

# Bayes Classifier, MAP, Naïve Bayes

# Contents

- MAP classification
  - Naïve Bayes

# Reminder

- Next Tuesday – Review for Exam 1
- Next Thursday – Exam 1

# Big Picture

- machine learning is about **learning functions from data**
- the inputs to the functions are called **features**

$$\mathbf{x} = \boxed{9}$$

- the true output is called a label
- $y \in \{0, 1, 2, \dots, 9\}$
- labeled examples of features make up *training* and *test* data

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

- useful characterization:
- $(\mathbf{x}_i, y_i) \stackrel{i.i.d.}{\sim} p(\mathbf{x}, y)$
- if we have a good approximation of  $p(\mathbf{x}, y)$ , we can design a good function for classification:

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

# Maximum A-posteriori Probability (MAP)

- simple idea: given the data or features  $\mathbf{x}$ , what is the most probable class  $y$ ?
- the *posterior* distribution  $p(y|\mathbf{x})$  answers this

$$\hat{y} = f(\mathbf{x}) = \arg \max_y p(y|\mathbf{x}) \quad (f(x) \text{ is a function, not a pdf})$$

- in words: MAP picks the most probable class given the observation  $\mathbf{x}$

- example:  $x = \boxed{\text{4}}$   $f(\mathbf{x}) = \arg \max_{i=0,\dots,9} p(y = i|\mathbf{x} = \boxed{\text{4}})$

$$p(y = 0|\mathbf{x} = \boxed{\text{4}}) =$$

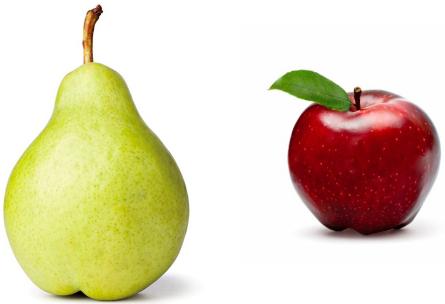
$$p(y = 1|\mathbf{x} = \boxed{\text{4}}) =$$

⋮

$$p(y = 9|\mathbf{x} = \boxed{\text{4}}) =$$

- challenge: estimating  $p(y|\mathbf{x})$  from data.

# Example – pear/apple classifier



$x$  is top to bottom ratio

$y = 0$  if pear,  $y = 1$  if apple

$x \in \{0.5, 1, 2\}$  and  $y \in \{0,1\}$

- $x$  is discrete.

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$



		$x$		
		0.5	1	2
$y$	0	10	30	0
	1	0	50	10

$$p(\mathbf{x}, y) = \frac{\sum_{i=1}^n \mathbb{I}_{\{\mathbf{x}_i = \mathbf{x}, y_i = y\}}}{n}$$

- classification: given  $x$ , what is the best guess for  $y$ ?
- MAP rule:

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

		$x$		
		0.5	1	2
$y$	0	0.1	0.3	0
	1	0	0.5	0.1

# Example – pear/apple classifier



$x$  is top to bottom ratio

$y = 0$  if pear,  $y = 1$  if apple

		$x$		
		0.5	1	2
$y$	0	0.1	0.3	0
	1	0	0.5	0.1

- new fruit with  $x = 0.5$ ?

- new test fruit with  $x = 1$ ?

- MAP rule:

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

- new test fruit with  $x = 2$ ?

# Maximum-a-posteriori (MAP)

- MAP estimate (Bayes optimal classification rule):

- $x, y$  discrete

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

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})}$$

$$\hat{y} = \arg \max_y \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})}$$

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

- $y$  discrete,  $x$  continuous

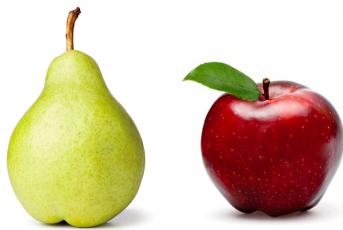
$$\hat{y} = \arg \max_y p(y|\mathbf{x})$$

$$p(y|\mathbf{x}) = \frac{f(\mathbf{x}|y)p(y)}{f(\mathbf{x})}$$

$$\hat{y} = \arg \max_y \frac{f(\mathbf{x}|y)p(y)}{f(\mathbf{x})}$$

$$\hat{y} = \arg \max_y f(\mathbf{x}|y)p(y)$$

# Example – pear/apple classifier



$x$  is top to bottom ratio

$y = 0$  if pear,  $y = 1$  if apple

$x$  is continuous

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$



$$p(y = 0) = p(y = 1) = 0.5$$

- classification: given  $x$ , what is the best guess for  $y$ ?
- MAP rule:

