

DSC 40A

Theoretical Foundations of Data Science I

Announcements

- Homework 7 was released and due 12/3 – no slip days.
- Gal's OH today at 4:15 instead of Wednesday.

Question

Answer at q.dsc40a.com

Remember, you can always ask questions at
q.dsc40a.com!

If the direct link doesn't work, click the "Lecture Questions" link in the top right corner of dsc40a.com.

Last Week

- We defined Bayes' Theorem:

$$P(\underline{B}|\underline{A}) = \frac{P(A|B) * P(\underline{B})}{P(A)} = \frac{P(A|B) P(B)}{P(A|B)P(B) + P(A|\bar{B})P(\bar{B})}$$

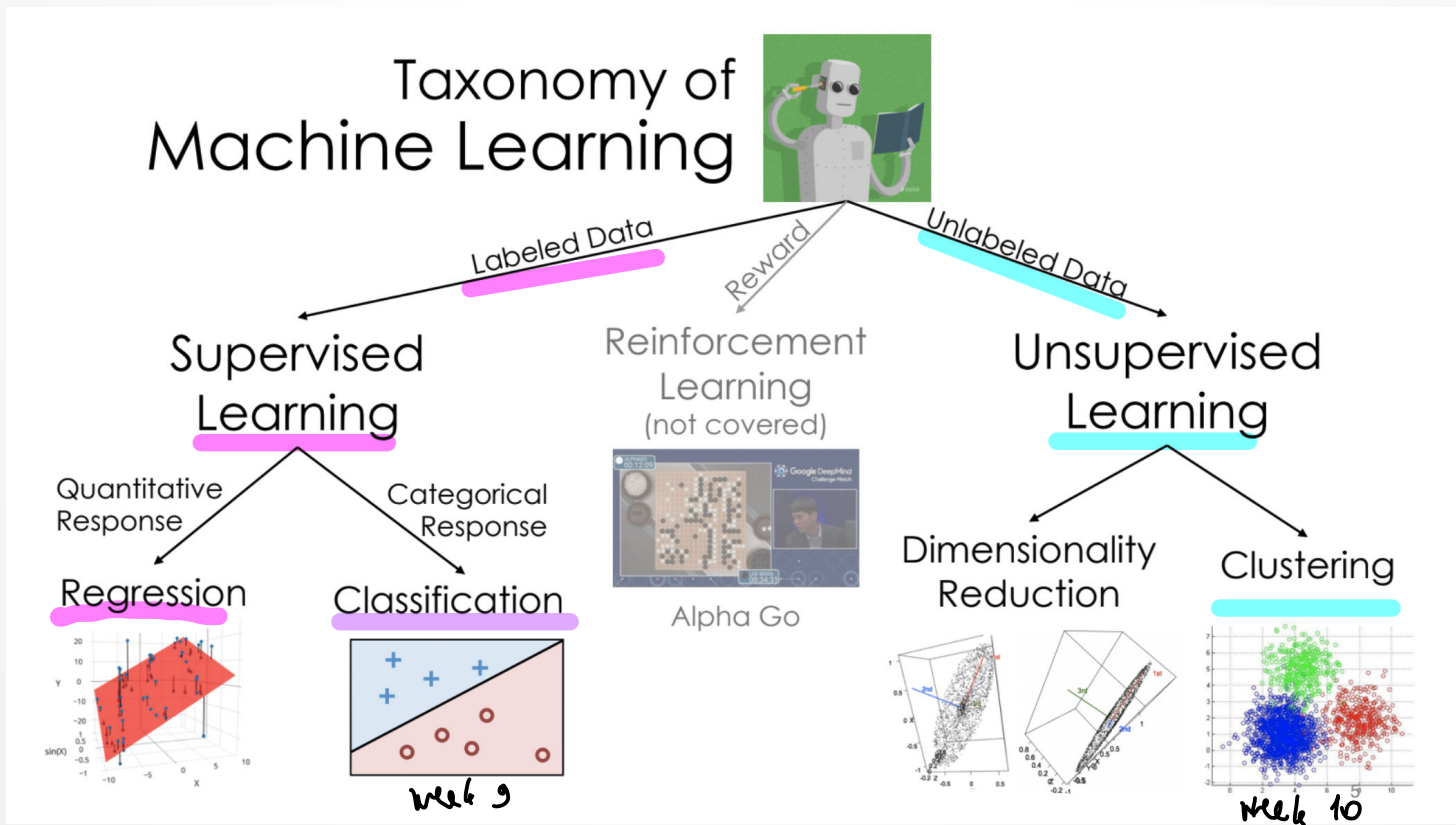
- Bayes' Theorem describes how to update the probability of one event given that another has occurred.

