

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Spring 2025



Haverford
COLLEGE

Admin

- Lab 5 due **Wednesday** (tomorrow)
- Lab 6 posted
- **Midterm 1** returned today
- No office hours on Wednesday (tomorrow)
- **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 or start Lab 6

Lab 5 implementation

Partition contains:

- Features dictionary F:

$F = \{\text{age: [Senior, Middle-age, Mid-adult, Young-adult, Child]}, \text{workclass: [Private, Local-gov...]} \dots \}$

- List of Examples

- Each example contains

- features = {age: Senior, workclass: Private ... }

- label = 1 (Female)

defaultdict

```
from collections import defaultdict

# Create a defaultdict with int as the default factory (default value is 0)
count_dict = defaultdict(int)

# Increment the count for some items
items = ['apple', 'banana', 'apple', 'orange', 'banana', 'apple']
for item in items:
    count_dict[item] += 1

# Print the counts
print(count_dict)
# Expected output: defaultdict(<class 'int'>, {'apple': 3, 'banana': 2, 'orange': 1})

# Access a non-existent key
print(count_dict['grape'])
# Expected output: 0 (because int() returns 0)
```

Outline

- Entropy and Shannon encoding
- Information gain for selecting features
- Go over Midterm 1
- Continuous features (if time)

Outline

- 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



[Journal of Medical Systems](#)
October 2002, Volume 26, [Issue 5](#), pp 445–463 | [Cite as](#)

Decision Trees: An Overview and Their Use in Medicine

Authors [Authors and affiliations](#)

Vili Podgorelec , Peter Kokol, Bruno Stiglic, Ivan Rozman

- Credit risk analysis



[Computational Economics](#)
April 2000, Volume 15, [Issue 1–2](#), pp 107–143 | [Cite as](#)

Credit Risk Assessment Using Statistical and Machine Learning: Basic Methodology and Risk Modeling Applications

Authors [Authors and affiliations](#)

J. Galindo, P. Tamayo

- 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✉

Nature **533**, 73–76 (05 May 2016) | [Download Citation](#) ↓

How do we choose the best feature?

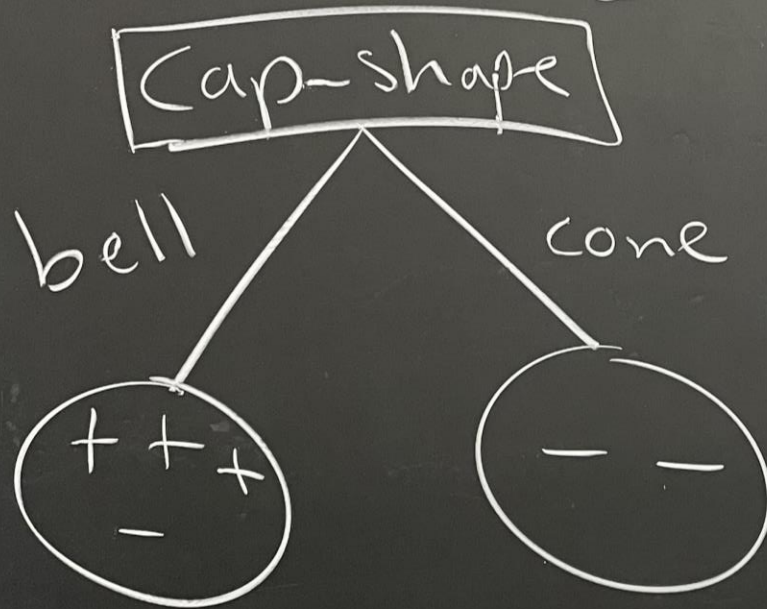
- Single feature model + evaluate with a ROC curve (**Lab 4**)
- What feature gives us the most information about the label? (**Lab 6**)

Information Theory

How do we choose the best feature?

Lab 4

~~+~~ edible
~~+~~ poisonous
- + -
+ = -

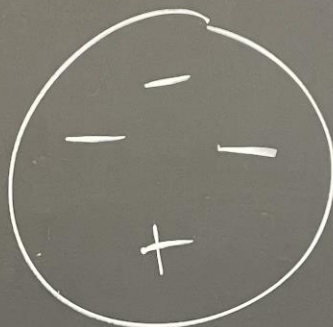


gained
information

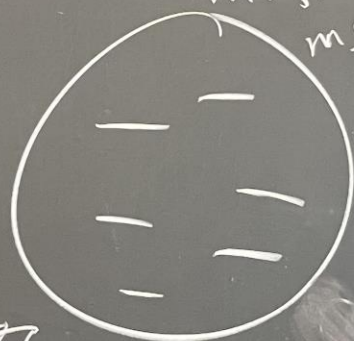
ature?

