# CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023

# Admin

- Lab 5 due **Wednesday** (tomorrow)
- Lab 6 posted tomorrow
- **Midterm 1** returned today

- **Lab today**: Lab 5 implementation advice and check-ins
  - If you're *completely* finished, don't need to attend, but please email me
  - Otherwise will check in about Lab 5

# Lab 5 implementation

*Partition contains:*

- Features dictionary F:

   F = {**age**: [Senior, Middle-age, Mid-adult, Young-adult, Child], **workclass**: [Private, Local-gov…] … }

- List of Examples
  - Each example contains

     features = {**age**: Senior, **workclass**: Private … }

     label = 1 (Female)

# Outline for October 24

- Entropy and Shannon encoding

- Information gain for selecting features

- Go over Midterm 1

- Continuous features (if time)

# Outline for October 24

- Entropy and Shannon encoding

- Information gain for selecting features

- Go over Midterm 1

- Continuous features (if time)

# Decision Trees use entropy to select best features

**Examples**

- ## Medical diagnostics

- ## Credit risk analysis

- ## Modeling calendar scheduling preferences

# Decision Trees in Chemistry reactions

- Example of decision trees in practice
- Use decision trees to interpret another ML algorithm (SVMs)

## Machine-learning-assisted materials discovery using failed experiments

Paul Raccuglia, Katherine C. Elbert, Philip D. F. Adler, Casey Falk, Malia B. Wenny, Aurelio Mollo, Matthias Zeller, Sorelle A. Friedler ✉, Joshua Schrier ✉ & Alexander J. Norquist ✉
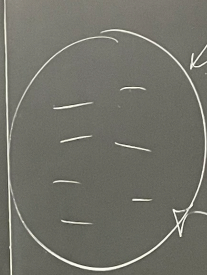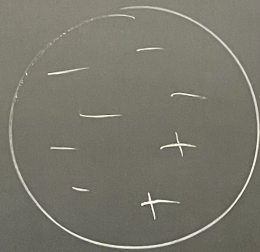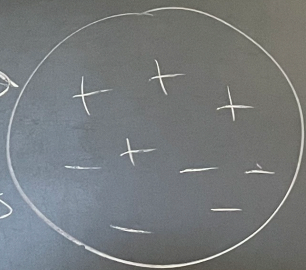
Optional Reading!

# How do we choose the best feature?

- Single feature model + evaluate with a ROC curve **(Lab 4)**

- What feature gives us the most information about the label? **(Lab 6)**

# Idea of Entropy : avg # of bits needed to send one datapoint



posinous + edible mushrooms

← all edible

0 entropy

(1 bit)

high entropy

low entropy (what we want)

(i.e. grocery store)

send information

1 1 0 0 1 0 0 0 0 0 1 1 0

↑ junior ↑ junior ↑ junior ↑ senior/junior soph?

---

binary ..... □·$2^2$ + □·$2^1$ + □·$2^0$ + □·$2^{-1}$ + □·$2^{-2}$ ...

$5 = 1·4 + 0·2 + 1·1$

⇒ boxed: 101 in binary

decimal point! ↑ $\frac{1}{2}$ ↑

---

$5.5 \Rightarrow$ 101.1

$\frac{1}{8}$ $\frac{1}{4}$ $\frac{1}{2}$

-1

-2

half as likely ≠ twice as many bits!

pos ed

**fixed len encoding**

00

01

10 ~~sort~~

11

works!

| | year | prob (p) | idea | cummulative prob | in binary | $\lceil -\log_2(p) \rceil$ | Shannon len encoding |
|---|------|----------|------|------------------|-----------|----------------------------|----------------------|
| | senior | 0.5 | 0 | 0 | 0.0̲0̲0̲0... | 1 | 0 |
| | junior | 0.25 | 1 | 0.5 | 0.1̲0̲0̲... | 2 | 10 |
| | soph | 0.125 | 01 | 0.75 | 0.1̲1̲0̲... | 3 | 110 |
| | first | 0.125 | 10 | 0.875 | 0.1̲1̲1̲0̲... | 3 | 111 |

Sum to 1

# bits

no code is the prefix of another

# Entropy

$$H(y) = -\sum_{c \in \text{vals}(y)} p(y=c) \log_2 (\underbrace{p(y=c)}_{\text{\# bits}})$$

