

# CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2019



**HAVERFORD**  
COLLEGE

# Admin

- Lab 2 **TODAY**
- Next office hours **Friday 3-5pm**
- **TA hours** (in H204)
  - Sunday 7-8pm (Pablo)
  - Monday 7-8pm (Charlie)

# Outline for September 12

- Reading quiz
- Continue entropy and information gain
- Continuous features
- Lab 2 implementation suggestions
- Learning problem so far + terminology

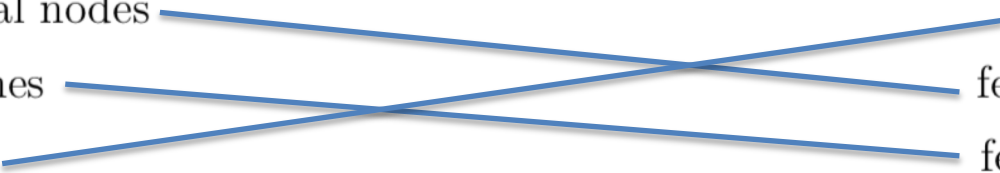
# Outline for September 12

- Reading quiz
- Continue entropy and information gain
- Continuous features
- Lab 2 implementation suggestions
- Learning problem so far + terminology



# Reading Quiz #2

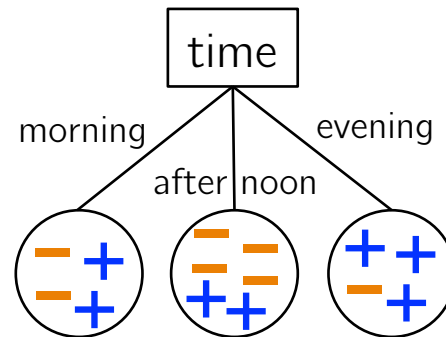
1)

- internal nodes
  - branches
  - leaves
- class labels
  - feature names
  - feature values
- 

# Reading Quiz #2

- 1)
- internal nodes
  - branches
  - leaves
- class labels  
feature names  
feature values

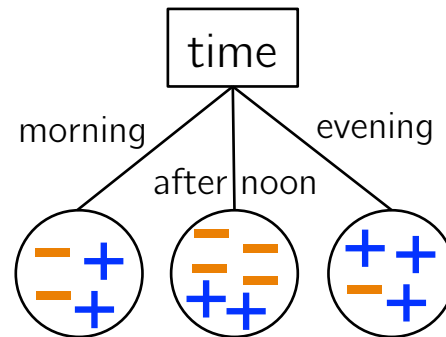
2) (a) +  
(b) 5/14



# Reading Quiz #2

- 1)
- internal nodes
  - branches
  - leaves
- class labels  
feature names  
feature values

- 2) (a) +  
(b) 5/14

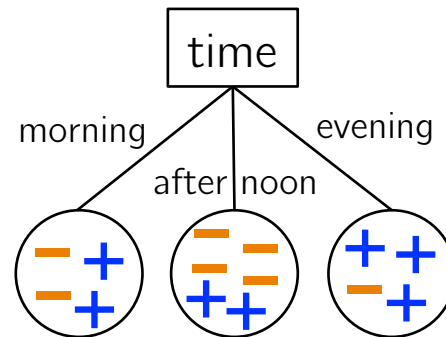


- 3) high

# Reading Quiz #2

- 1)
- internal nodes
  - branches
  - leaves
- class labels  
feature names  
feature values

- 2) (a) +  
(b) 5/14

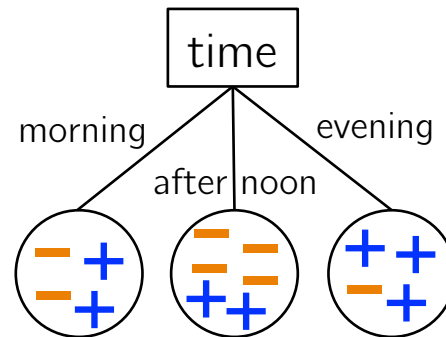


- 3) high  
4) zero/one loss

# Reading Quiz #2

- 1)
- internal nodes
  - branches
  - leaves
- class labels  
feature names  
feature values

