



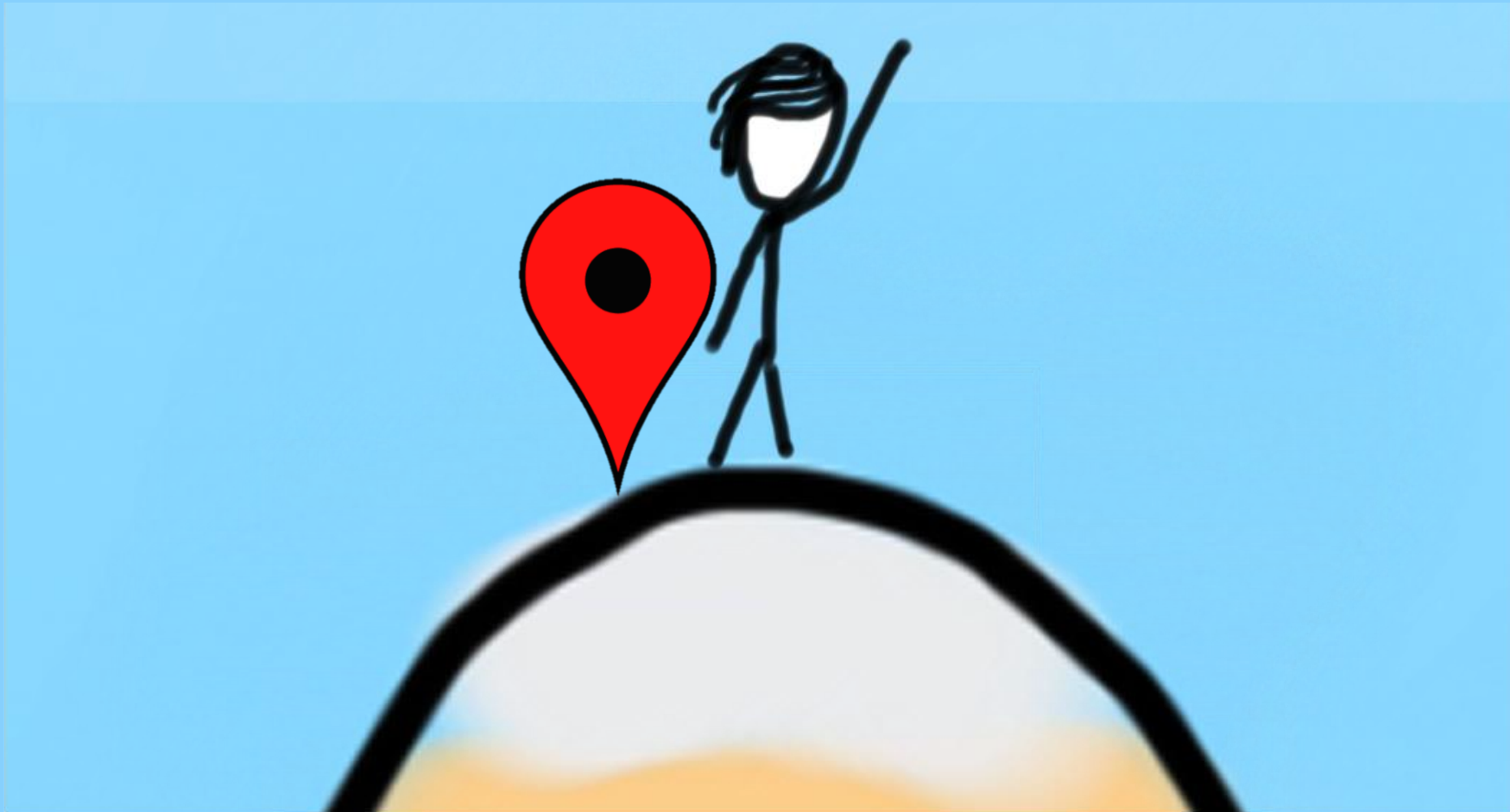
Information Theory

Chris Piech

CS109, Stanford University

Learning Goals

1. Calculate information gain
2. Make choices that maximize information gain
3. Numerically score how similar two distributions are