$\uparrow$
label

$$H(\text{year}) = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3$$

$$= \boxed{1.75} \text{ bits}$$

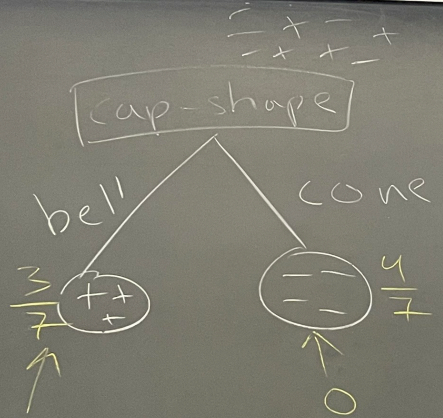$$\neq \frac{1+2+3+3}{4} = 2.25$$

# Outline for October 24

- Entropy and Shannon encoding

- Information gain for selecting features

- Go over Midterm 1

- Continuous features (if time)

$$\begin{array}{cccc} - & + & - & + \\ - & + & + & - \end{array}$$

$$\boxed{\text{cap-shape}}$$

bell        cone

$\frac{3}{7}$   (+ +)     (- -)   $\frac{4}{7}$

                0

$$H(Y|X=\text{bell}) = -\left( 0 \cdot \log 0 + 1 \cdot \log 1 \right)$$

$$= 0$$

single
feature
value

Conditional entropy

$$H(Y|X) = \sum_{v \in \text{vals}(X)} P(X=v) H(Y|X=v) = \frac{3}{7} \cdot 0 + \frac{4}{7} \cdot 0 = 0$$

ex
cap-shape

?

weighted avg

same

$$H(Y|X=v) = -\sum_{c \in \text{vals}(Y)} P(Y=c|X=v) \log_2 P(Y=c|X=v)$$

$y=0$

cap-shape = bell

all

$=0$

Handout 13

all values/leaves

Information gain

$$G(Y, X) = H(Y) - H(Y|X)$$

Want low

Want high

Want the feature that maximizes info gain

# Handout 13

Handout 13

$L_i =$

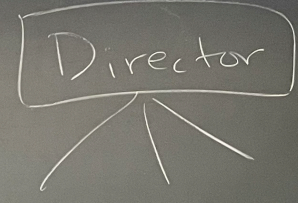① $P(\text{yes}) = \frac{6}{9} = \frac{2}{3}$

for Lab6
$\Rightarrow$ no
rounding

② $H(L_i) = -\left(\frac{2}{3}\log_2\frac{2}{3} + \frac{1}{3}\log_2\frac{1}{3}\right) \approx 0.92$

③ $\text{Gain}(L_i, T) = 0.92 - 0.61 = 0.31$

$0.92 - 0.61 \ldots$

Director

| Movie | Type | Length | Director | Famous actors | Liked? |
|-------|------|--------|----------|---------------|--------|
| m1 | Comedy | Short | Adamson | No | Yes |
| m2 | Animated | Short | Lasseter | No | No |
| m3 | Drama | Medium | Adamson | No | Yes |
| m4 | Animated | Long | Lasseter | Yes | No |
| m5 | Comedy | Long | Lasseter | Yes | No |
| m6 | Drama | Medium | Singer | Yes | Yes |
| m7 | Animated | Short | Singer | No | Yes |
| m8 | Comedy | Long | Adamson | Yes | Yes |
| m9 | Drama | Medium | Lasseter | No | Yes |

$P(Li = yes) =$ **2/3**

$H(Li) =$ **0.92**

$H(Li \mid T) = 0.61$
$H(Li \mid Le) = 0.61$
$\boxed{H(Li \mid D) = 0.36}$  MIN ENTROPY
$H(Li \mid F) = 0.85$

$Gain(Li, T) =$ **0.92 − 0.61 = 0.31**
$Gain(Li, Le) =$ **0.92 − 0.61 = 0.31**
$\boxed{Gain(Li, D) =}$ **0.92 − 0.36 = 0.56**  MAX INFO GAIN
$Gain(Li, F) =$ **0.92 − 0.85 = 0.07**

Director

Start of the tree

# Outline for October 24

- Entropy and Shannon encoding

- Information gain for selecting features

- Go over Midterm 1

- Continuous features (if time)

# Midterm 1 Grades

- 90-100%    A
- 80-89%     B
- 70-79%     C
- Below 70%: please meet with me
- Below 60%: not passing

- Any questions about the exam: bring to me within one week

Midterm solutions
not posted online

# Outline for October 24

- Entropy and Shannon encoding

- Information gain for selecting features

- Go over Midterm 1

- Continuous features (if time)