Idea of entropy: avg # of bits
needed to send one data pt

high impurity
"chaos"



all edible
mushrooms



entropy
○
(○ bits)

How to find # bits?

year	prob (p)	idea	Cumulative prob	in binary!	# digits to take $\lceil -\log_2(p) \rceil$	code (Shannon)
Senior	0.5	0	0	0. 0 00...	1	0
junior	0.25	1	0.5	0. 1 000...	2	10
soph	0.125	01	0.75	0. 11 00...	3	110
first-year	0.125	10	0.875	0. 111 0...	3	111
Sorted high \rightarrow low			1			

example first? decode!

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

↑ junior ↑ junior?

instead prefix encoding

no code is the prefix of another code

binary

$$\dots \square \cdot 2^2 + \square \cdot 2^1 + \square \cdot 2^0 + \square \cdot 2^{-1} + \square \cdot 2^{-2} \dots$$

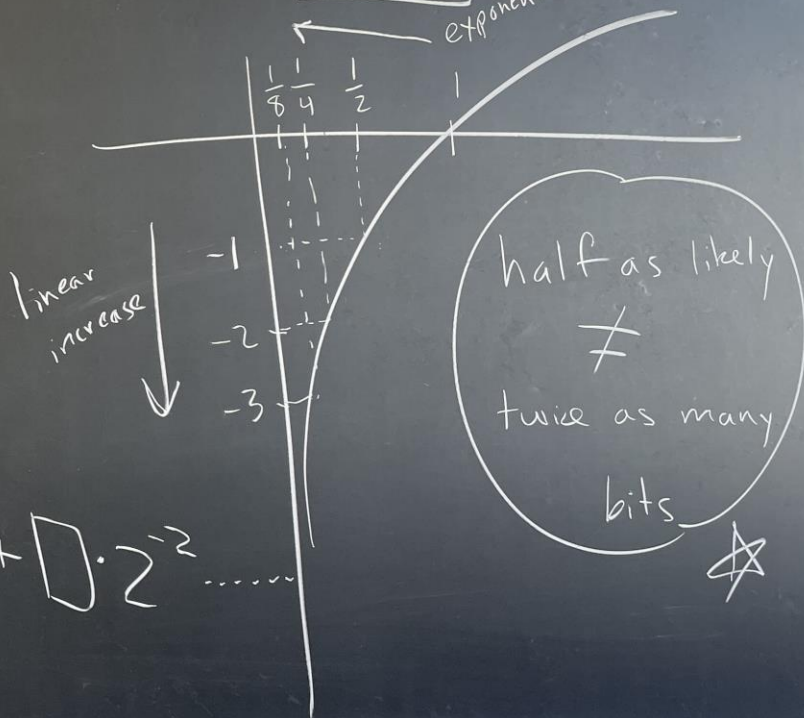
$$5 = 1 \cdot 4 + 0 \cdot 2 + 1 \cdot 1$$

$$\Rightarrow \boxed{101} \text{ in binary}$$

5.5

$\Rightarrow \boxed{101.1}$

exponential decay



Outline

- Entropy and Shannon encoding
- Information gain for selecting features
- Go over Midterm 1
- Continuous features (if time)

Entropy

$$H(Y) = - \sum \underbrace{p(Y=c)}_{\substack{\text{(Evals(Y))} \\ \text{prob or} \\ \text{freq of} \\ c}} \underbrace{\log_2 p(Y=c)}_{\substack{\text{\# bits} \\ \text{for } c}}$$

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

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

not

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

conditional entropy

$$H(Y|X) = \sum_{v \in \text{vals}(X)} p(x=v) \underbrace{H(Y|X=v)}_?$$

label
feature

$v \in \text{vals}(X)$

conditional entropy of one feature value (leaf)

