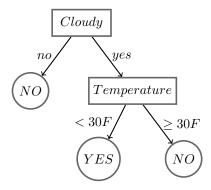
1 Definitions Oh My!

(a) **Decision Trees**: Popular representation and technique for classification and prediction. Each node represents an attribute that is tested (or a decision), each branch represents an outcome, and each leaf node represents a resulting label.

Example: Will it be rainy today?



(b) **Entropy**: A measurement of the uncertainty in a random variable. Hence, high entropy means we are less confident about the outcome.

Note: A random variable assigns numerical values to outcomes of a random experiment or process. For example, a coin flip (F) is a random variable where $P(F = heads) = \frac{1}{2}$.

$$H(Y) = -\sum_y P(Y=y) \log_2 P(Y=y)$$

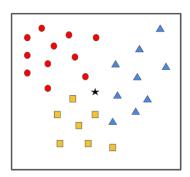
(c) Conditional-Entropy: The expected value of the entropy of Y given X, over all values of X. This lets us quantify the entropy of Y given that we know X.

$$H(Y\mid X) = \sum_{x} P(X=x) H(Y\mid X=x)$$

(d) **Mutual Information**: Measurement of how much uncertainty about variable Y decreases if we know and have observed X.

$$I(Y;X) = H(Y) - H(Y \mid X)$$

(e) **K-Nearest Neighbors (KNN)**: Classification algorithm in which the predicted label for a particular data point is determined by the labels of the k nearest neighbors to that point.



2 Decision Trees

	2.1	Mutual	Information	and Entropy	Practice
--	-----	--------	-------------	-------------	----------

1.	Calculate the entropy of tossing a fair coin.
2.	Calculate the entropy of tossing a coin that lands only on tails. $\underline{note}: 0 \cdot \log_2 0 = 0.$
3.	Calculate the entropy of a fair dice roll.
4.	When is the mutual information $I(Y;X) = 0$?

2.2 Decision Tree Runthrough

Fun fact: Pat loves to go on runs! However, sometimes things come up which deter him from running outside. The following attributes have affected his running in the past:

- How the sky looks the day of his run (sunny, overcast, rainy)
- The temperature outside (hot, mild, cool)
- The humidity 30 minutes before he plans on running (high, normal)

As a member of the 10-315 running club, you have gathered the following data and decided that you want to make a decision tree to predict whether Pat will go on a run.

Outlook (X_1)	Temperature (X_2)	Humidity (X_3)	Go on run? (Y)
sunny	hot	high	no
overcast	hot	high	yes
rainy	mild	high	yes
rainy	cool	normal	yes
sunny	mild	high	no
sunny	mild	normal	yes
rainy	mild	normal	yes
overcast	hot	normal	yes

To start, let's first determine the attribute for the root node of the decision tree. To do this, we need to determine the mutual information for each attribute.

1.	Calculate $H(Y)$, the entropy of Pat going on a run.
0	
2.	Calculate $I(Y; X_1)$, the mutual information when splitting on the outlook.

3.	Calculate $I(Y; X_2)$, the mutual information when splitting on temperature.
4.	Calculate $I(Y; X_3)$, the mutual information when splitting on humidity.
5.	Which attribute do we split on?
6.	What does the decision tree look like so far?
	l

Now that we've split on the first attribute, the remaining data points which still cannot be classified are shown in the table below.

Outlook (X_1)	Temperature (X_2)	Humidity (X_3)	Go on run? (Y)
sunny	hot	high	no
sunny	mild	high	no
sunny	mild	normal	yes

From now on, we will be strictly considering the branch of the decision tree where the outlook is sunny.

7.	Calculate $H(Y)$, the entropy of going on a run given that it is sunny.
8.	Calculate $I(Y; X_2)$.
9.	Calculate $I(Y; X_3)$.
10.	Which attribute do we split on now?
11.	Draw the final decision tree.

2.3 Discussion Questions

Explore different maximum depths of a decision tree by running through the DT Jupyter Notebook. Let d be the max-depth of the decision tree

1.	at is the best value of d for this dataset?		
2.	At which value of d do the training and validation errors begin to diverge?		
3.	How is the value of d related to overfitting?		

3 K Nearest Neighbor

Explore different values for "k" by running through the kNN Jupiter Notebook (see recitation materials).

3.1	Discussion Questions
1.	What is the best value of k for this dataset?
2.	Why do you think the error increases as k gets larger?
3.	How is the value of k related to overfitting?
4.	Why can we not display the decision boundaries like we saw in lecture?
5.	For which values of k would the decision boundaries be the most complex?

4 Putting it Together

Which methods can achieve zero-training error on this dataset?

- (a) Decision Trees
- (b) 1-Nearest Neighbor
- (c) Both
- (d) Neither

Bonus: If zero error, what could the decision boundary look like?

