

Ex 2.2

a) The prosecutor claims that

$$P(\text{Btype} = \text{Crime} | \text{Innocent}) = 0,01$$

and confuses this with

$$P(\text{Innocent} | \text{Btype} = \text{Crime}) = 0,01$$

and because this probability is low than he concludes that the probability of being guilty knowing that possesses the guilty blood type is high.

But we know by Bayes theorem that

$P(\text{Btype} = \text{Crime} | \text{Innocent}) \neq P(\text{Innocent} | \text{Btype} = \text{Crime})$
so he is wrong.

b) The defender claims that the probability of being guilty knowing that the defendant has the guilty blood type is only $1/8000$ and that it is not relevant, though it is relevant! The introduction of the blood type knowledge was able to shrink the guilty space from 800000 to 8000 .

Ex. 2.7

lets consider the example of the toss of a balanced coin twice, therefore the probability of heads (H) and tails (T) is $0,5$.

lets define the following events:

A - having the first toss to be heads, $\{HH, HT\}$

B - having the second toss to be heads, $\{HH, TH\}$

C - having both tosses the same, $\{HH, TT\}$

The probability of each event is:

(2)

$$P(A) = P(B) = P(C) = \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{2}$$

The probability of $A \cap B$, $B \cap C$ and $A \cap C$ are:

$$P(A \cap B) = \frac{1}{2} \cdot \frac{1}{2} = P(A) \cdot P(B) = \frac{1}{4}$$

$$P(B \cap C) = \frac{1}{2} \cdot \frac{1}{2} = P(B) \cdot P(C) = \frac{1}{4}$$

$$P(A \cap C) = \frac{1}{2} \cdot \frac{1}{2} = P(A) \cdot P(C) = \frac{1}{4}$$

This means the events are pair wise independent.

Now let's check for mutual independence:

$$P(A \cap B \cap C) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \neq P(A) \cdot P(B) \cdot P(C) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8}$$

So they are not mutually independent though they are pair wise!

Ex 4.2

$$a) P(Y < a) = E\{P(Y < a)\} = E\{P(WX < a)\} = \sum_{\bar{w} \in \mathcal{W}} P(\bar{w}X < a) \cdot P_{\bar{w}}(\bar{w})$$

$$= P(-X < a) \cdot P(-1) + P(X < a) \cdot P_{\bar{w}}(1) =$$

$$= \frac{1}{2} \cdot P(X < a) + \frac{1}{2} P(X < a) = P(X < a), \forall a \in \mathbb{R}$$

therefore both follow the same distribution.

$$\begin{aligned} b) \quad \text{cov}[X, Y] &= E[XY] - \overset{=0}{E[X]}E[Y] = E[E[XY|W]] = \\ &= E[E[X^2 W|W]] = E\left[\sum_{\bar{w} \in \mathcal{W}} x^2 \bar{w} \cdot P_{\bar{w}}(\bar{w})\right] = \\ &= E[-x^2 \cdot 0,5 + x^2 \cdot 0,5] = E[0] = 0 \end{aligned}$$

$$\begin{aligned}
 \hat{\sigma}(x) &= \arg \min_{a \in A} P(a|x) = \arg \min_{a \in A} \sum_{y \in \{1, \dots, c+1\}} L(y, a) \cdot P(y|x) = \\
 &= \arg \min_{a \in A} L(1, a) \cdot P(1|x) + \dots + L(j, a) \cdot P(j|x) + \dots + \\
 &\quad + L(c, a) \cdot P(c|x)
 \end{aligned}$$

