(\frac{2}{8}\right) \left(\frac{2}{2}\right) \left(\frac{2}{2}\right) \left(\frac{1}{2}\right) \left(\frac{0}{2}\right) = 0 \end{aligned}$$

$$\begin{aligned} \rightarrow G=1: & P(G=1) P(O=1|G=1) P(B=1|G=1) P(C=1|G=1) P(A=1|G=1) \\ & = \left(\frac{6}{8}\right) \left(\frac{3}{6}\right) \left(\frac{2}{6}\right) \left(\frac{3}{6}\right) \left(\frac{5}{6}\right) = \frac{5}{96} \end{aligned}$$

Since $\frac{5}{96} > 0$, $G=1$

Sample #9: $G=1$

Sample #10: $G=1$

c) Redo part (a) using Laplace Smoothing for $P(X|G)$.

$$\begin{aligned} \cdot P(O=1|G=1) &= \frac{3+1}{6+2} = \frac{4}{8} & \therefore P(O=1|G=0) &= \frac{2+1}{2+2} = \frac{3}{4} \\ \cdot P(B=1|G=1) &= \frac{2+1}{6+2} = \frac{3}{8} & \therefore P(B=1|G=0) &= \frac{2+1}{2+2} = \frac{3}{4} \\ \cdot P(C=1|G=1) &= \frac{3+1}{6+2} = \frac{4}{8} & \therefore P(C=1|G=0) &= \frac{1+1}{2+2} = \frac{2}{4} \\ \cdot P(A=1|G=1) &= \frac{5+1}{6+2} = \frac{6}{8} & \therefore P(A=1|G=0) &= \frac{0+1}{2+2} = \frac{1}{4} \end{aligned}$$

d) Repeat (b) with data from (c).

• For sample #9...

$$\begin{aligned} \rightarrow G=0: & P(G=0) P(O=0|G=0) P(B=1|G=0) P(C=0|G=0) P(A=1|G=0) \\ &= \left(\frac{2}{8}\right) \left(\frac{1}{4}\right) \left(\frac{3}{4}\right) \left(\frac{2}{4}\right) \left(\frac{1}{4}\right) = \frac{3}{512} \end{aligned}$$

$$\begin{aligned} \rightarrow G=1: & P(G=1) P(O=0|G=1) P(B=1|G=1) P(C=0|G=1) P(A=1|G=1) \\ &= \left(\frac{6}{8}\right) \left(\frac{4}{8}\right) \left(\frac{3}{8}\right) \left(\frac{4}{8}\right) \left(\frac{6}{8}\right) = \frac{27}{512} \end{aligned}$$

$$\text{Since } \frac{27}{512} > \frac{3}{512}, \quad G = 1$$

• For sample #10...

$$\begin{aligned} \rightarrow G=0: & P(G=0) P(O=1|G=0) P(B=1|G=0) P(C=1|G=0) P(A=1|G=0) \\ &= \left(\frac{2}{8}\right) \left(\frac{3}{4}\right) \left(\frac{3}{4}\right) \left(\frac{3}{4}\right) \left(\frac{1}{4}\right) = \frac{9}{512} \end{aligned}$$

$$\begin{aligned} \rightarrow G=1: & P(G=1) P(O=1|G=1) P(B=1|G=1) P(C=1|G=1) P(A=1|G=1) \\ &= \left(\frac{6}{8}\right) \left(\frac{4}{8}\right) \left(\frac{3}{8}\right) \left(\frac{4}{8}\right) \left(\frac{6}{8}\right) = \frac{27}{512} \end{aligned}$$

$$\text{Since } \frac{27}{512} > \frac{9}{512}, \quad G = 1$$

Sample #9: $G = 1$
Sample #10: $G = 1$

6] Extend Naive Bayes to binary feature values.

a) Use Naive Bayes Assumption to write joint probability of data, $P(x^1, \dots, x^m, y^1, \dots, y^m)$ in terms of θ s.