Uncertainty Theory

Beta
Distributions

Thompson
Sampling

Adding
Random Vars

Central Limit
Theorem

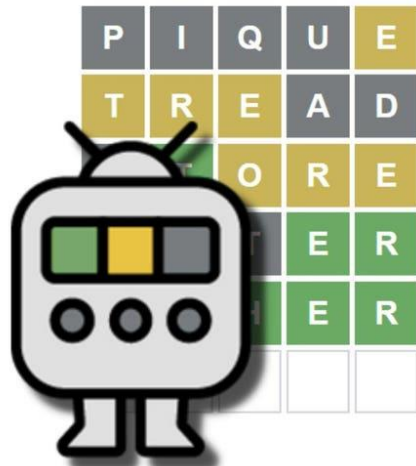
Sampling

Bootstrapping

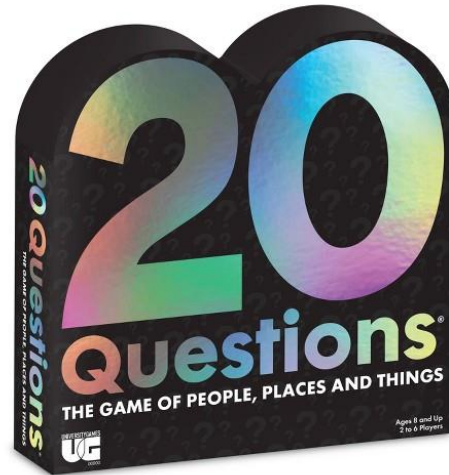
Algorithmic
Analysis

Information
Theory

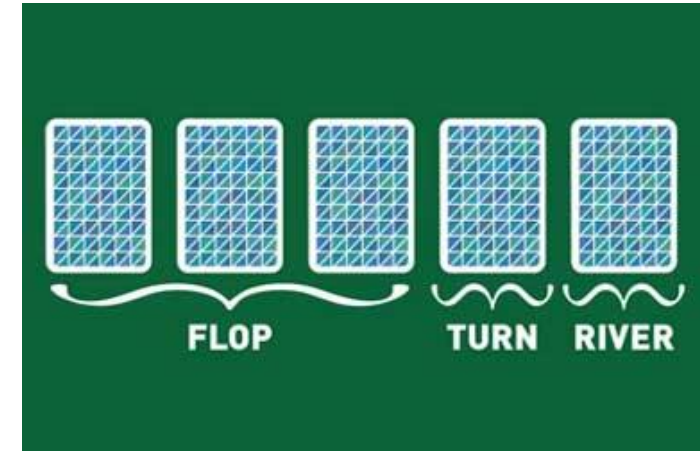
WorldeBot



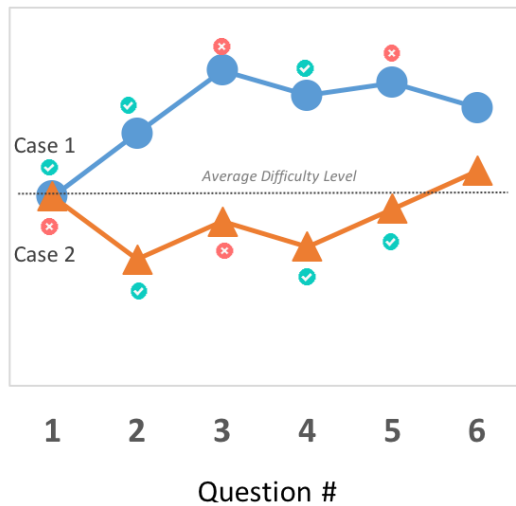
Decision Trees



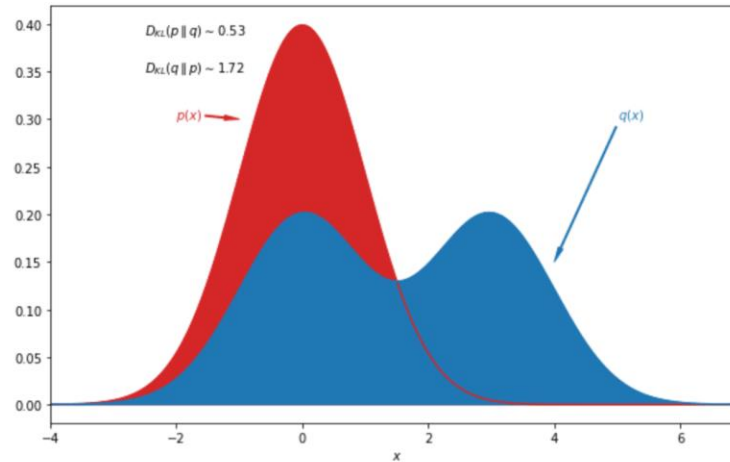
Value of Info in Poker



Adaptive Tests



Comparing Distributions



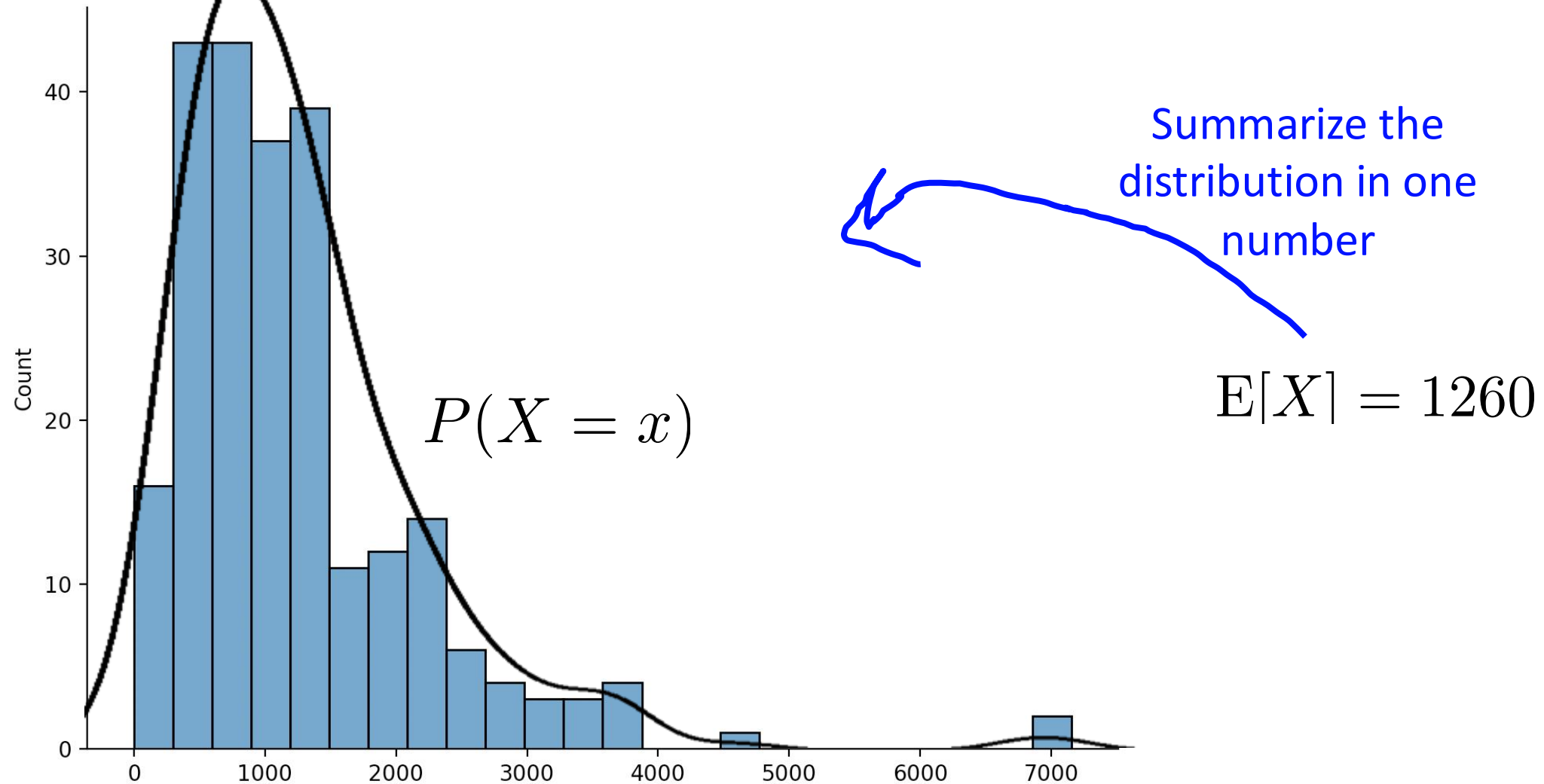
Compression of Data



Review

Limitation of Expectation

X = time to complete the medical diagnosis problem (in seconds)



Expectation

$$E[X] = \sum_x x \cdot P(X = x)$$

$$E[X] = \sum_i E[Y_i] \quad \text{if } X = \sum_i Y_i$$

$$E[X] = \sum_y E[X|Y = y]P(Y = y) \quad \text{For any } Y$$

Expectation of a Function

Law of unconscious statistician

$$\mathbf{E}[g(X)] = \sum_x g(x) \cdot P(X = x)$$

So for example...

$$\mathbf{E}[X^2] = \sum_x x^2 \cdot P(X = x)$$

End Review

Plot Line

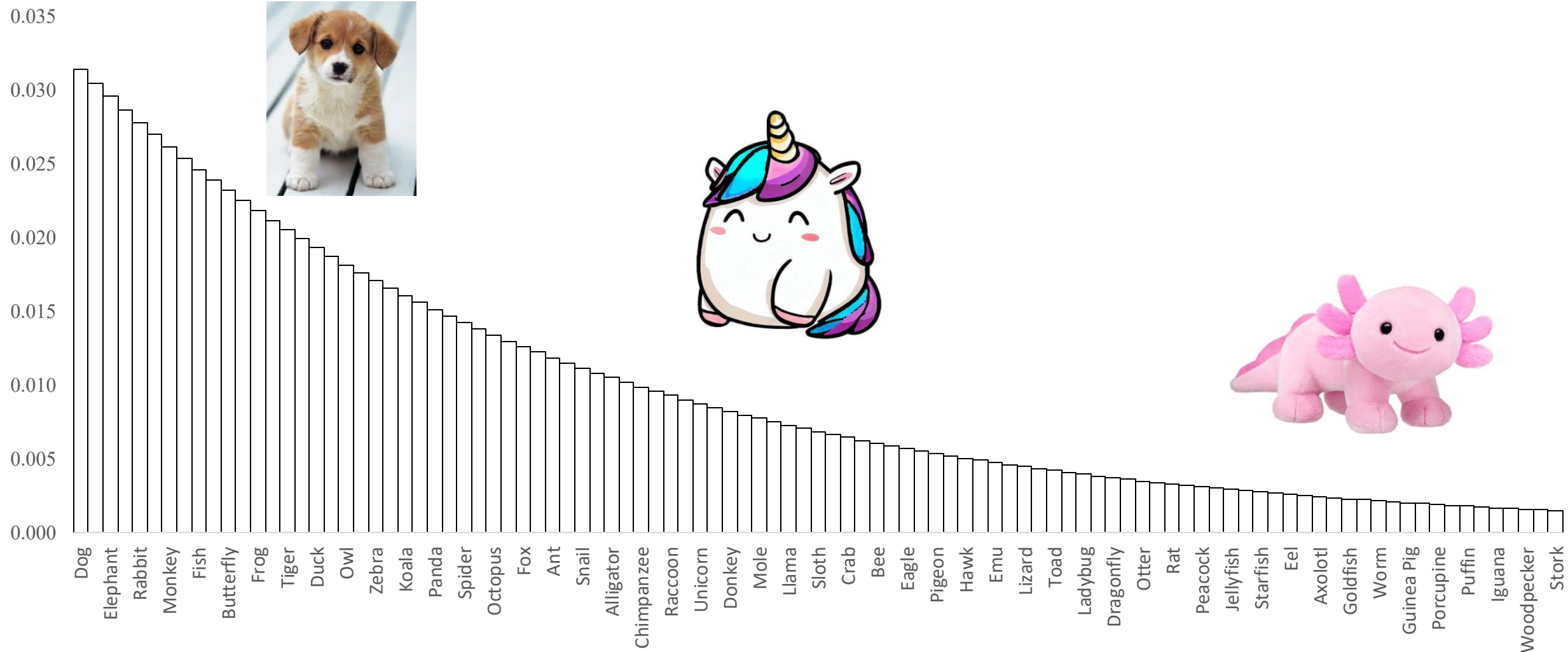


1. We want to **chose questions** in **think of an animal** (Let X be the animal random var)
2. Idea! Select the question which most **reduces our expected “uncertainty”** in X
3. We can **measure “uncertainty”** as “expected amount of surprise when we find out X ”
4. We can **measure “surprise”** in an assignment as $\log 1/P(X=x)$
5. This measure of **uncertainty of X** is **super helpful** for lots of problems!

I am thinking of an **animal**...



I am thinking of an animal



What is the **best question** to ask?

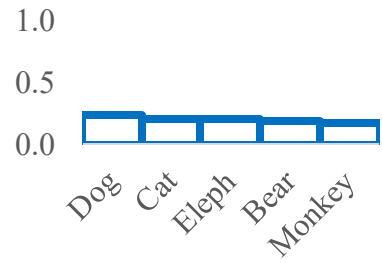
Question Choosing Algorithms

Algorithm	Average Questions	Standard Error of the Mean
Random Questions	61.34	1.41
Binary Search	7.79	0.02
Information Theory Search	5.33	0.04

Hey what's that?