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

ex

$$H(Y| \overset{x}{\text{cap-shape}} = \overset{v}{\text{bell}})$$

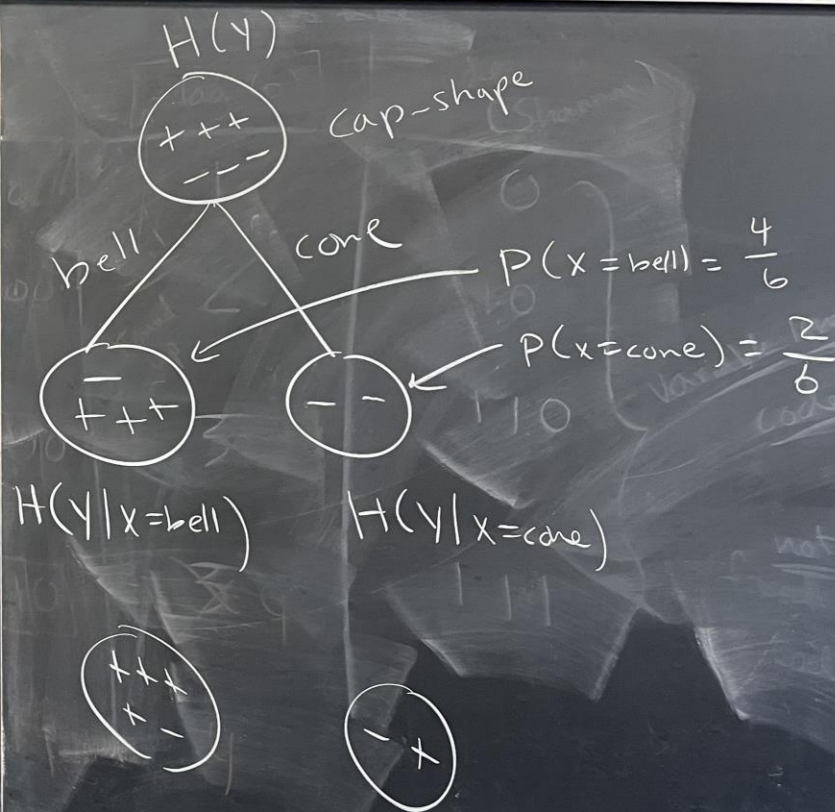
$$= - \left(\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} \right)$$

$$\approx \boxed{0.811}$$

Information Gain

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

Want high
Want low



Handout 13

Handout 13

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	
m9	Drama	Medium	Lasseter	No	

$$P(Li = \text{yes}) = 2/3$$

$$H(Li) = 0.92$$

$$H(Li | T) = 0.61$$

$$H(Li | Le) = 0.61$$

$$H(Li | D) = 0.36 \quad \text{MIN ENTROPY}$$

$$H(Li | F) = 0.85$$

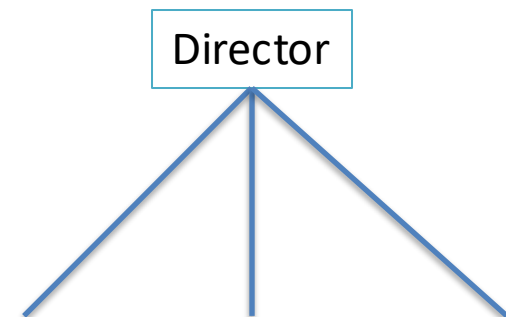
$$\text{Gain}(Li, T) = 0.92 - 0.61 = 0.31$$

$$\text{Gain}(Li, Le) = 0.92 - 0.61 = 0.31$$

$$\text{Gain}(Li, D) = 0.92 - 0.36 = 0.56 \quad \text{MAX INFO GAIN}$$

$$\text{Gain}(Li, F) = 0.92 - 0.85 = 0.07$$

$$H(Li) = -\left(\frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3}\right) = 0.92$$



Start of the tree

Outline

- Entropy and Shannon encoding
- Information gain for selecting features
- **Go over Midterm 1**
- Continuous features (if time)

Midterm 1 Grades

Median: 88

- 88-100% A
- 78-87% B
- 68-77% C
- 58-67% D
- Below 58%: not passing, please meet with me
- Note: as per the syllabus, you must pass at least one exam to pass the course
- Any questions about the exam: bring to me within one week

Quote of the week

“Enjoy the little things, for one day you may look back and realize they were the big things.”

--Robert Breault