↳ Naive Bayes Assumption: class conditional independence

The likelihood of one example, $P(x_i, y_i)$ is given by

$$P(y_i) P(x_i | y_i) = P(y_i) P(x_{i1}, \dots, x_{in} | y_i) = P(y_i) \prod_{j=1}^n P(x_{ij} | y_i)$$

Then, extending to m examples,

$$\begin{aligned} P(x^1, \dots, x^m, y^1, \dots, y^m) &= \prod_{i=1}^m P(x_i, y_i) \text{ by independence of data points.} \\ &= \prod_{i=1}^m P(y_i) \prod_{j=1}^n P(x_{ij} | y_i) \end{aligned}$$

Using the Θ notation and indicator on $P(x^1, y^1)$, we find

$$\begin{aligned} P(x^1, y^1) &= P(y^1) \prod_{j=1}^n P(x_{ij} | y^1) \\ &= \Theta_0^{I[Y=0]} (1 - \Theta_0)^{I[Y=1]} \prod_{j=1}^n \prod_{k=1}^s \Theta_{j,k}^{I[Y=0, x_{ij}=k]} \Theta_{j,k}^{I[Y=1, x_{ij}=k]} \end{aligned}$$

$$\therefore P(x^1, \dots, x^m, y^1, \dots, y^m) = \prod_{i=1}^m \Theta_0^{I[Y_i=0]} (1 - \Theta_0)^{I[Y_i=1]} \prod_{j=1}^n \prod_{k=1}^s \Theta_{j,k}^{I[Y_i=0, x_{ij}=k]} \Theta_{j,k}^{I[Y_i=1, x_{ij}=k]}$$

Note, $\log(P(x^1, \dots, x^m, y^1, \dots, y^m))$

$$\begin{aligned} &= \sum_{i=1}^m \left[I[Y_i=0] \log(\Theta_0) + I[Y_i=1] \log(1 - \Theta_0) \right. \\ &\quad \left. + \sum_{j=1}^n \sum_{k=1}^s (I[Y_i=0, x_{ij}=k] \log(\Theta_{j,k}) + I[Y_i=1, x_{ij}=k] \log(\Theta_{j,k})) \right] \end{aligned}$$

Let the boxed expression be denoted as \star .

b) Maximize the joint from (a) wrt $\theta_0, \theta_{j|k|y=0}, \theta_{j|k|y=1}$. Explain the meaning of the results.

$$\frac{\partial \ell}{\partial \theta_0} = \frac{\partial}{\partial \theta_0} \left(\sum_{i=1}^m I[Y_i=0] \log \theta_0 + I[Y_i=1] \log (1-\theta_0) \right)$$

$$= \frac{1}{\theta_0} \sum_{i=1}^m I[Y_i=0] - \frac{1}{1-\theta_0} \sum_{i=1}^m I[Y_i=1]$$

$$= 0$$

$$\Rightarrow \frac{\sum_{i=1}^m I[Y_i=0]}{\theta_0} = \frac{\sum_{i=1}^m I[Y_i=1]}{1-\theta_0}$$

$$\Rightarrow \sum_{i=1}^m I[Y_i=0] = \theta_0 \sum_{i=1}^m I[Y_i=1] + (1-\theta_0) \sum_{i=1}^m I[Y_i=0]$$

$$\Rightarrow \theta_0 = \frac{\sum_{i=1}^m I[Y_i=0]}{\sum_{i=1}^m I[Y_i=0] + \sum_{i=1}^m I[Y_i=1]}$$

$$\Rightarrow \theta_0 = \frac{\sum_{i=1}^m I[Y_i=0]}{m}$$

$$\frac{\partial \ell}{\partial \theta_{j|k|y=0}} = \frac{\partial}{\partial \theta_{j|k|y=0}} \left(\sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^{S_j} I[Y_i=0, X_{ij}=k] \log \theta_{j|k|y=0} + I[Y_i=0, X_{ij}=S_j] \log (1 - \sum_{k=1}^{S_j} \theta_{j|k|y=0}) \right)$$