• We choose $a = \alpha_j$ if $P(\alpha_j|x) \leq P(\alpha_i|x)$
 $\forall i \in \{1, \dots, c\} \setminus \{j\}$

$$\begin{aligned}
 P(\alpha_j|x) \leq P(\alpha_i|x) &\Leftrightarrow \\
 \Leftrightarrow L(1, \alpha_j) \cdot P(1|x) + \dots + L(j, \alpha_j) \cdot P(j|x) + \dots + L(c, \alpha_j) \cdot P(c|x) &\leq L(1, \alpha_i) \cdot P(1|x) + \dots + \\
 + \dots + L(c, \alpha_i) \cdot P(c|x) &\Leftrightarrow L(1, \alpha_j) \cdot P(1|x) + \dots + L(j, \alpha_j) \cdot P(j|x) + \dots + L(c, \alpha_j) \cdot P(c|x) \leq \\
 L(1, \alpha_i) \cdot P(1|x) + \dots + L(j, \alpha_i) \cdot P(j|x) + \dots + L(c, \alpha_i) \cdot P(c|x) &\Leftrightarrow \\
 \Leftrightarrow Z_s \cdot P(c|x) \leq Z_s \cdot P(j|x) &\Leftrightarrow P(j|x) \geq P(c|x)
 \end{aligned}$$

• And also if $P(\alpha_j|x) \leq P(\alpha_{c+1}|x)$

$$\begin{aligned}
 P(\alpha_j|x) \leq P(\alpha_{c+1}|x) &\Leftrightarrow \\
 \Leftrightarrow L(1, \alpha_j) \cdot P(1|x) + \dots + L(j, \alpha_j) \cdot P(j|x) + \dots + L(c, \alpha_j) \cdot P(c|x) &\leq Z_r \\
 + \dots + L(c+1, \alpha_j) \cdot P(c+1|x) &\leq Z_r \Leftrightarrow \\
 \Leftrightarrow Z_s \cdot \sum_{\substack{i=1 \\ i \neq j}}^c P(i|x) \leq Z_r &\Leftrightarrow Z_s (1 - P(j|x)) \leq Z_r \Leftrightarrow \\
 \Leftrightarrow P(j|x) \geq 1 - \frac{Z_r}{Z_s}
 \end{aligned}$$

b) As Z_r/Z_s increases the risk of rejection grows proportionally to the risk of misclassification.

when $Z_r/Z_s = 0$, i.e., $Z_r = 0$ there is no risk in rejecting so the hypothesis is always rejected ($P(x=j|x)$ can't be bigger than 1)

when $Z_r/Z_s = 1$, then $1 - \frac{Z_r}{Z_s} = 0$ and no hypothesis is ever rejected

a) Multinomial distribution

$$f(x, y; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \cdot \exp\left(-\frac{1}{2} (v - \mu)^T \Sigma^{-1} (v - \mu)\right)$$

where:

$$v = \begin{bmatrix} y \\ x \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_y \\ \mu_x \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{bmatrix}$$

Maximum Likelihood

of μ and Σ

• maximize distribution of N points

$$\begin{aligned} f(x_1, y_1, \dots, x_N, y_N; \mu, \Sigma) &= \prod_{i=1}^N f(x_i, y_i; \mu, \Sigma) = \\ &= \left(\frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \right)^N \cdot \exp\left(-\frac{1}{2} \sum_{i=1}^N (v_i - \mu)^T \Sigma^{-1} (v_i - \mu)\right) \end{aligned}$$

• $L = \log(\cdot)$

$$L = -\frac{N}{2} \left[\log(2\pi)^n + \log |\Sigma| \right] - \frac{1}{2} \sum_{i=1}^N (v_i - \mu)^T \Sigma^{-1} (v_i - \mu)$$

• $\hat{\mu}$ - first order condition

$$\frac{\partial L}{\partial \hat{\mu}} = 0 \Leftrightarrow \sum_{i=1}^N (v_i - \hat{\mu})^T \Sigma^{-1} = 0 \Leftrightarrow \sum_{i=1}^N (v_i) - N \cdot \hat{\mu} = 0 \Leftrightarrow$$

$$\Leftrightarrow \hat{\mu} = \frac{1}{N} \sum_{i=1}^N v_i$$

$$\Rightarrow \hat{\mu} = \begin{bmatrix} \frac{1}{N} \sum_{i=1}^N y_i \\ \frac{1}{N} \sum_{i=1}^N x_i \end{bmatrix}$$