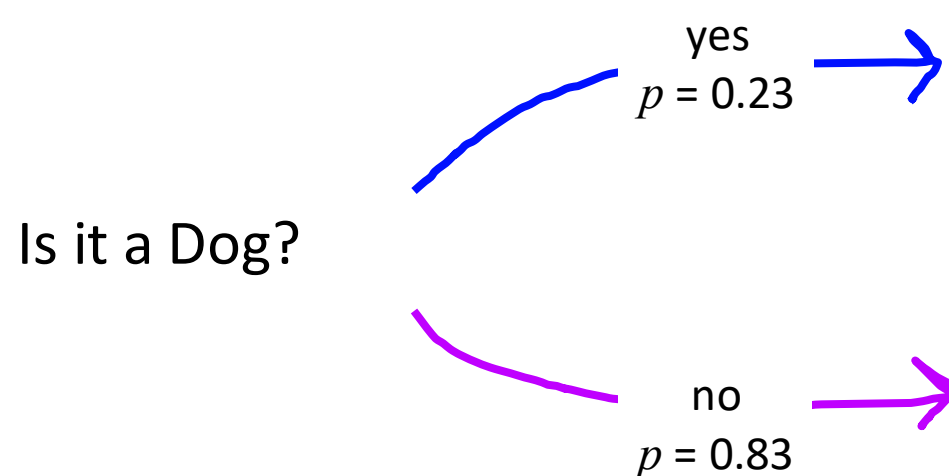
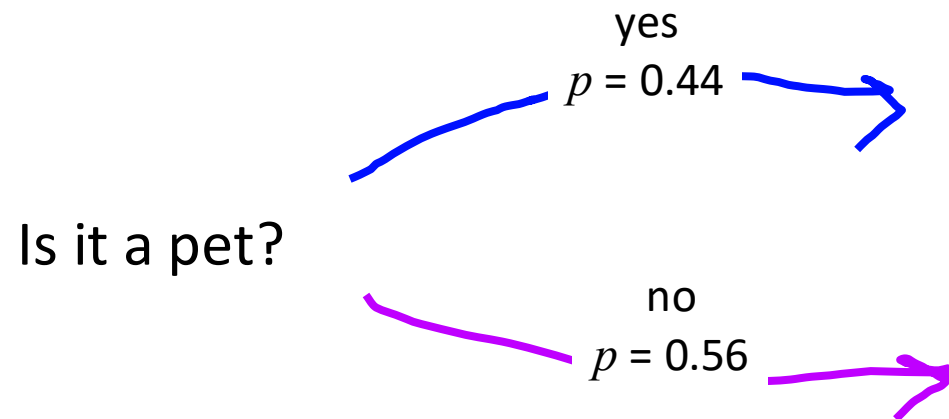
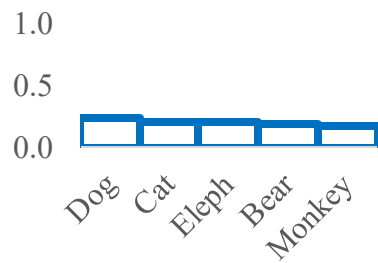
Which Question is Better?



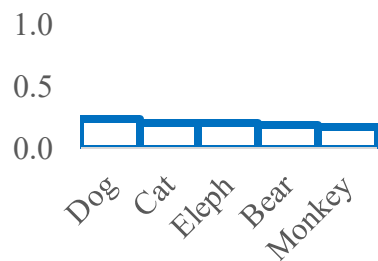
Is it a pet?

Is it a Dog?

Which Question is Better?



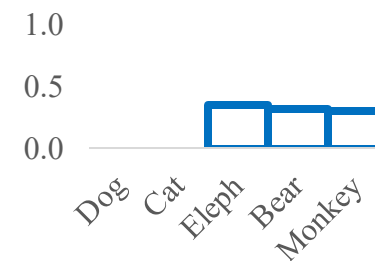
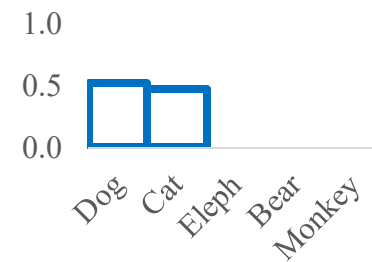
Which Question is Better?



Is it a pet?

yes
 $p = 0.44$

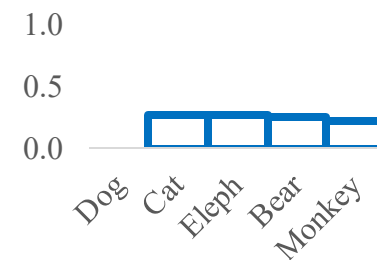
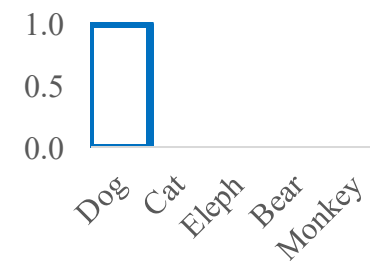
no
 $p = 0.56$



Is it a Dog?

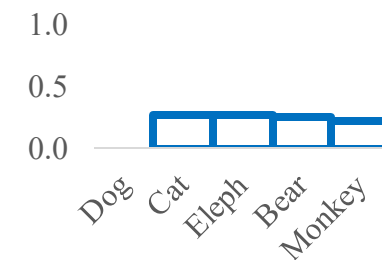
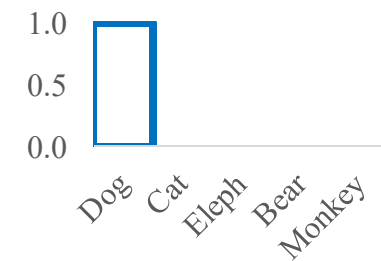
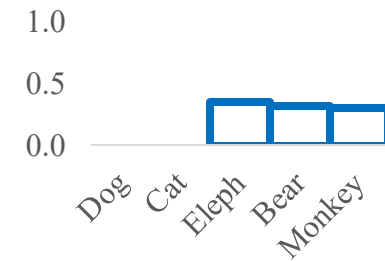
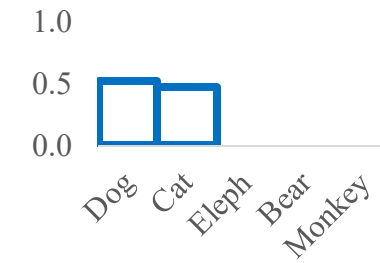
yes
 $p = 0.23$

no
 ~~$p = 0.83$~~
77



Which Question is Better

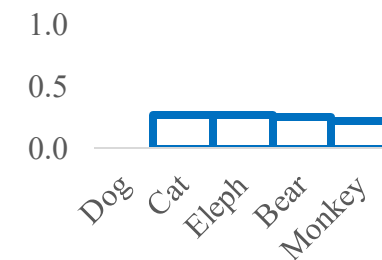
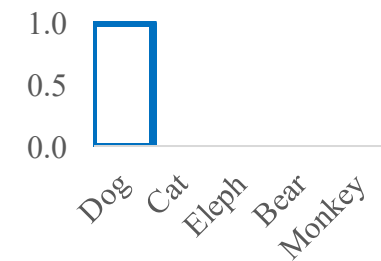
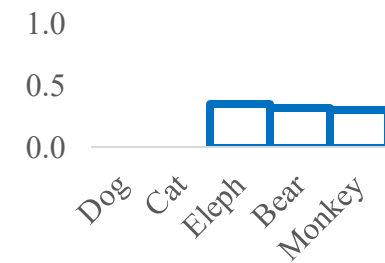
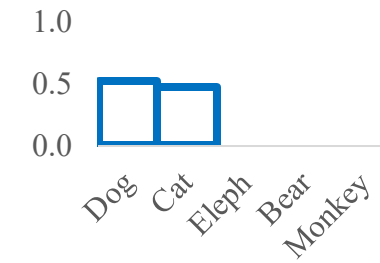
Can you “score” how **bad** each of these PMFs are?



Which Question is Better

uncertain

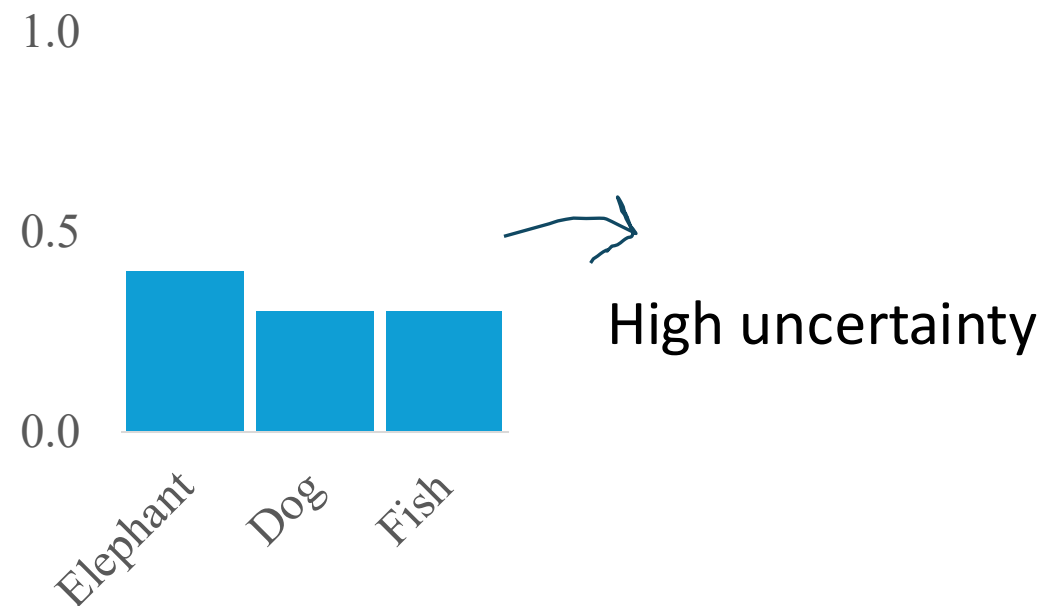
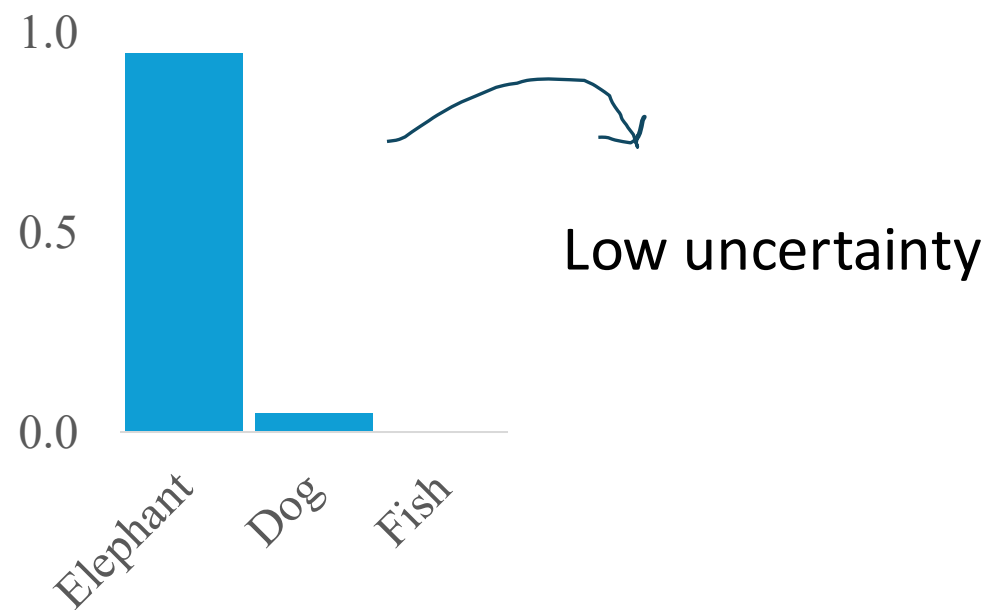
Can you “score” how ~~bad~~
each of these PMFs are?



