

WEEK 2

NAIVE BAYES

BAYESIAN CLASSIFICATION

→ is part of a family of **probabilistic** classification

- First, compute the bayesian probability of each class

$$P(y|x) = \frac{P(x|y) P(y)}{\sum_{y'} P(x|y') P(y')}$$

P(x|y)
P(y)
P(x)
Prior
Normalizer factor

(class model)

- Then the most probable class given observation

$$\hat{y} = \arg \max_y P(y|x)$$

COMPONENTS

- i.e. y ... UK patient has Ebola
 x ... observed symptoms

If we present a Bayesian classifier w/
an instance that is very unlike any instance
it trained on, the value of the normalizer
 $\sum_y P(x|y) P(y)$ will be small

$P(y)$: prior probability of each class

$P(x|y)$: class-conditional model

↳ describes how likely to see observation x
for class y

$P(x)$: normalize probabilities across observations

↳ does not affect which class is most likely
 $(\arg \max)$

↳ makes probabilities more comparable to
other data points/instances in your data

NAIVE BAYES: A GENERATIVE MODEL

→ Computes its predictions by modelling each class (not just by looking at boundary)



- $P(y|x) \propto P(x|y) P(y)$

→ All generative classifiers are probabilistic, BUT not all probabilistic classifiers are generative

- It is possible to ~~estimate~~ estimate $P(y|x)$ directly
 - ↳ i.e. logistic regression

→ For a probabilistic classifier, the decision boundary represents the points where classes are equally probable

INDEPENDENCE ASSUMPTION

→ We assume observations $x_1 \dots x_d$ are conditionally independent given y

$$\begin{aligned} P(x|y) &= P(x_1, x_2, \dots, x_d | y) \\ &= [P(x_1 | x_2, \dots, x_d, y) P(x_2 | x_3, \dots, x_d, y) \dots P(x_d | y)] \xrightarrow{\text{CHAIN RULE}} \end{aligned}$$

↓ SIMPLIFIED TO

$$\begin{aligned} P(x|y) &= P(x_1|y) P(x_2|y) \dots P(x_d|y) \xrightarrow{\text{INDEPENDENCE}} \\ &= \prod_{i=1}^d P(x_i|y) \end{aligned}$$

CONDITIONAL INDEPENDENCE

- 2 random events A & B are conditionally indep. given a third event C if the occurrence of A and occurrence of B are indep. events in their conditional probability distribution given C.

i.e. Probabilities of going to the beach and getting a heat stroke are NOT independent: $P(B, HS) > P(B)P(HS)$

↳ May be independent if we know the weather is hot

$$P(B, HS | \text{hot}) = P(B|\text{hot}) P(HS|\text{hot})$$

↳ Hot weather 'explains' all the dependence between beach & heatstroke

- In classification, class value explains all the dependence between attributes

GAUSSIAN NAIVE BAYES CLASSIFIER [For real-valued data]

→ Suppose we're trying to distinguish children from adults based on size

- Classes: $\{a, c\}$
- Attributes: height (cm), weight (kg)
- Training examples: $\{h_i, w_i, y_i\}$ for 4 adults and 12 children
the class: are they an adult/child?

(1)

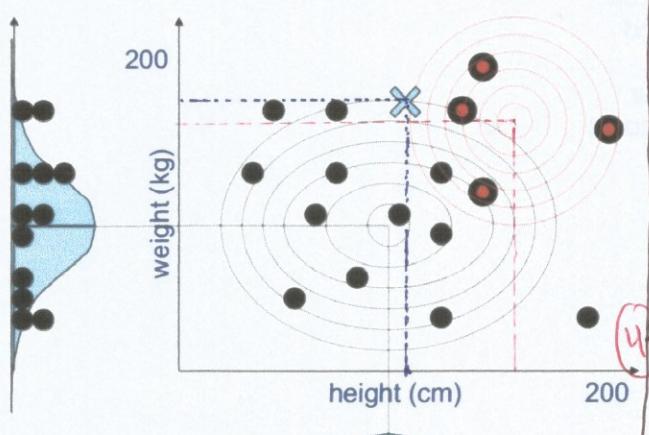
Estimate prior: $P(a) = \frac{4}{4+12} = 0.25, P(c) = 0.75$