Naïve Bayes Classifier

The background of the slide features abstract, overlapping green geometric shapes, primarily triangles and polygons, in various shades of green, creating a modern and dynamic visual effect.

Today

- Using Bayes' Theorem to solve the classification problem



Preview: Bayes' Theorem for Classification

Bayes' Theorem is very useful for classification problems, where we want to predict a class based on some features.

$$P(B|A) = \frac{P(A|B) * P(B)}{P(A)}$$

B = belonging to a certain class

A = having certain features

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

Classification

- Making predictions based on examples (training data)
- Response variable is categorical
- Categories are called *classes*
- Examples:
 - decide whether patient has kidney disease *binary - yes, no*
 - identify handwritten digits *{0, 1, 2, 3, 4, 5, 6, 7, 8, 9} MNIST*
 - determine whether an avocado is ripe
 - predict whether credit card activity is fraudulent

Example

Color	Ripeness
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	unripe
purple-black	ripe
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	ripe
green-black	unripe
purple-black	ripe

You have a green-black avocado. Based on this data, would you predict that your avocado is ripe or unripe?

Which class would you predict?



A. ripe

B. unripe

Example

Color	Ripeness
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	unripe
purple-black	ripe
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	ripe
green-black	unripe
purple-black	ripe

You have a green-black avocado. Based on this data, would you predict that your avocado is ripe or unripe?

Strategy: Calculate two probabilities:

$$P(\text{ripe} \mid \text{green-black}) = \frac{\# \text{ripe and green black}}{\# \text{green black}} = \frac{3}{5}$$

$$P(\text{unripe} \mid \text{green-black}) = \frac{2}{5}$$

Then choose the class according to the **larger** of these two probabilities.

$$\frac{3}{5} > \frac{2}{5} \Rightarrow \text{ripe}$$

Bayes' Theorem for Classification

Bayes' Theorem gives another strategy for predicting the class given features.

$$P(B|A) = \frac{P(A|B) * P(B)}{P(A)}$$

A = have certain feature
B = belong to a class

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

*estimate this
data*

if not

Bayes' Theorem for Classification

Bayes' Theorem gives another strategy for predicting the class given features.

$$P(B|A) = \frac{P(A|B) * P(B)}{P(A)}$$

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

Can all be
estimated
from the
training data

Avocado Ripeness

Color	Ripeness
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	unripe
purple-black	ripe
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	ripe
green-black	unripe
purple-black	ripe

You have a green-black avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

$$P(\text{green black} | \text{ripe}) = \frac{3}{7}$$

$$P(\text{ripe}) = \frac{7 \text{ ripe avo.}}{11 \text{ total avo.}}$$

$$P(\text{green-black}) = \frac{5}{11}$$

$$P(\text{ripe} | \text{green-black}) =$$

$$\frac{P(\text{green-black} | \text{ripe}) P(\text{ripe})}{P(\text{green-black})} = \frac{\frac{3}{7} \cdot \frac{7}{11}}{\frac{5}{11}} = \frac{3}{5}$$

$$P(\text{unripe} | \text{green-black}) = \dots = \frac{2}{5}$$

Avocado Ripeness

Color	Ripeness
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	unripe
purple-black	ripe
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	ripe
green-black	unripe
purple-black	ripe

You have a green-black avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

$$P(\text{ripe} | \text{green-black})$$

v / ^ ?

$$P(\text{unripe} | \text{green-black})$$

} have same denominator
we don't need to calc.
just compare numerators

Avocado Ripeness

Color	Ripeness
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	unripe
purple-black	ripe
bright green	unripe
green-black	ripe
purple-black	ripe
green-black	ripe
green-black	unripe
purple-black	ripe

You have a green-black avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

Shortcut: Both probabilities have same denominator. To find larger one, choose one with larger numerator.

$$P(\text{ripe} | \text{green-black}) \propto \frac{3}{7} \cdot \frac{7}{11} = \frac{3}{11}$$

(proportional)

$$P(\text{unripe} | \text{green-black}) \propto \frac{2}{4} \cdot \frac{4}{11} = \frac{2}{11}$$

multiple features

response variable (categorical)

More Features

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

Strategy: Calculate two probabilities:

$P(\text{ripe} \mid \text{firm, green-black, Zutano})$

Problem: no firm, gb, Zutano avo in data

$P(\text{unripe} \mid \text{firm, green-black, Zutano})$

Then choose the class according to the **larger** of these two probabilities.

Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

Problem: We have not seen an avocado with all these features. Both probabilities will be undefined.

$$\begin{aligned}
 &P(\text{ripe} \mid \text{firm, green-black, Zutano}) \\
 &= \frac{P(\text{ripe} \ \& \ \text{firm, gb, zutano})}{P(\text{firm, gb, zutano})} \longrightarrow \text{undefined} \Rightarrow \frac{0}{0} \\
 &P(\text{unripe} \mid \text{firm, green-black, Zutano})
 \end{aligned}$$

Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

Solution: Use Bayes' Theorem, plus a simplifying assumption, to calculate the two numerators.

Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

Simplifying assumption: Within a given class, the features are independent.

conditional independence assumption

$P(\text{firm, green-black, Zutano} \mid \text{ripe}) =$

$P(\text{firm} \mid \text{ripe}) * P(\text{green-black} \mid \text{ripe}) * P(\text{Zutano} \mid \text{ripe})$

Can calculate even if we don't have an avocado with all 3 features

Conditional Independence

- Recall that A and B are independent if

$$P(A \text{ and } B) = P(A) * P(B)$$

- A and B are conditionally independent given C if

$$P((A \text{ and } B)|C) = P(A|C) * P(B|C)$$

feature *class* *feat 1* *class* *feat 2* *class*

- Given that C occurs, this says that A and B are independent of one another.

Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

$P(\text{ripe} | \text{green-black, firm, Zutano})$

$P(\text{features}|\text{class})$ (conditional independence)

$P(\text{class})$ (prior)

$$P(\text{firm}|\text{ripe}) \cdot P(\text{gb}|\text{ripe}) \cdot P(\text{Zutano}|\text{ripe})$$

$$= \frac{1}{7} \cdot \frac{3}{7} \cdot \frac{2}{7}$$

$$P(\text{ripe}|\text{features}) \propto \frac{1 \cdot 3 \cdot 2}{7 \cdot 7 \cdot 7} \cdot \frac{7}{11} = \frac{6}{539}$$

Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

common mistake

$$P(A|B) = 1 - P(\bar{A}|B) \neq 1 - P(A|\bar{B})$$

Assuming conditional independence of features given the class, calculate $P(\text{firm, green-black, Zutano} | \text{unripe})$.

- A. 0
- B. $1/4$
- C. $3/16$
- ☒ D. $1 - (1/7 * 3/7 * 2/7) \neq 1 - P(\text{firm, gb, zutano} | \text{ripe})$

Avocado Ripeness

Color	Softness	Variety	Ripeness
bright green	firm	Zutano	unripe
green-black	medium	Hass	ripe
purple-black	firm	Hass	ripe
green-black	medium	Hass	unripe
purple-black	soft	Hass	ripe
bright green	firm	Zutano	unripe
green-black	soft	Zutano	ripe
purple-black	soft	Hass	ripe
green-black	soft	Zutano	ripe
green-black	firm	Hass	unripe
purple-black	medium	Hass	ripe

You have a firm green-black Zutano avocado. Based on this data, would you predict that your avocado is ripe or unripe?

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

Naïve Bayes Algorithm

- Bayes' Theorem shows how to calculate $P(\text{class} \mid \text{features})$.

$$P(\text{class} \mid \text{features}) = \frac{P(\text{features} \mid \text{class}) * P(\text{class})}{P(\text{features})}$$

- Rewrite the numerator, using the naïve assumption of conditional independence of features given the class.
- Estimate each term in the numerator based on the training data.
- Select class based on whichever has the larger numerator.

Summary

- The Naïve Bayes algorithm gives a strategy for classifying data according to its features.
 - It relies on an assumption of conditional independence of the features.
 - **Next time:** application to text classification
-