What we really need is a **measure** of our **uncertainty** in a random variable

Uncertainty of a Random Variable

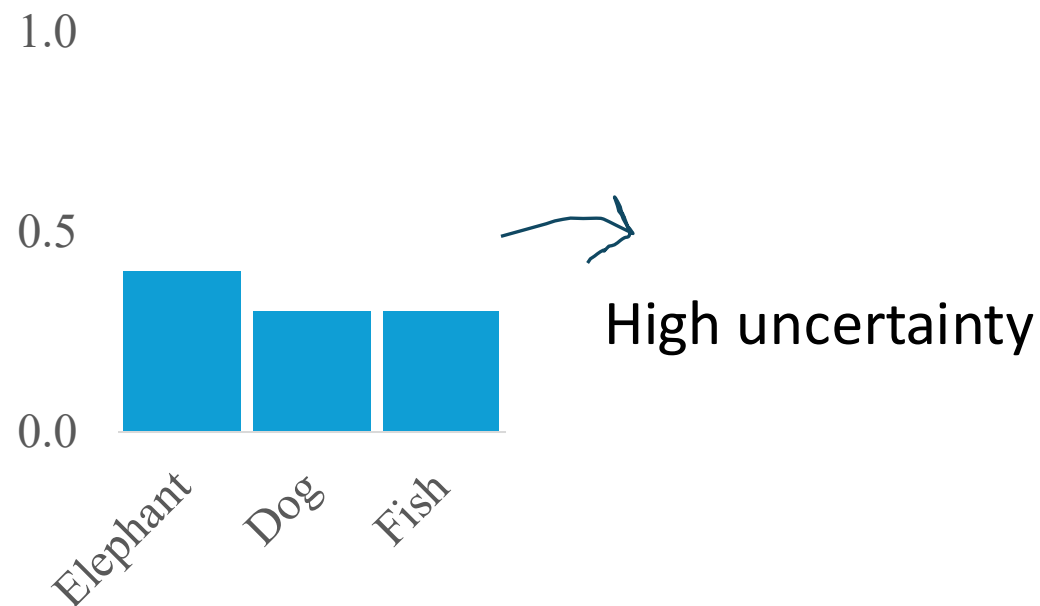
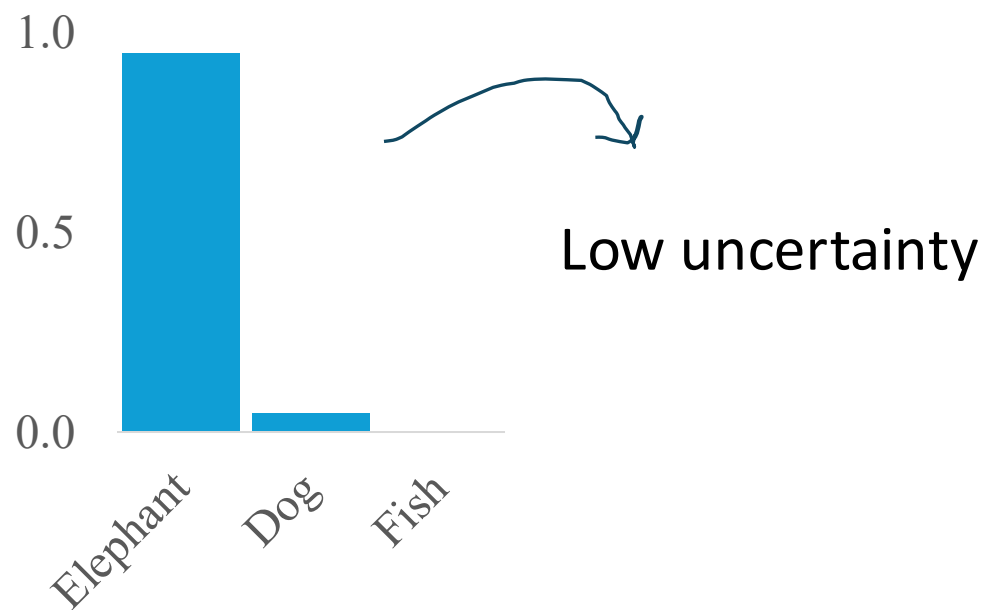
Let X be any random variable. We can calculate a statistic, “**Uncertainty**” to express how much we don’t know about X



Uncertainty of a Random Variable

Let X be any random variable. We can calculate a statistic, “**Uncertainty**” to express how much we don’t know about X

$\text{Uncertainty}(X) =$ Expected “Surprise” when I observe X

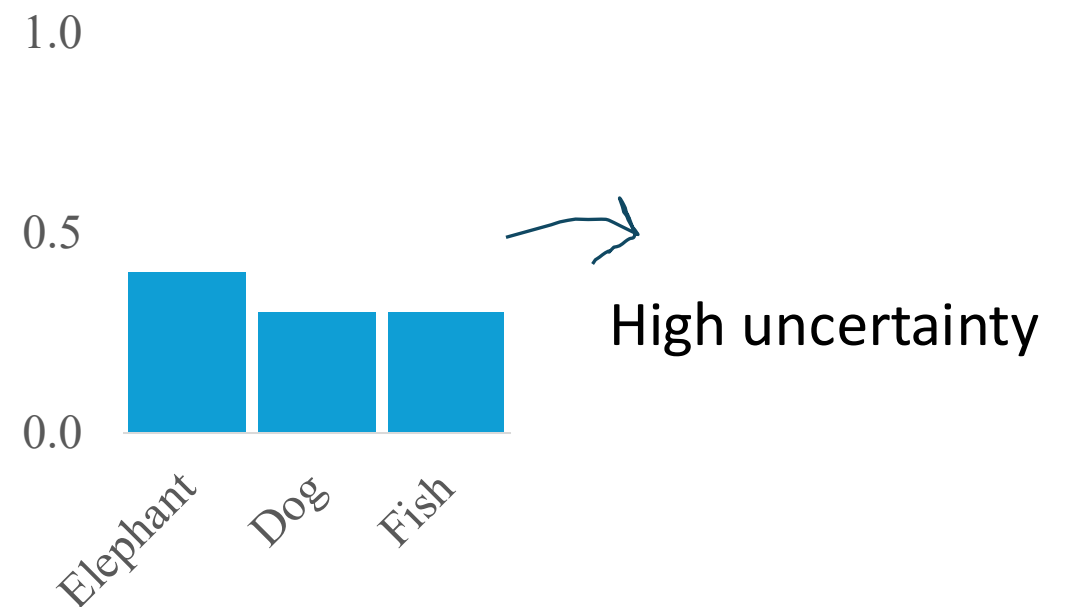
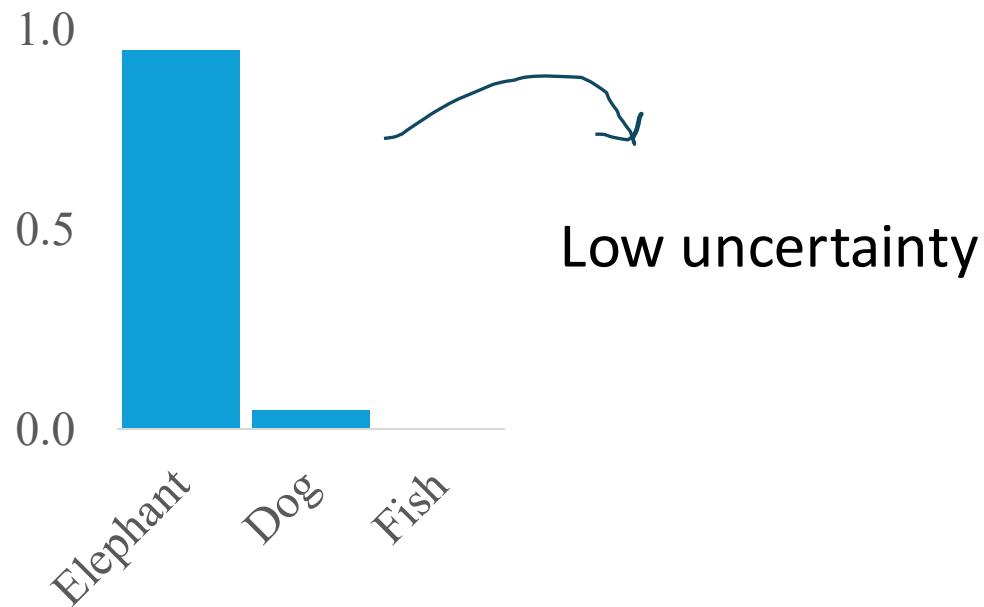


Uncertainty of a Random Variable

Let X be any random variable. We can calculate a statistic, “**Uncertainty**” to express how much we don’t know about X

$$\text{Uncertainty}(X) = \sum_{x \in X} \text{Surprise}(X = x) \cdot P(X = x)$$

Uncertainty is
expected Surprise



Ok, but then what is our measure of
Surprise?