$$= 0$$

$$\Rightarrow 0 = \frac{1}{\theta_{j|k|y=0}} \sum_{i=1}^m I[Y_i=0, X_{ij}=k] - \frac{1}{1-\sum_{k=1}^{S_j} \theta_{j|k|y=0}} \sum_{i=1}^m I[Y_i=0, X_{ij}=S_j]$$

$$\Rightarrow \theta_{j|k|y=0} \sum_{i=1}^m I[Y_i=0, X_{ij}=k] = (1 - \sum_{k=1}^{S_j} \theta_{j|k|y=0}) \sum_{i=1}^m I[Y_i=0, X_{ij}=S_j]$$

$$\Rightarrow \theta_{j|k|y=0} = \frac{\sum_{i=1}^m I[Y_i=0, X_{ij}=k]}{\sum_{i=1}^m I[Y_i=0, X_{ij}=S_j]} \quad (1)$$



→ Now, must find $\Theta_{j's|y=0}$.

Take sum of both sides of equation iterating k from 1 to $s-1$.

$$\sum_{k=1}^{s-1} \Theta_{j'k|y=0} = \frac{\sum_{k=1}^s \sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=k]}{\sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=s]} \quad \Theta_{j's|y=0}$$

$$1 - \Theta_{j's|y=0} = \frac{\sum_{k=1}^s \sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=k]}{\sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=s]} \quad \Theta_{j's|y=0}$$

$$1 = \left[\frac{\sum_{i=1}^m \sum_{k=1}^{s-1} \mathbb{1}[y_i=0, x_{ij'}=k]}{\sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=s]} + \frac{\sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=s]}{\sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=s]} \right] \Theta_{j's|y=0}$$

$$= \frac{\sum_{i=1}^m \sum_{k=1}^{s-1} \mathbb{1}[y_i=0, x_{ij'}=k]}{\sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=s]} \quad \Theta_{j's|y=0}$$

$$= \left(\frac{\sum_{i=1}^m \mathbb{1}[y_i=0]}{\sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=s]} \right) \Theta_{j's|y=0}$$

$$\Rightarrow \Theta_{j's|y=0} = \frac{\sum_{i=1}^m \mathbb{1}[y_i=0, x_{ij'}=s]}{\sum_{i=1}^m \mathbb{1}[y_i=0]} \quad (2)$$



Combining (1) and (2),

$$\Theta_{j|K|y=0} = \frac{\sum_{i=1}^m \mathbb{I}[Y_i = 0, X_{ij} = k]}{\sum_{i=1}^m \mathbb{I}[Y_i = 0]} = \frac{\sum_{i=1}^m \mathbb{I}[Y_i = 0, X_{ij} = k]}{\sum_{i=1}^m \mathbb{I}[Y_i = 0]}$$

We can find the estimation of the remaining parameters by symmetry, with a summary and explanation of the results as follows:

$$\cdot \Theta_0 = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[Y_i = 0]$$

↳ Relative frequency of datapoints class $y=0$, i.e. the ratio of points with class $y=0$ divided by total # of data points.

$$\cdot \Theta_{j|k|y=0} = \frac{\sum_{i=1}^m \mathbb{I}[Y_i = 0, X_{ij} = k]}{\sum_{i=1}^m \mathbb{I}[Y_i = 0]}$$

↳ Relative frequency of datapoints with class $y=0$ whose j^{th} feature equals the k^{th} feature value.

$$\cdot \Theta_{j|k|y=1} = \frac{\sum_{i=1}^m \mathbb{I}[Y_i = 1, X_{ij} = k]}{\sum_{i=1}^m \mathbb{I}[Y_i = 1]}$$

↳ Relative frequency of datapoints with class $y=1$ whose j^{th} feature equals the k^{th} feature value.