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

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

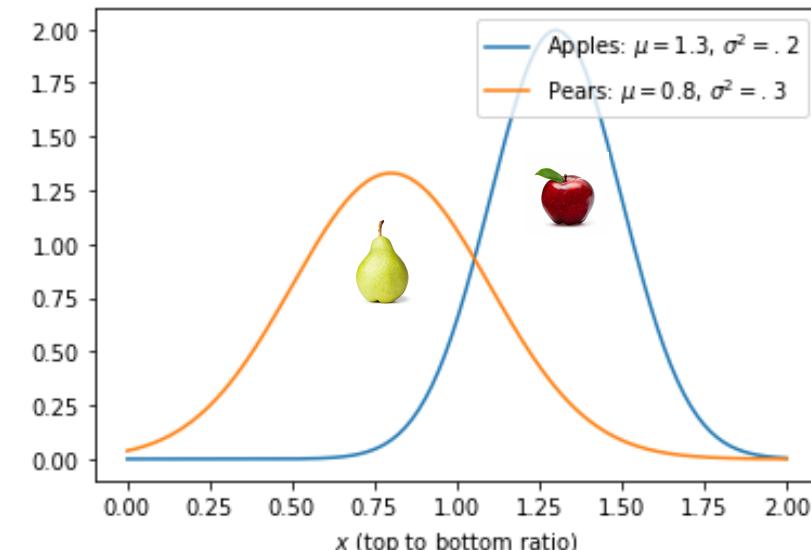
$$X|Y = 1 \sim \mathcal{N}(1.3, 0.2)$$

$$f(x|y = 1) = \frac{1}{\sqrt{2\pi(0.2)^2}} e^{-\frac{(x-1.3)^2}{2(0.2)^2}}$$



$$X|Y = 0 \sim \mathcal{N}(0.8, 0.3)$$

$$f(x|y = 0) = \frac{1}{\sqrt{2\pi(0.3)^2}} e^{-\frac{(x-0.8)^2}{2(0.3)^2}}$$



# Example – pear/apple classifier



$x$  is top to bottom ratio

$y = 0$  if pear,  $y = 1$  if apple

$p(y = 0) = p(y = 1) = 0.5$

- if

$$\frac{1}{\sqrt{2\pi(0.2)^2}} e^{-\frac{(x-1.3)^2}{2(0.2)^2}} \geq \frac{1}{\sqrt{2\pi(0.3)^2}} e^{-\frac{(x-0.8)^2}{2(0.3)^2}}$$

then  $\hat{y} = 1$  else  $\hat{y} = 0$ .

- MAP rule:

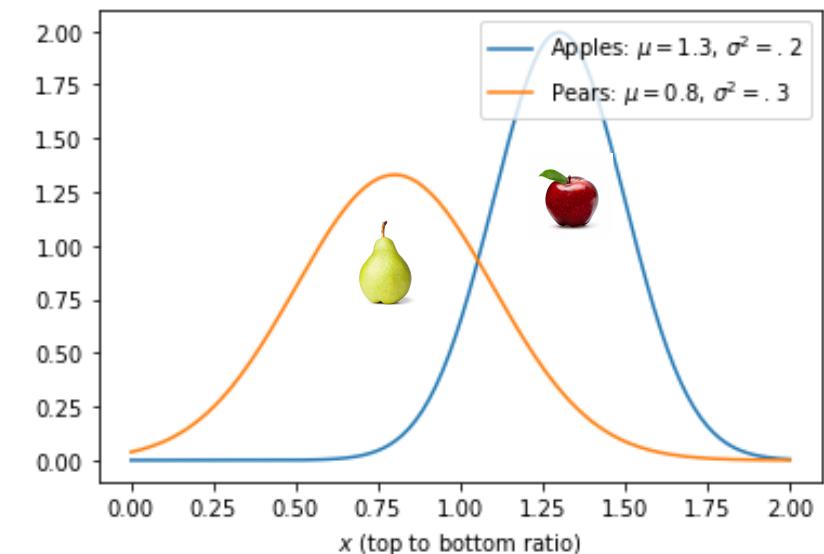
$$\hat{y} = \arg \max_y p(y|x)$$

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

- if

$$f(x|y = 1)p(y = 1) \geq f(x|y = 0)p(y = 0)$$

then  $\hat{y} = 1$  else  $\hat{y} = 0$ .



- new test fruit with  $x = 1.2$ :

# Learning in Low Dimensions

- how do we estimate  $p(\mathbf{x}|y)$  from examples?

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

- if  $\mathbf{x}$  low dimensional (i.e, apples vs. pears), no problem.

- continuous: approximate it with a suitable distribution (i.e, Gaussian)
- discrete: histogram estimate

$$p(\mathbf{x}|y) = \frac{\sum_{i=1}^n \mathbb{I}_{\{\mathbf{x}_i=\mathbf{x}, y_i=y\}}}{\sum_{i=1}^n \mathbb{I}_{\{y_i=y\}}}$$

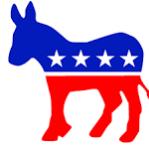
$$p(\mathbf{x}, y) = \frac{\sum_{i=1}^n \mathbb{I}_{\{\mathbf{x}_i=\mathbf{x}, y_i=y\}}}{n}$$

# Example: Digital Advertising and Tracking Data

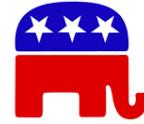
- example - user visited: cars.com, espn.com, hiking.com, stackexchange.com, ...

$$\mathbf{x} \in \mathbb{R}^{100,000} \quad \mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

- classification: given  $x$ , is user a democrat or a republican?