- 2) (a) +  
(b) 5/14



- 3) high
- 4) zero/one loss
- 5) Case when no training examples have the feature value

# Outline for September 12

- Reading quiz
- Continue entropy and information gain
- Continuous features
- Lab 2 implementation suggestions
- Learning problem so far + terminology

# Real-World 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](#) ↓



Note: we can put the negative sign inside or outside the ceiling function, since the ceiling function rounds up in an absolute value sense. i.e.  $\text{ceil}(2.25) = 3$  and  $\text{ceil}(-2.25) = -3$ .

Year	prob (p)	cumulative prob	in binary	$\lceil -\log_2(p) \rceil$	code
Senior	0.5	0	0.000...	1	0
junior	0.25	0.5	0.100...	2	10
soph	0.125	0.75	0.110...	3 ceiling (round up)	110
first	0.125	0.875	0.1110...	3	111

entropy: avg # of bits needed to transmit one example.

inside or outside

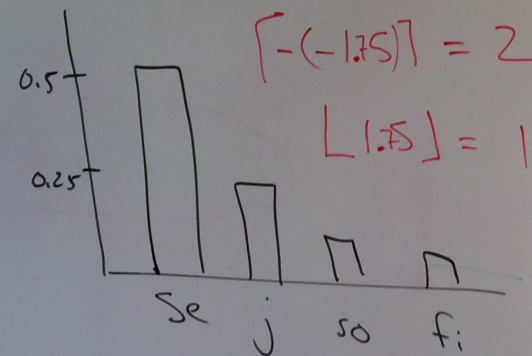
110 110 100 01110

Shannon # digits encoding.

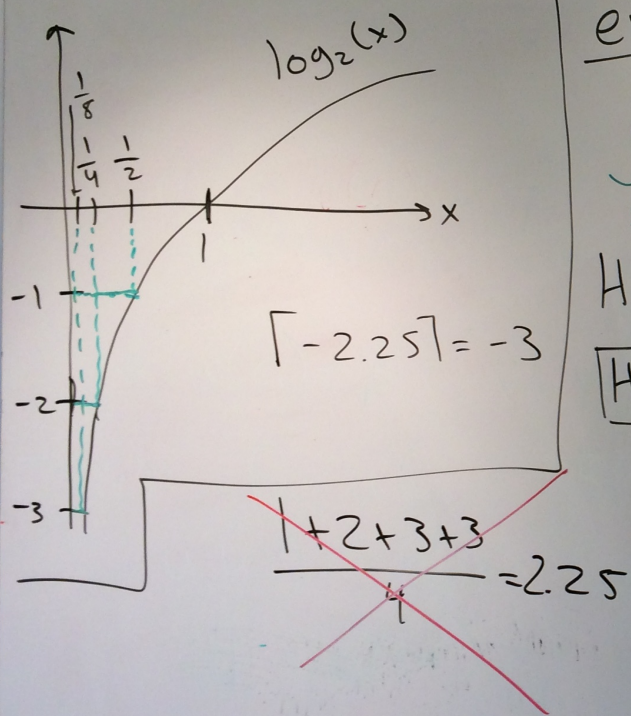
binary:  $\square \cdot 2^2 + \square \cdot 2^1 + \square \cdot 2^0 + \square \cdot 2^{-1} + \square \cdot 2^{-2} \dots$

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

1948







## entropy

$$H(Y) = - \sum_{c \in \text{Evals}(Y)} P(Y=c) \cdot \log_2(P(Y=c))$$

function 2  $c \in \text{Evals}(Y)$       function 1  $P(Y=c)$       # of bits  $\log_2(P(Y=c))$

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

$$H(\text{year}) = 1.75$$

## Conditional Entropy

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

class  $Y$       one feature  $X=v$       function 3  $P(Y=c|X=v)$

$$P(Y=\text{yes} | X=\text{sunny}) = \frac{P(\overset{A}{\text{yes}} \text{ AND } \overset{B}{\text{sunny}})}{P(\overset{B}{\text{sunny}})} = \frac{2}{5}$$