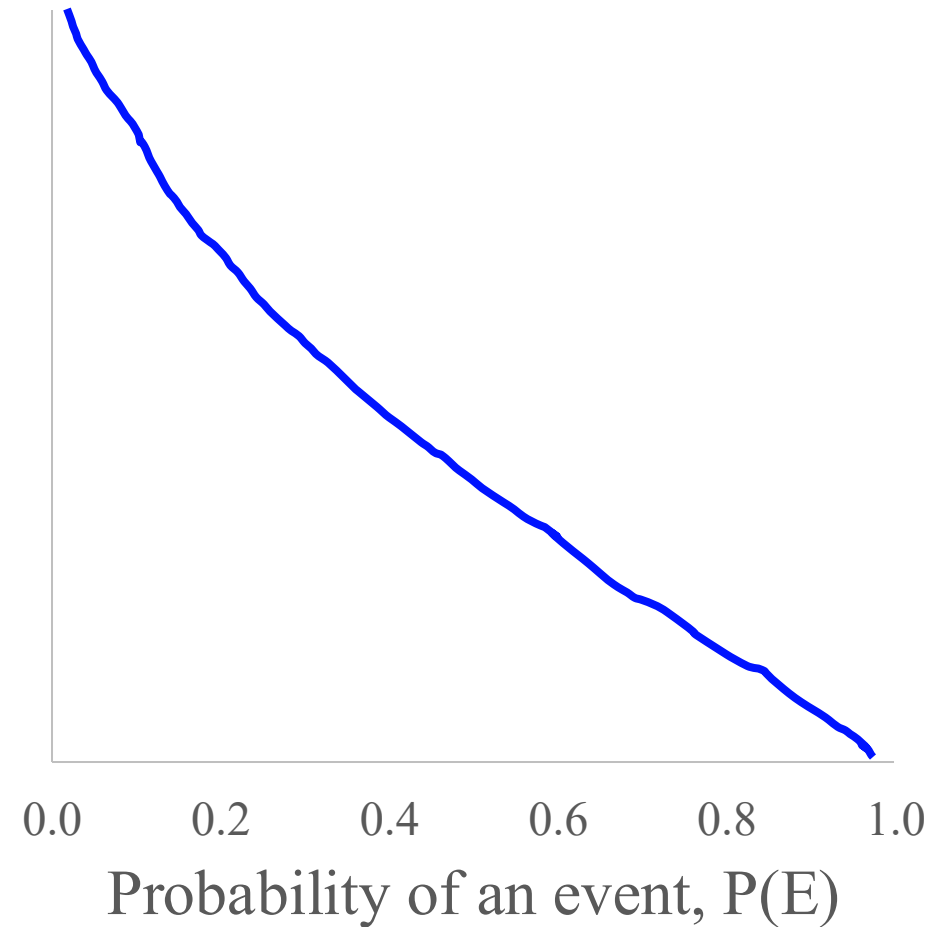
Surprise of an Event

High probability events are **not surprising**

Low probability events are **surprising**

Relationship should be monotonic

How surprising is the event?



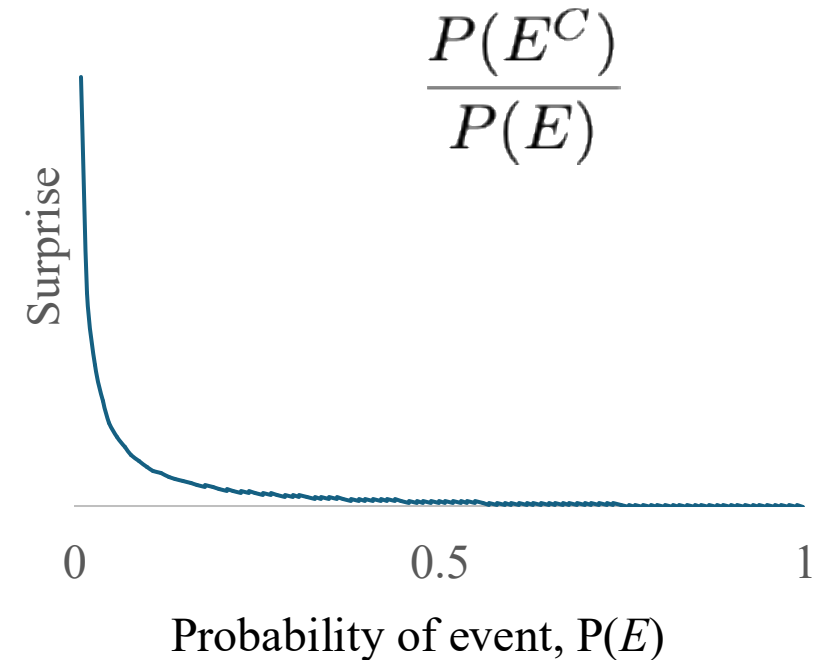
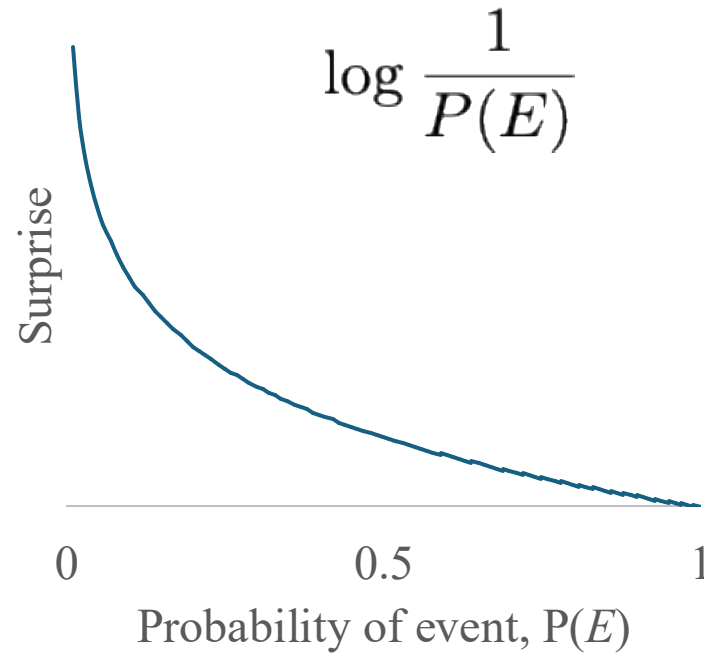
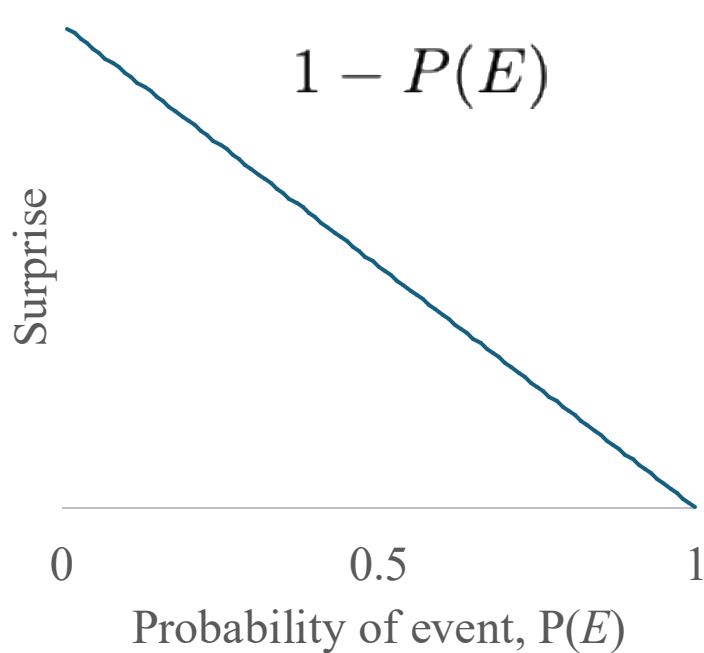
Surprise of an Event