• $\hat{\Sigma}$ - first order conditions

$$\begin{aligned} \text{using: } \frac{\partial}{\partial A} x^T A x &= \frac{\partial}{\partial A} \text{tr}(x^T A x) \\ &= \frac{\partial}{\partial A} \text{tr}(x x^T A) \\ &= (x x^T)^T = x x^T \end{aligned}$$

$$\frac{\partial}{\partial A} \log |A| = A^{-1}$$

$$\frac{\partial L}{\partial \hat{\Sigma}} = 0 \Leftrightarrow \frac{N}{2} \hat{\Sigma} - \frac{1}{2} \sum_{i=1}^N (v_i - \mu)(v_i - \mu)^T = 0 \Leftrightarrow$$

$$\Leftrightarrow \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (v_i - \mu)(v_i - \mu)^T$$

$$\Leftrightarrow \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} y_i - \mu_y \\ x_i - \mu_x \end{bmatrix} \begin{bmatrix} y_i - \mu_y \\ x_i - \mu_x \end{bmatrix}^T$$

$$\Leftrightarrow \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} (y_i - \mu_y)(y_i - \mu_y)^T & (y_i - \mu_y)(x_i - \mu_x)^T \\ (x_i - \mu_x)(y_i - \mu_y)^T & (x_i - \mu_x)(x_i - \mu_x)^T \end{bmatrix}$$

$$\hat{\Sigma}_{xx} = \frac{1}{N} X_c^T \cdot X_c \quad \hat{\Sigma}_{yy} = \frac{1}{N} Y_c^T Y_c$$

$$\hat{\Sigma}_{xy} = \frac{1}{N} X_c^T Y_c \quad \hat{\Sigma}_{yx} = \frac{1}{N} Y_c^T X_c$$

$$\hat{\mu}_x = \frac{1}{N} X^T \cdot \mathbf{1}$$

$$\hat{\mu}_y = \frac{1}{N} Y^T \cdot \mathbf{1}$$

$$P(Y|X) = \mathcal{N}(Y | \mu_{Y|X}, \Sigma_{Y|X}) \quad \mu_{Y|X} = E\{Y|X\}$$

$$E\{Y|X\} = \hat{\mu}_y + \Sigma_{yx} \Sigma_{xx}^{-1} (x - \hat{\mu}_x)$$

$$= \Sigma_{yx} \Sigma_{xx}^{-1} x + \hat{\mu}_y - \Sigma_{yx} \Sigma_{xx}^{-1} \hat{\mu}_x$$

$$= \frac{1}{N} \cdot Y_c^T X_c \cdot \mathcal{N}(X_c^T X_c)^{-1} x + \hat{\mu}_y - \frac{1}{N} \cdot Y_c^T X_c \cdot \mathcal{N}(X_c^T X_c)^{-1} \hat{\mu}_x + \hat{\mu}_y$$

$$= \underbrace{Y_c^T X_c (X_c^T X_c)^{-1}}_{w^T} x + \hat{\mu}_y - \underbrace{Y_c^T X_c (X_c^T X_c)^{-1} \hat{\mu}_x}_{w_0^T}$$

↳

②

My answer will evaluate the advantages and disadvantages of generative and discriminative approaches generally and not only to linear regression.

- Easy to fit?

It is usually easier to fit generative classifiers.

Ex: Naive Bayes model can be fitted by only counting while logistic regression requires solving a complex optimization problem.

- Fit classes separately?

In generative classifiers, the parameters of each class are estimated independently, so there is no need to retrain the model when adding more classes.

In discriminative models all the parameters interact, so the whole model must be retrained when inserting a new class.

- Handle missing features easily?

In the discriminative models there is no solution to deal with missing features while in generative ones there is.

- Can handle unlabeled training data?

Semi-supervised learning is much easier to do with generative models than discriminative.

- Symmetric in inputs and outputs?

A generative model can be run "backwards", and infer probable inputs given the output by computing $P(X|y)$. Not possible with the discriminative model.

- Can handle feature preprocessing?

A big advantage of discriminative methods is that they allow the preprocess of the data in arbitrary ways.

It is hard to define a generative model on preprocessed data, since the new features are correlated in complex ways.

- Well calibrated probabilities?

Some generative models, such as the naive Bayes, make strong assumption which are often not valid. This can lead to extreme posterior class probabilities.

Discriminative models are usually better calibrated in terms of their probability estimates