$\frac{2/14}{5/14}$

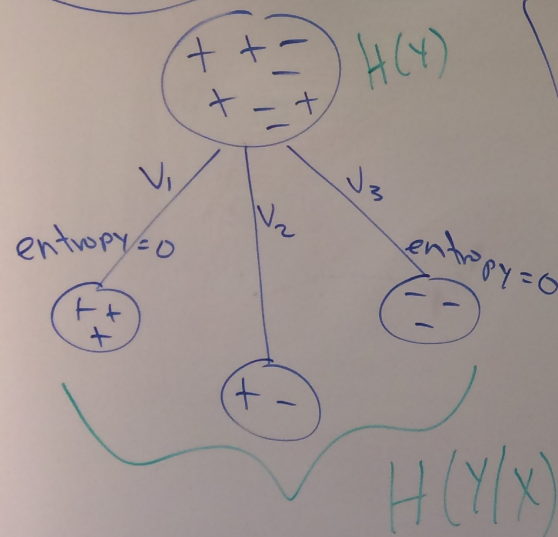


# Practice Problem

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

function 4

Information Gain



$$\text{Gain}(Y, X) = H(Y) - H(Y|X)$$

want to be high! starts high want low!

function 5

start at root  
Goal:  $H(Li|FA)$   
famous actors

start  
 $P(Li = \text{Yes} | FA = \text{No})$

$P(Li = \text{No} | FA = \text{No})$

0.85

Handout 4: work with your group!

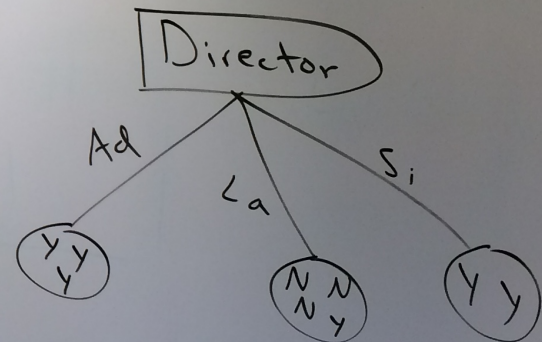
# Handout 4

Movie	Type	Len	Director	Famous	Liked?
$m_1$	Com	S	Ad	N	Y
$m_2$	Ani	S	La	N	<del>Y</del> N
$m_3$	Dra	M	Ad	N	Y
$m_4$	Ani	L	La	Y	<del>Y</del> N
$m_5$	Com	L	La	Y	<del>Y</del> N
$m_6$	Dra	M	Si	Y	Y
$m_7$	Ani	S	Si	N	Y
$m_8$	Com	L	Ad	Y	Y
$m_9$	Dra	M	La	N	Y

$$H(Y) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right)$$

$$\approx 0.92$$

$$\text{Gain}(\tilde{L}_i, T) = 0.92 - 0.61 = 0.31$$





# Outline for September 12

- Reading quiz
- Continue entropy and information gain
- **Continuous features**
- Lab 2 implementation suggestions
- Learning problem so far + terminology

Continuous Features

Continuous

$x_i$	$y$ ← label
10	Y
7	Y
8	N
3	Y
7	N
12	Y
2	Y

① Sort feature values

2	3	7	7	8	10	12
Y	Y	Y	N	N	Y	Y

②

2	3	7	8	10	12	} unique values
Y	Y	None	N	Y	Y	

③ make a feature whenever the label changes

~~~~~\*~~~~~

|                |
|----------------|
| $x_j \leq 5$   |
| $x_j \leq 7.5$ |
| $x_j \leq 9$   |

# Outline for September 12

- Reading quiz
- Continue entropy and information gain
- Continuous features
- **Lab 2 implementation suggestions**
- Learning problem so far + terminology



# Implementation Suggestions

- Start slow with **entropy**! Build up function by function
- Think back to **trees in data structures**
- Distinguish between **data** (X,y) and **options for data** (values for each feature, classes for y)

# Outline for September 12

- Reading quiz
- Continue entropy and information gain
- Continuous features
- Lab 2 implementation suggestions
- Learning problem so far + terminology

*Next time!*