

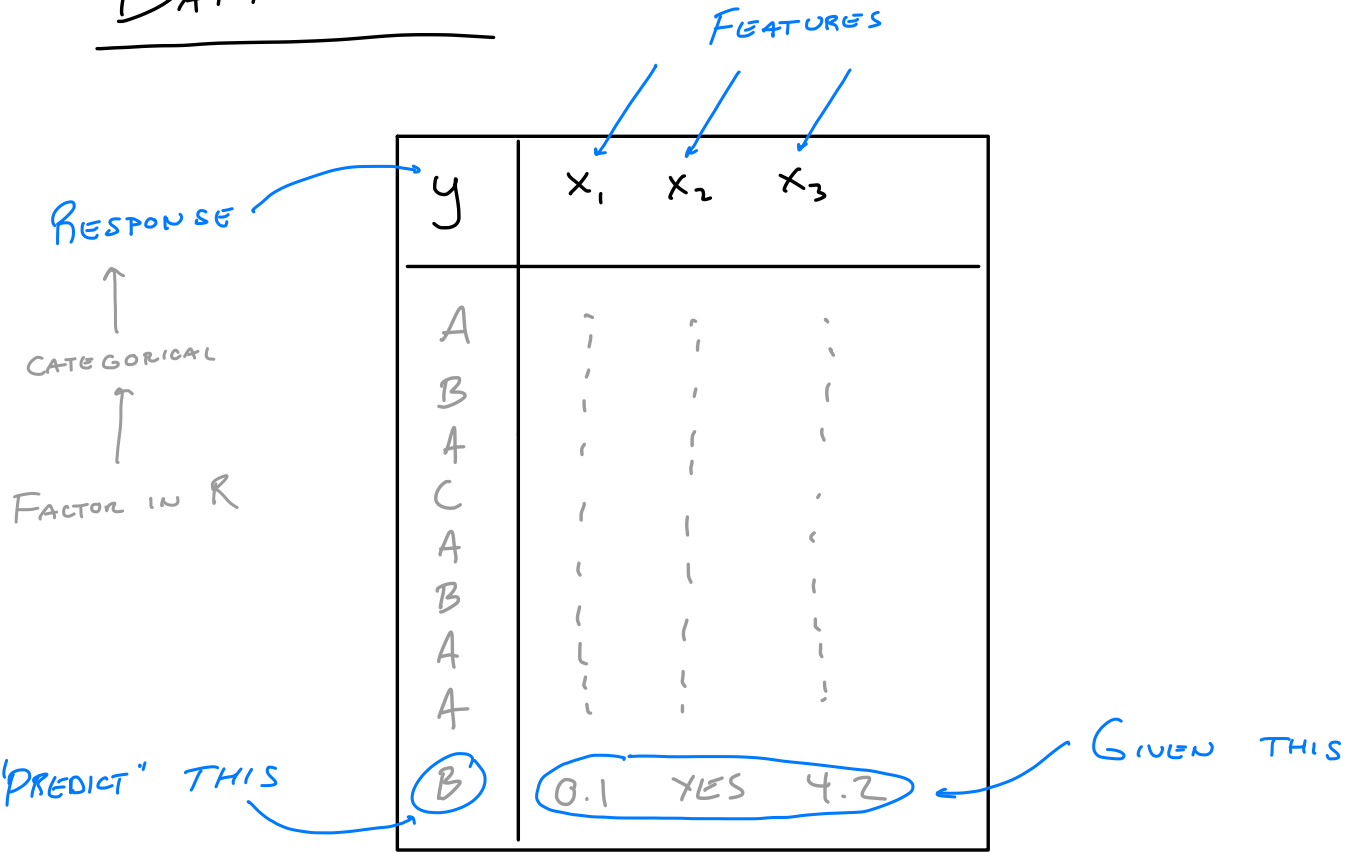
CLASSIFICATION

↳ AN INTRODUCTION

DALPIAZ

02/17/2020

DATA VIEW



PROBABILITY VIEW

$$(X, Y) \in \mathbb{R}^P \times \{1, 2, \dots, K\}$$

↑ FEATURES ↑ RESPONSE

FIND A CLASSIFIER $C(x)$ THAT MINIMIZES
OUTPUT: ONE OF THE K CATEGORIES

$$P[C(x) \neq Y] \rightarrow \text{MISCLASSIFICATION RATE}$$

WHERE $C: \mathbb{R}^P \rightarrow \{1, 2, 3, \dots, K\}$

BAYES CLASSIFIER

← MINIMIZES MISCLASSIFICATION

$$C^B(x) \triangleq \underset{k \in \{1, \dots, K\}}{\operatorname{ARGMAX}} P[y=k | X=x]$$

GIVEN FEATURE VECTOR x , CLASSIFY OBSERVATION

AS THE CATEGORY WITH THE HIGHEST PROBABILITY.