High probability events are **not surprising**

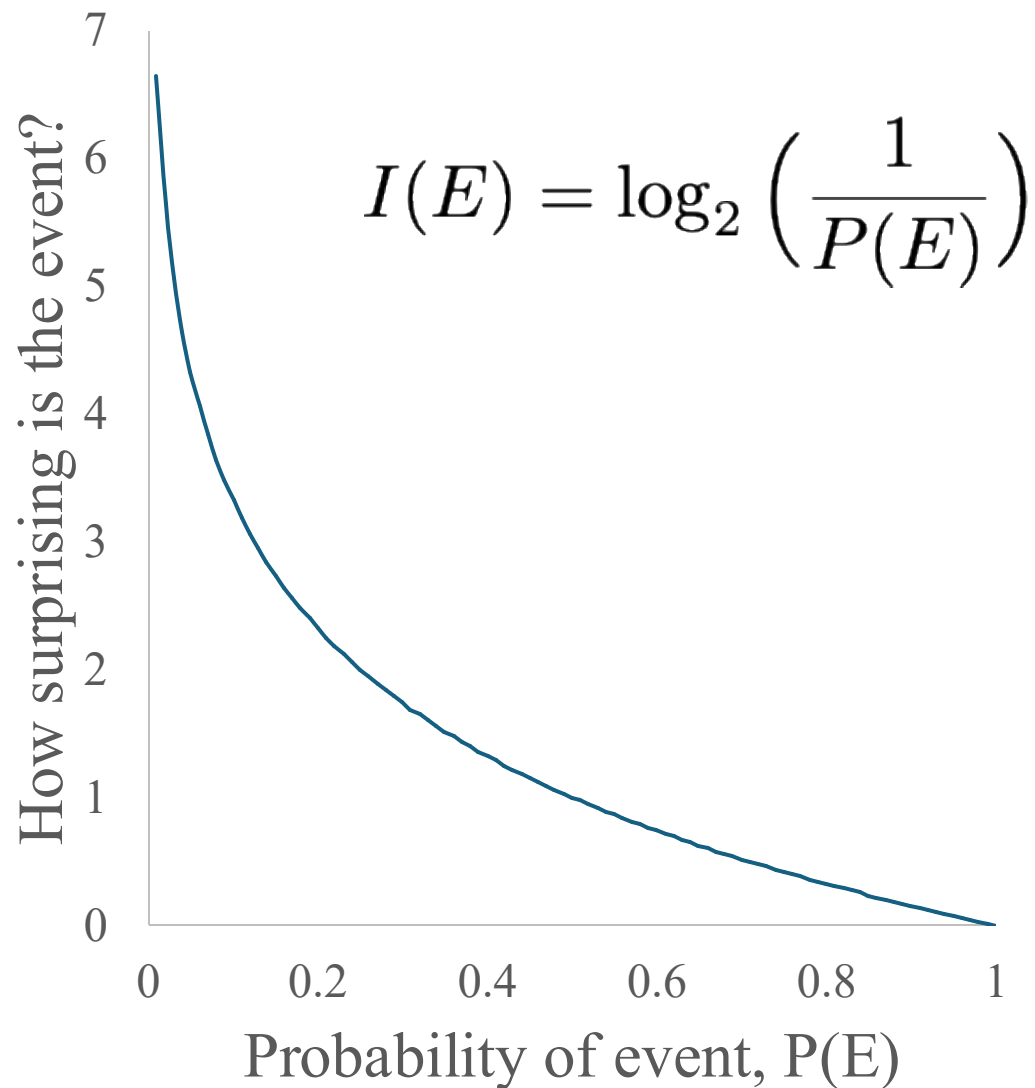
Low probability events are **surprising**

Relationship should be monotonic

Here are three reasonable options



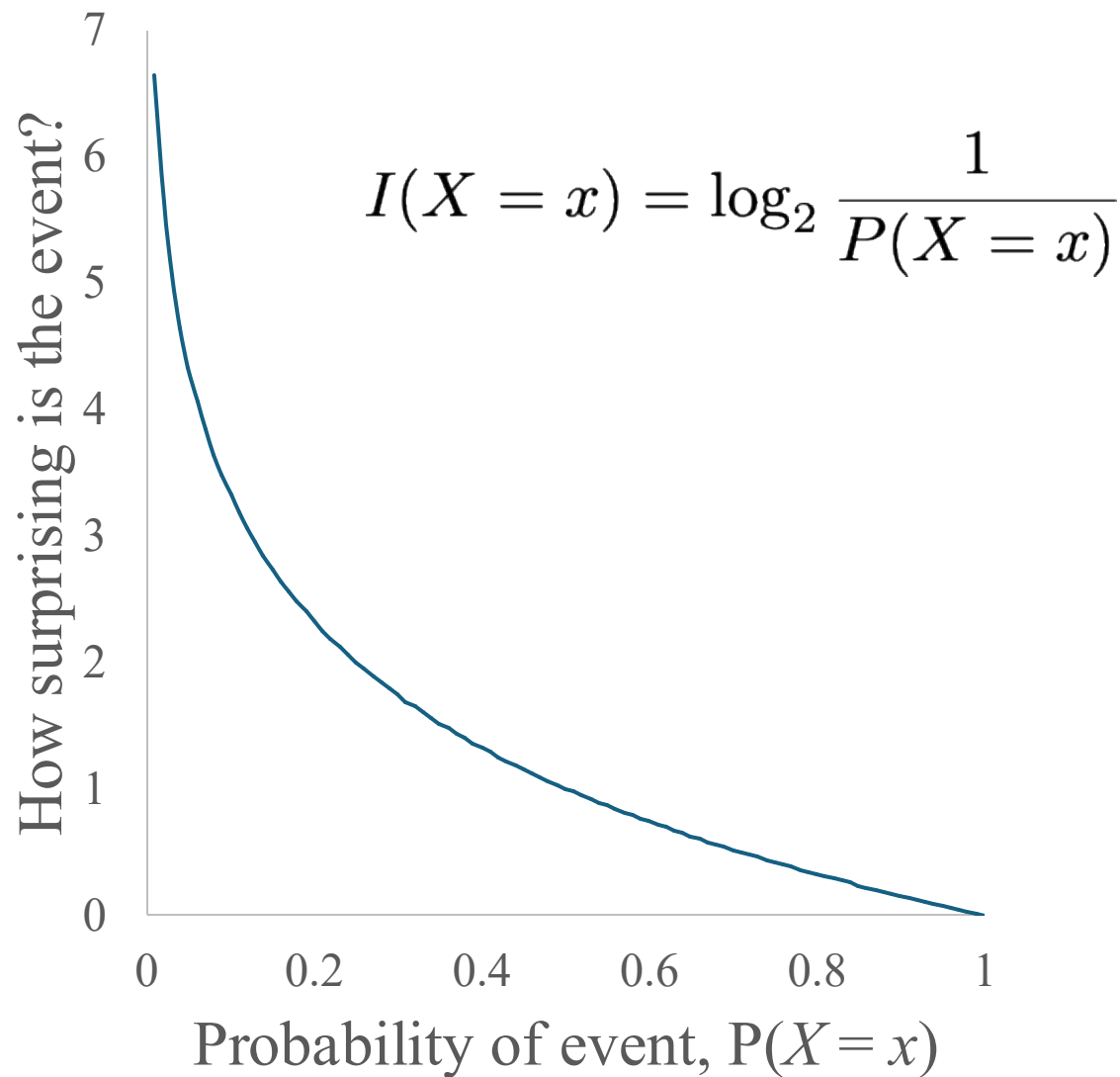
Surprise of an Event, $I(E)$



Probability of Event $P(E)$	Surprise of Event $I(E)$
1	0
1/2	1
1/4	2
1/8	3
1/16	4
1/32	5
1/64	6

$I(E)$ stands for “Information Content” aka “Surprisal” aka “Self-Information”

Surprise of an Event, $I(X = x)$



Probability of Event $P(X = x)$	Surprise of Event $I(X = x)$
1	0
$\frac{1}{2}$	1
$\frac{1}{4}$	2
$\frac{1}{8}$	3
$\frac{1}{16}$	4
$\frac{1}{32}$	5
$\frac{1}{64}$	6

$I(X = x)$ stands for
“Information Content” aka
“Surprisal” aka
“Self-Information”