(2)

Assume heights and weights are modelled by gaussian distribution:

Model for adults <ul style="list-style-type: none"> - Height - Gaussian w/ mean, variance $\mu_{h,a} = \frac{1}{4} \sum_{i:y_i=a} h_i$ - Weight - Gaussian ($\mu_{w,a}, \sigma^2_{w,a}$) - Assume height & weight are independent. 	$\sigma^2_{h,a} = \frac{1}{4} \sum_{i:y_i=a} (h_i - \mu_{h,a})^2$
Model for children <ul style="list-style-type: none"> - Height ($\mu_{h,c}, \sigma^2_{h,c}$) - Weight ($\mu_{w,c}, \sigma^2_{w,c}$) 	<p>R Maximum likelihood Estimate (MLE) of the mean and variance</p>

$$P(a) = \frac{4}{4+12} = 0.25; P(c) = 0.75$$



Copyright © Victor Lavrenko, 2014

$(\mu_{h,c}, \sigma^2_{h,c})$
Not a very good fit; but
then again we don't have
very many points to begin w/

This individual has a height that looks a lot like a child's height; it falls very close to the gaussian distribution for children's height

(3).

$$p(h_x|c) = \frac{1}{\sqrt{2\pi \sigma_{h,c}^2}} \exp -\frac{1}{2} \left(\frac{(h_x - \mu_{h,c})^2}{\sigma_{h,c}^2} \right) = \text{going to be a big number}$$

$$p(w_x|c) = \frac{1}{\sqrt{2\pi \sigma_{w,c}^2}} \exp -\frac{1}{2} \left(\frac{(w_x - \mu_{w,c})^2}{\sigma_{w,c}^2} \right) = \text{going to be small number}$$

$$p(h_x|a) = \frac{1}{\sqrt{2\pi \sigma_{h,a}^2}} \exp -\frac{1}{2} \left(\frac{(h_x - \mu_{h,a})^2}{\sigma_{h,a}^2} \right) = \text{going to be a small number}$$

$$p(w_x|a) = \frac{1}{\sqrt{2\pi \sigma_{w,a}^2}} \exp -\frac{1}{2} \left(\frac{(w_x - \mu_{w,a})^2}{\sigma_{w,a}^2} \right) = \text{going to be a big number}$$

