

1. Above is some data similar to today's slides, except that shape has used to describe the classes rather than the blue/orange. Start by covering up the validation data (right) and then answer the questions below.

- Draw freehand an estimate of a good separating line for the points in the training data.
- Now, uncover the validation data and create the same line from (a) on the right.
- Next, compute a confusion matrix for the validation data in the space below:

		Predicted	
		●	▽
Actual	●	9	1
	▽	2	8

[your line and results may be different]

- Calculate the error rate for the training data and the validation data.

$$\text{training} = \frac{3}{20} = .15$$

$$\text{validation} = \frac{3}{20} = .15$$

[features are just possible options; signs in parenthesis]

2. Five examples of predictive modelling tasks are described. For each, identify (i) the classes, (ii) what an observation is, (iii) two possible numeric features, and (iv) pick one class to be the reference class (the class you are trying to predict) and estimate the signs you would expect to be associated with each of the features.

a. Create a model to predict whether a basketball player makes a free-throw shot.

Classes: made, missed

Observation: shot

Features: distance to basket (+)

player height (+)

b. Predict whether it will rain tomorrow.

Classes: rain, dry

Observation: day

Features: mm of rain this year (+)

mm of rain last year (+)

[could be neg]

c. Estimate whether a user's cell phone will still be working one year from now.

Classes: works, broken

Observation: cell phone  
or user

Features: current cell age (-)

current cell price (+)

d. Predict whether a tweet is associated with positive or negative feelings.

Classes: positive, negative

Observation: one tweet

Feature: # smiling emoji (+)

# frowning emoji (-)

e. Predict whether a music song is classified as "rock" or "classical" music.

Classes: classical, rock

Observation: song

Feature: length of song (+)

loudness of song (-)

f. Create a model to determine if an image was taken indoors or outside.

Classes: indoor, outdoor

Observation: image

Feature: brightness (-)

# people (+)

3. Assume that we have built a logistic regression model to predict whether a text message is spam using two features: the count of the number of numeric digits 0-9 (we will call this  $X$ ) and the count of the number of words with two or more letters that are written all in capitals (we will call this  $Y$ ). Our training data yielded the following model:

$$\text{Prob}(\text{spam}) = F(-1.2 + 0.3 \times X + 0.2 \times Y)$$

Answer the following questions based on this model.

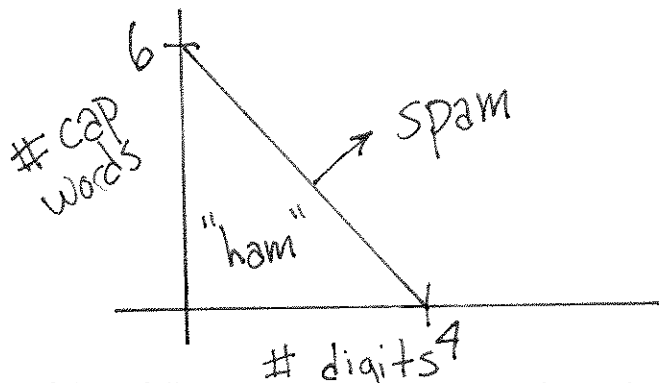
- a. If a message contains no words with two or more letters that are written all in capitals, after how many numerical digits is there a predicted probability of at least 0.5 that it is spam?

Need  $0.3X = 1.2 \Rightarrow 4 \text{ digits needed}$

- b. If a message contains no numerical digits, after how many words with two or more letters that are written all in capitals is there a predicted probability of at least 0.5 that it is spam?

Need  $0.2Y = 1.2 \Rightarrow 6 \text{ words needed}$

- c. Based on (a) and (b), draw a sketch of the boundary line implied by the logistic regression model.



- d. For each of three following messages, compute the probability that the text message is spam.

- "Congratulations, you have won €12.000. Text us back ASAP to claim!"
- "Hi Joe. I owe you £30.50 for dinner last night. Do you have Venmo?"
- "I am taking MATH389 this semester. It is by far the WORST class I have ever taken."

i.  $F(-1.2 + 0.3 \times 5 + 0.2 \times 1) = F(0.5) = 0.623$

ii.  $F(-1.2 + 0.3 \times 4 + 0.2 \times 0) = F(0) = 1/2$

iii.  $F(-1.2 + 0.3 \times 3 + 0.2 \times 1) = F(-0.1) = 0.475$

could be 2 if  
you count MATH389

4. I promised we would not do much in the way of formal mathematical calculations, but this one is relatively easy and conceptually important. We defined the logistic function  $F(x) = \exp(x) / (\exp(x) + 1)$ . Use this to answer the two questions below.

- a. Show that  $1 - F(x)$  is equal to  $F(-x)$  for any number  $x$ .

$$\begin{aligned} 1 - F(x) &= 1 - \frac{\exp(x)}{\exp(x) + 1} = \frac{\exp(x) + 1}{\exp(x) + 1} - \frac{\exp(x)}{\exp(x) + 1} \\ &= \frac{1}{\exp(x) + 1} = \frac{1}{\exp(x) + 1} \cdot \frac{\exp(-x)}{\exp(-x)} = \frac{\exp(-x)}{1 + \exp(-x)} = F(-x) \end{aligned}$$

- b. From today's notes, assume we have built a two-feature logistic regression model  $a + b \times X + c \times Y$  to predict the probability that a point is orange. Use the result in (a) to calculate the logistic regression model to predict that probability that a point is blue. From the result, you should see that it is not particularly important which class we set as the reference class.

If

$$\text{Probability [orange]} = F(a + b \times X + c \times Y)$$

Then

$$\begin{aligned} \text{Probability [blue]} &= 1 - \text{Probability [orange]} \\ &= 1 - F(a + b \times X + c \times Y) \end{aligned}$$

$$\text{from (a)} \rightarrow = F(-a - b \times X - c \times Y)$$

So same model for probability orange but with signs of parameters flipped.