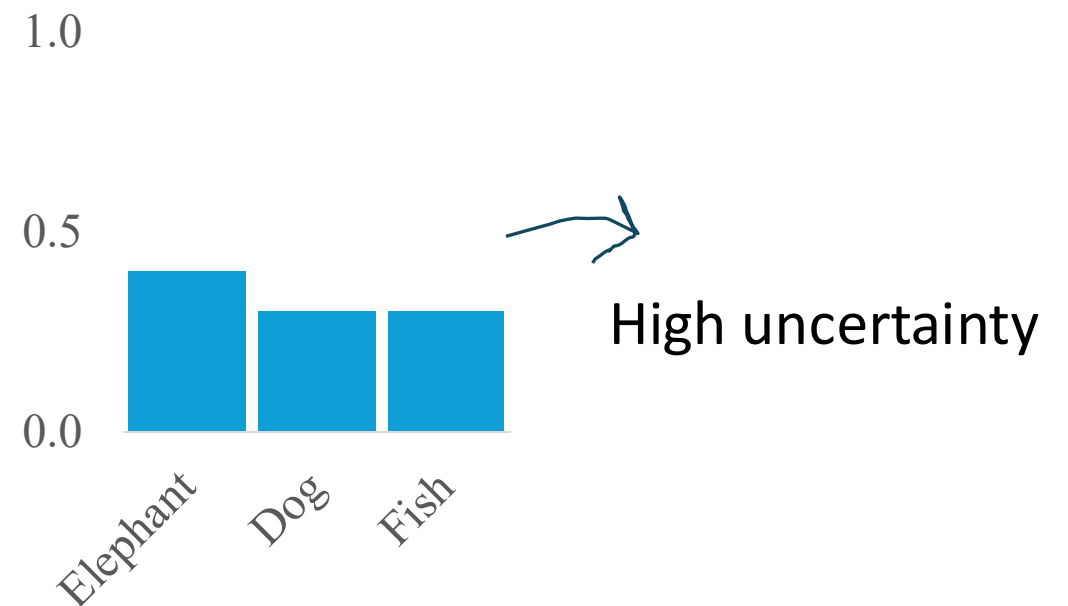
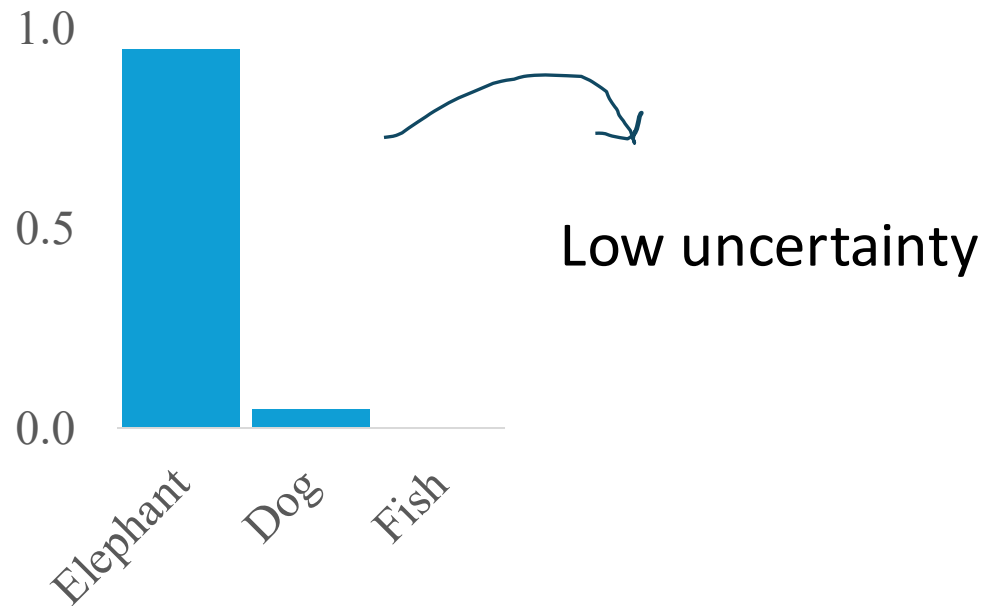
Back to the **measure** of **uncertainty** in
the outcome of a random variable

Uncertainty of a Random Variable

Let X be any random variable. We can calculate a statistic, “**Uncertainty**” to express how much we don’t know about X

$$\text{Uncertainty}(X) = \sum_{x \in X} \text{Surprise}(X = x) \cdot P(X = x)$$

Uncertainty is expected Surprise



Uncertainty of a Random Variable

Let X be any random variable. We can calculate a statistic, “**Uncertainty**” to express how much we don’t know about X

$$\text{Uncertainty}(X) = \sum_{x \in X} \text{Surprise}(X = x) \cdot P(X = x)$$

Uncertainty is expected Surprise

Uncertainty of a Random Variable

Let X be any random variable. We can calculate a statistic, “**Uncertainty**” to express how much we don’t know about X

$$\text{Uncertainty}(X) = \sum_{x \in X} \text{Surprise}(X = x) \cdot P(X = x)$$

$$= \sum_{x \in X} \log_2 \frac{1}{P(X = x)} \cdot P(X = x)$$

$$= \sum_{x \in X} \log_2 P(X = x)^{-1} \cdot P(X = x)$$

$$= \sum_{x \in X} -\log_2 P(X = x) \cdot P(X = x)$$

$$= - \sum_{x \in X} \log_2 P(X = x) \cdot P(X = x)$$

Uncertainty is expected Surprise

Our favorite measure of Surprise

$1/x$ is the same as x^{-1}

Log of a power (here -1)

Pull the negative out

Uncertainty of a Random Variable (Entropy)

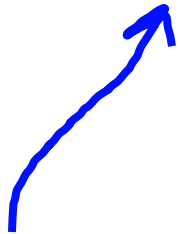
Let X be any random variable. We can calculate a statistic, “**Uncertainty**” to express how much we don’t know about X

Calculates expected surprise

$$\text{Uncertainty}(X) = \sum_{x \in X} \log_2 \frac{1}{P(X = x)} \cdot P(X = x)$$

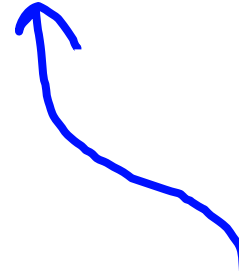
H

~~Uncertainty~~



My preferred name for “entropy” aka $H(X)$

$$\text{Surprise}(X = x) = \log_2 \frac{1}{P(X = x)}$$



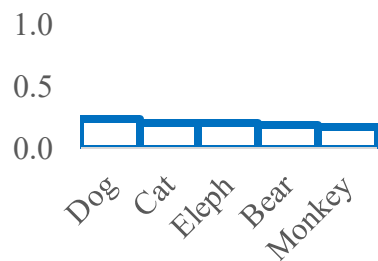
By the way...

$$H(X)$$

Is called **Shannon Entropy** to scare students and impress Physics people

Back to
I am thinking of an **animal**...

Which Question is Better?



$$H(X) = 2.3$$

Is it a pet?

yes

$$p = 0.44$$

no

$$p = 0.56$$

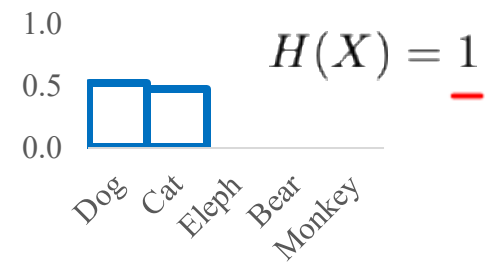
Is it a Dog?

yes

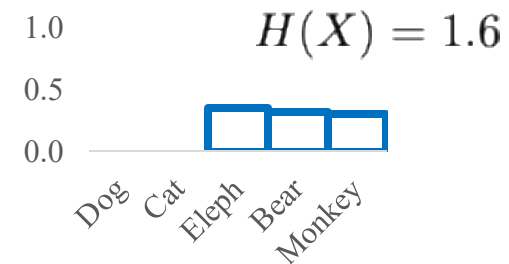
$$p = 0.23$$

no

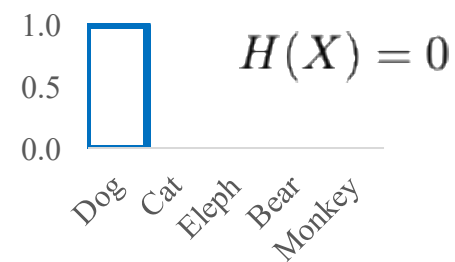
$$p = 0.83$$



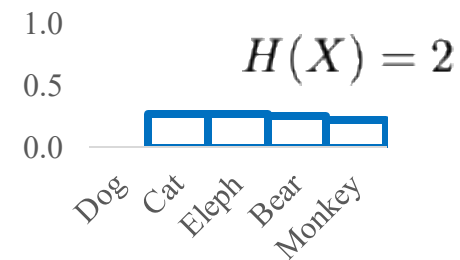
$$H(X) = 1$$



$$H(X) = 1.6$$

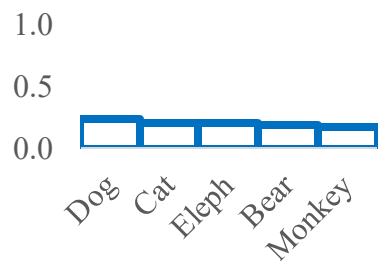


$$H(X) = 0$$



$$H(X) = 2$$

Which Question is Better?



$$H(X) = 2.3$$

Is it a pet?

$$E[H(X)] = \underline{1.3}$$

yes

$$p = 0.44$$

no

$$p = 0.56$$

yes

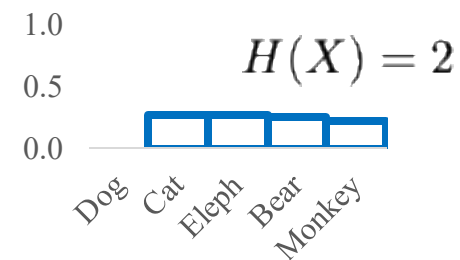
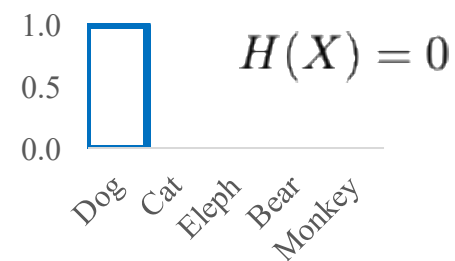
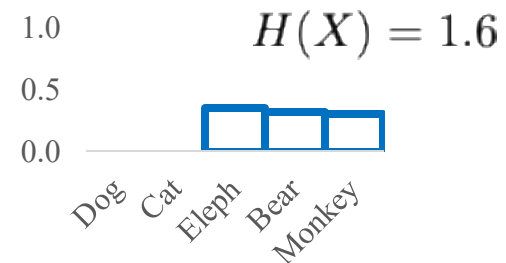
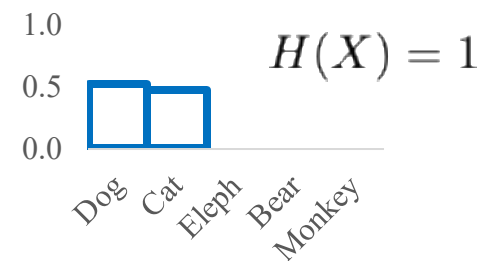
$$p = 0.23$$