$$P(x|a) = p(h_x|a)p(w_x|a)$$

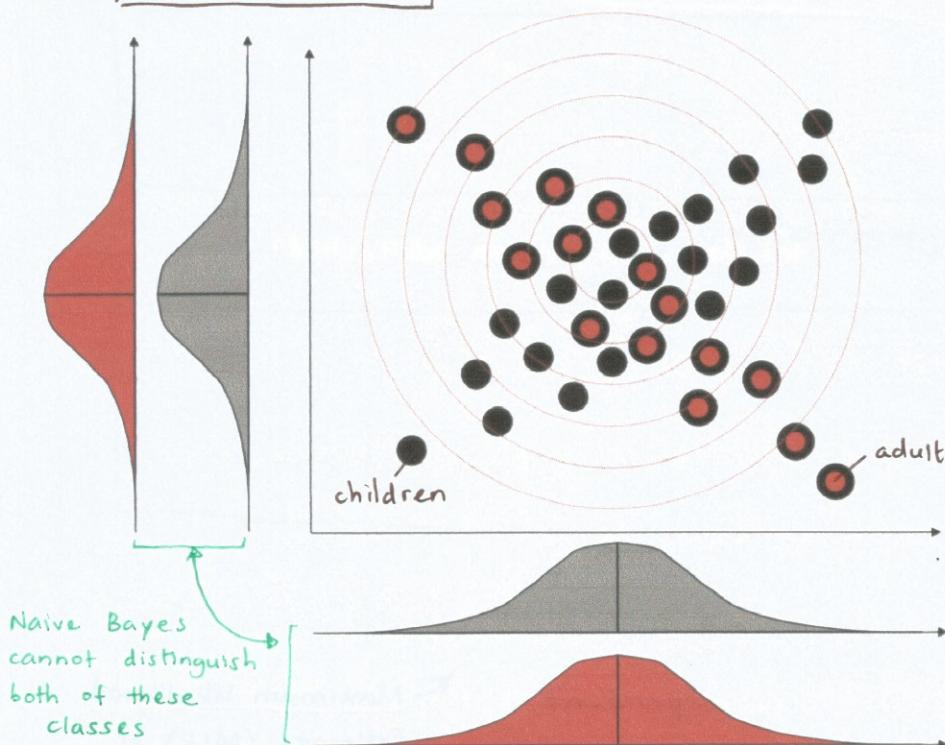
$$P(x|c) = p(h_x|c)p(w_x|c)$$

$$P(a|x) = \frac{P(x|a)P(a)}{P(x|a)P(a)+P(x|c)P(c)}$$

(4).

PROBLEMS W/ NAIVE BAYES

1) Correlation



- What differentiates the 2 classes is not the mean height/mean weight/variance, but it is the co-variance.

~~black & red classes are effectively correlated~~

- There is (+)ve correlation for black, (-)ve correlation for red class; and that is the only difference

↓
Naive Bayes cannot do correlation!

2) Misclassification + Assuming word independence

- Separate spam from valid email, attributes = words

D1: "send us your password"
 D2: "send us your review"
 D3: "review your password"
 D4: "review us"
 D5: "send your password"
 D6: "send us your account"

new email: "review us now"

spam
ham
ham
spam
spam
spam

		P (spam) = 4/6 P (ham) = 2/6	
spam	ham	spam	ham
2/4	1/2	password	
1/4	2/2	review	
3/4	1/2	send	
3/4	1/2	us	
3/4	1/2	your	
1/4	0/2	account	

$$\begin{aligned} \rightarrow P(\text{spam} | \text{"review us"}) &\propto P(\text{spam}) P(0, 1, 0, 1, 0, 0 | \text{spam}) \\ &\propto \frac{4}{6} \cdot (1 - \frac{2}{4})(\frac{1}{4})(1 - \frac{3}{4})(\frac{3}{4})(1 - \frac{3}{4})(1 - \frac{1}{4}) \\ &\approx 0.00293 \end{aligned}$$

$$\begin{aligned} \rightarrow P(\text{ham} | \text{"review us"}) &\propto \frac{2}{6} \cdot (1 - \frac{1}{2})(\frac{2}{2})(1 - \frac{1}{2})(\frac{1}{2})(1 - \frac{1}{2})(1 - \frac{0}{2}) \\ &\approx 0.02083 > 0.00293 \end{aligned}$$

→ Classify the email as ham. However, notice the email is same as D4 which is a ~~spam message~~ training example.

↳ The only thing we did wrong is assume independence.

→ Naive Bayes allow this to happen.

3) Zero-frequency problem

$$P(\text{'account'} | \text{ham}) = \frac{0}{2}$$

→ The word 'account' has zero probability in ham emails

→ SOLUTION: Laplace Smoothing

- Add a small (1)ve no. to all counts

$$\epsilon = \frac{\text{num}(w)}{\text{num}}$$

MISSING VALUES IN NAIVE BAYES

→ Say we want to compute $P(X_1 = x_1, \dots, X_J = ?, \dots, X_d = x_d | y)$
but we don't know the value for some attribute X_J

- W/ Naive Bayes, you can simply ignore this attribute, hence

$$P(x_1 \dots \boxed{X_J} \dots x_d | y) = \prod_{i \neq J}^d P(x_i | y)$$

- This is based on conditional independence between attributes.