DUH?

EXAMPLE

$$C^B(x=0) = ?$$

	X		
	0	1	
A	0.1	0.1	0.2
B	0.2	0.1	0.3
C	0.1	0.4	0.5

JOINT
DISTRIBUTION
OF (X, Y)

$$P[Y | X=0] = \begin{cases} 0.25 & y = A \\ 0.50 & y = B \\ 0.25 & y = C \end{cases} = \frac{P[X=0 \cap Y=A]}{P[X=0]}$$

CONDITIONAL DISTRIBUTION
OF Y | X=0

0.4 0.6

MARGINAL DISTRIBUTION

OF X

$$C^B(x=0) = B$$

$$C^B(x=1) = C$$

BAYES ERROR

← AVERAGE MISCLASSIFICATION
USING BAYES CLASSIFIER

$$1 - E_x \left[\max_k P[Y=k | X=x] \right]$$

"(IRREDUCIBLE ERROR)"

	X		
	0	1	
A	0.1	0.1	0.2
B	0.2	0.1	0.3
C	0.1	0.4	0.5
	0.4	0.6	

$$= 1 - \left[\max_k P[Y=k | X=0] P[X=0] + \max_k P[Y=k | X=1] P[X=1] \right]$$

$$= 1 - \left[\left(\frac{0.2}{0.4} \right) (0.4) + \left(\frac{0.4}{0.6} \right) (0.6) \right]$$

$$= 1 - [0.2 + 0.4] = \underline{0.4}$$

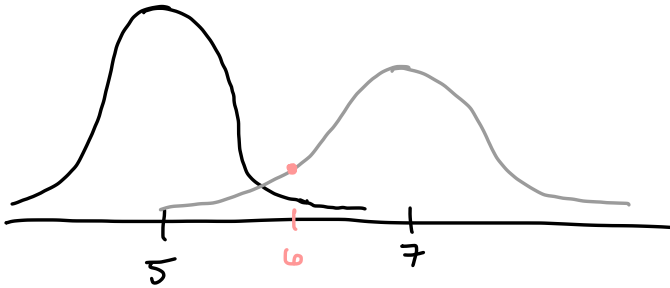
EXAMPLE

$$X | Y=0 \sim \mathcal{N}(\mu=5, \sigma=1) \quad | \quad f_0(x)$$

$$X | Y=1 \sim \mathcal{N}(\mu=7, \sigma=2) \quad | \quad f_1(x)$$

$$\pi_0 = P[Y=0] = 0.6$$

$$\pi_1 = P[Y=1] = 0.4$$



$$C^B(x=6) = ?$$

$$\text{CALCULATE } P[Y=0 | X=6] = \frac{\pi_0 f_0(6)}{\pi_0 f_0(6) + \pi_1 f_1(6)}$$

BINARY CLASSIFICATION

$$Y = 0$$

↑
"NEGATIVE"

OR

$$Y = 1$$

↑
"POSITIVE"

LATER

FP/TP

FN/TN

ETC

$$p(x) \triangleq P[Y=1 | X=x]$$

$$1-p(x) = P[Y=0 | X=x]$$

$$C^B(x) = \begin{cases} 1 & p(x) \geq 0.5 \\ 0 & \text{ELSE} \end{cases}$$

IN PRACTICE

Don't know $P[Y=k | X=x]$!!!

ESTIMATE IT!

CLASSIFIER



$$C(x) = \underset{k}{\text{ARGMAX}} \hat{P}[Y=k | X=x]$$

ESTIMATE OF CONDITIONAL DISTRIBUTION

A "GUESS" FOR



$$C^B(x)$$

How?

ESTIMATING CONDITIONAL DISTRIBUTIONS

KNN w/ `caret::knn3()`

TREES w/ `rpart::rpart()`

LINEAR MODELS w/ `glm()` ?
`nnet::nnet()` ?

} `predict(mod, data, type)`
≡

MAKE SURE RESPONSE VARIABLE IS A FACTOR!