$y = 0$  if democrat



$y = 1$  if republican

- MAP rule:

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

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

# Shortcomings ...

- how do we estimate  $p(\mathbf{x}|y)$  from examples?

$$p(\mathbf{x}|y) = \frac{\sum_{i=1}^n \mathbb{I}_{\{\mathbf{x}_i=\mathbf{x}, y_i=y\}}}{\sum_{i=1}^n \mathbb{I}_{\{y_i=y\}}}$$

- what if  $\mathbf{x}$  is a  $28 \times 28$  black and white image?

$$\mathbf{x} = \boxed{9}$$

- $2^{28 \times 28}$  possible values for  $\mathbf{x}$
- MNIST: 60,000 examples

- what if  $\mathbf{x}$  browsing data, i.e.,  $\mathbf{x} \in \{0, 1\}^{100,000}$  ?

- $2^{100,000}$  possible values for  $\mathbf{x}$
- training data: 50,000 examples at most

- nearly all the values  $\mathbf{x}$  won't appear in training data

$$p(\mathbf{x}, y) = \frac{\sum_{i=1}^n \mathbb{I}_{\{\mathbf{x}_i=\mathbf{x}, y_i=y\}}}{n}$$

# Naive Bayes (Assume *Conditional* Independence)

- how can we estimate  $p(\mathbf{x}|y)$  when  $\mathbf{x}$  is high dimensional and we have limited training data?
- naïve Bayes: features  $x_1, x_2 \dots$  are independent.

$$p(\mathbf{x}|y) \approx p(x_1|y) \times p(x_2|y) \dots p(x_n|y) = \prod_{j=1}^n p(x_j|y)$$

- much easier to estimate marginals  $p(x_j|y)$  than joint  $p(\mathbf{x}|y)$

$$p(x_j|y) = \frac{\sum_{i=1}^n \mathbb{I}\{[\mathbf{x}_i]_j = x_j, y_i = y\}}{\sum_{i=1}^n \mathbb{I}\{y_i = y\}}$$

- MNIST: percentage of times pixel  $j$  is black (or white) when class is  $y = 9$ .



$$p(\mathbf{x}) \approx p(x_1) \times p(x_2) \times \dots \times p(x_n)$$

- last step - use MAP:

$$\hat{y} = \arg \max_y p(y|\mathbf{x}) = \arg \max_y \prod_j p(x_j|y)$$

# Shortcomings of NB

- independence assumption is often really bad

$$x = \boxed{9}$$

- example: MNIST even conditioned on  $y$ , pixel values are highly correlated!

- still works reasonably well for classification, in spite of this shortcoming

- sometimes there are still no examples of a particular value

- simple solution: add a 1 to both numerator and denominator of

$$p(x_j|y) = \frac{\sum_{i=1}^n \mathbb{I}\{[\boldsymbol{x}_i]_j = x_j\} + 1}{n + 1} \quad \text{Laplace smoothing}$$

- more precisely, add 1 to numerator and  $k$  to denominator, where  $|\mathcal{X}| = k$

- numerically, better to work in with  $\log p(\boldsymbol{x})$

$$\log p(\boldsymbol{x}) \approx \log \left( \prod_{j=1}^n p(x_j) \right) = \sum_{j=1}^n \log(p(x_j))$$

# MAP is Bayes Optimal

- simple idea: given the data or features  $x$ , what is the most likely class  $y$ ?
- MAP estimate:

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

- **claim:** MAP minimizes risk:

- *risk* is the expected *loss*

- misclassification loss:  $\ell(\hat{y}, y) = \mathbb{1}_{\{\hat{y} \neq y\}}$

$$\ell(\hat{y}, y) = \begin{cases} \hat{y} \neq y & 1 \\ \text{else} & 0 \end{cases}$$

$$E[\ell(\hat{y}, y)] = E[\mathbb{1}_{\{\hat{y} \neq y\}}]$$

$$= \mathbb{P}(\hat{y} \neq y)$$

$$= 1 - \mathbb{P}(\hat{y} = y)$$

$$= 1 - \sum_x p(x) \mathbb{P}(\hat{y} = y | x)$$

$$\geq 1 - \sum_x p(x) \max_y p(y | x)$$

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

$$\mathbb{P}(\hat{y} = y | x) \leq \max_y p(y | x) \text{ for any classifier } \hat{y}$$

where the last holds with equality if and only if *MAP*