Is it a Dog?

$$E[H(X)] = \underline{1.7}$$

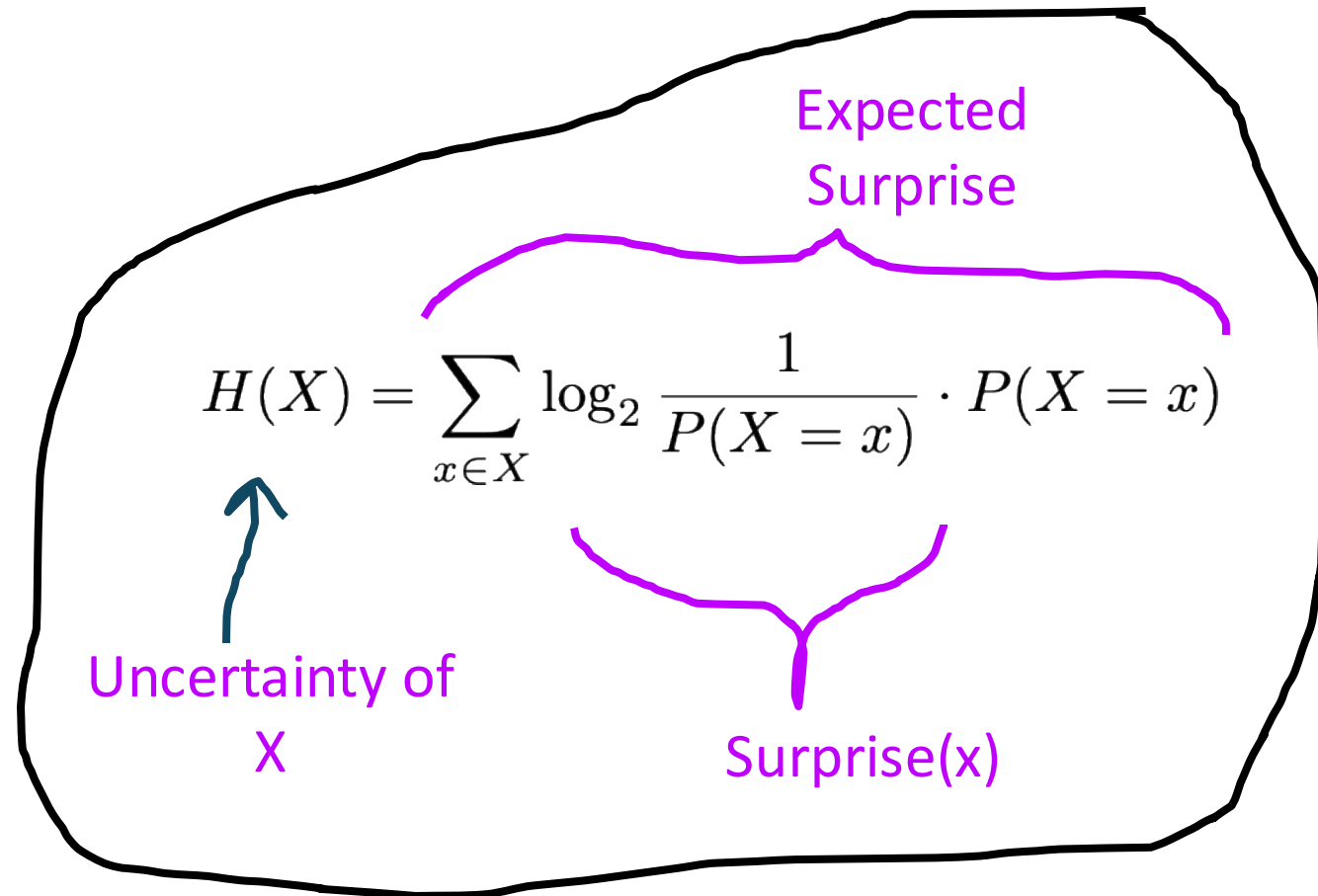
no

$$p = 0.83$$



Uncertainty (aka Entropy) in code

```
def calc_uncertainty(pmf):  
    # this calculates the entropy of the distribution  
    # aka the uncertainty  
    uncertainty = 0  
    for x in pmf: ←  
        p_x = pmf[x]  
        # skip zero probabilities  
        if p_x == 0: continue  
        surprise_x = np.log2(1/p_x)  
        uncertainty += surprise_x * p_x  
    return uncertainty
```



The diagram shows the entropy formula $H(X) = \sum_{x \in X} \log_2 \frac{1}{P(X=x)} \cdot P(X=x)$ enclosed in a hand-drawn black oval. A blue arrow points from the text 'Uncertainty of X' to the $H(X)$ term. A purple bracket above the formula spans the entire sum, labeled 'Expected Surprise'. A purple bracket below the formula spans the $\log_2 \frac{1}{P(X=x)}$ term, labeled 'Surprise(x)'.

$$H(X) = \sum_{x \in X} \log_2 \frac{1}{P(X=x)} \cdot P(X=x)$$

Expected Surprise

Uncertainty of X

Surprise(x)

To the code!

Limitations of Expected Entropy for Decision Making?

Entropy in the sum of two dice.

What is more informative:

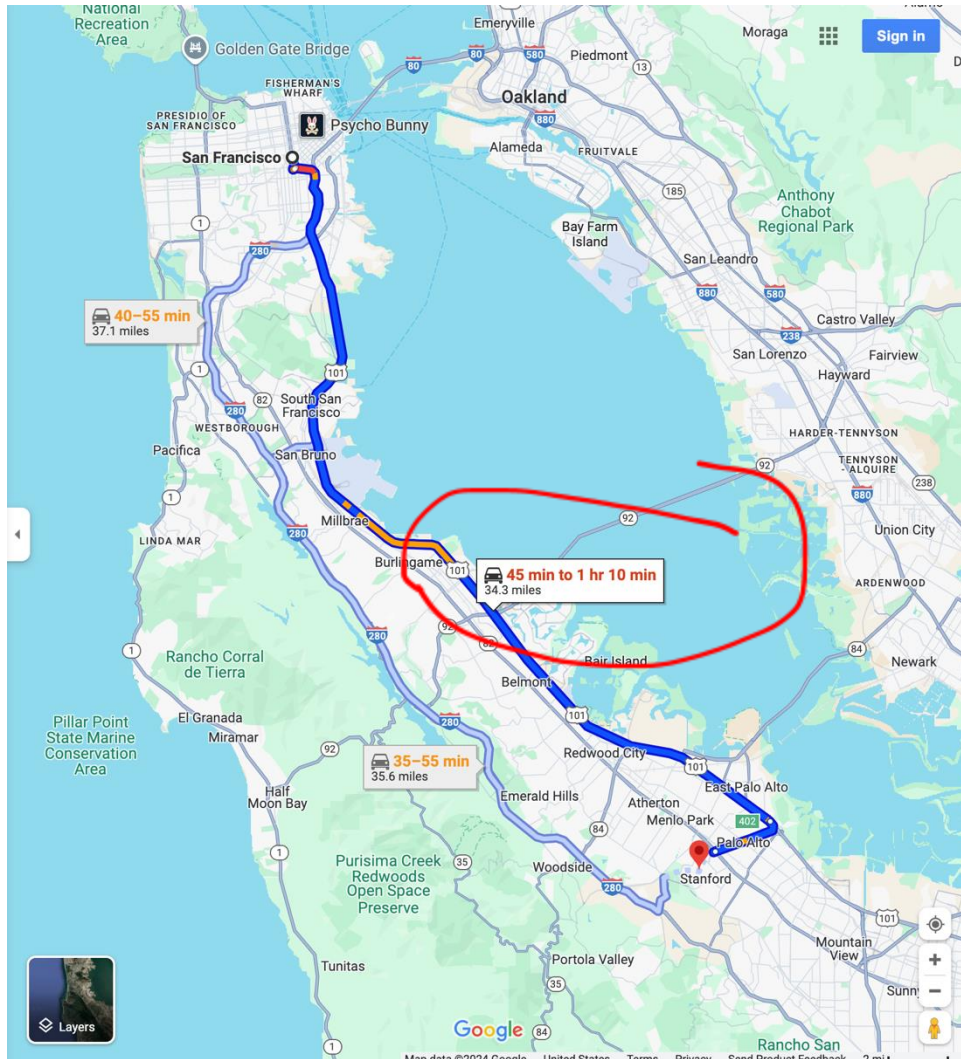
I tell you that the **sum is odd**

I tell you the value of the **first dice is 1**

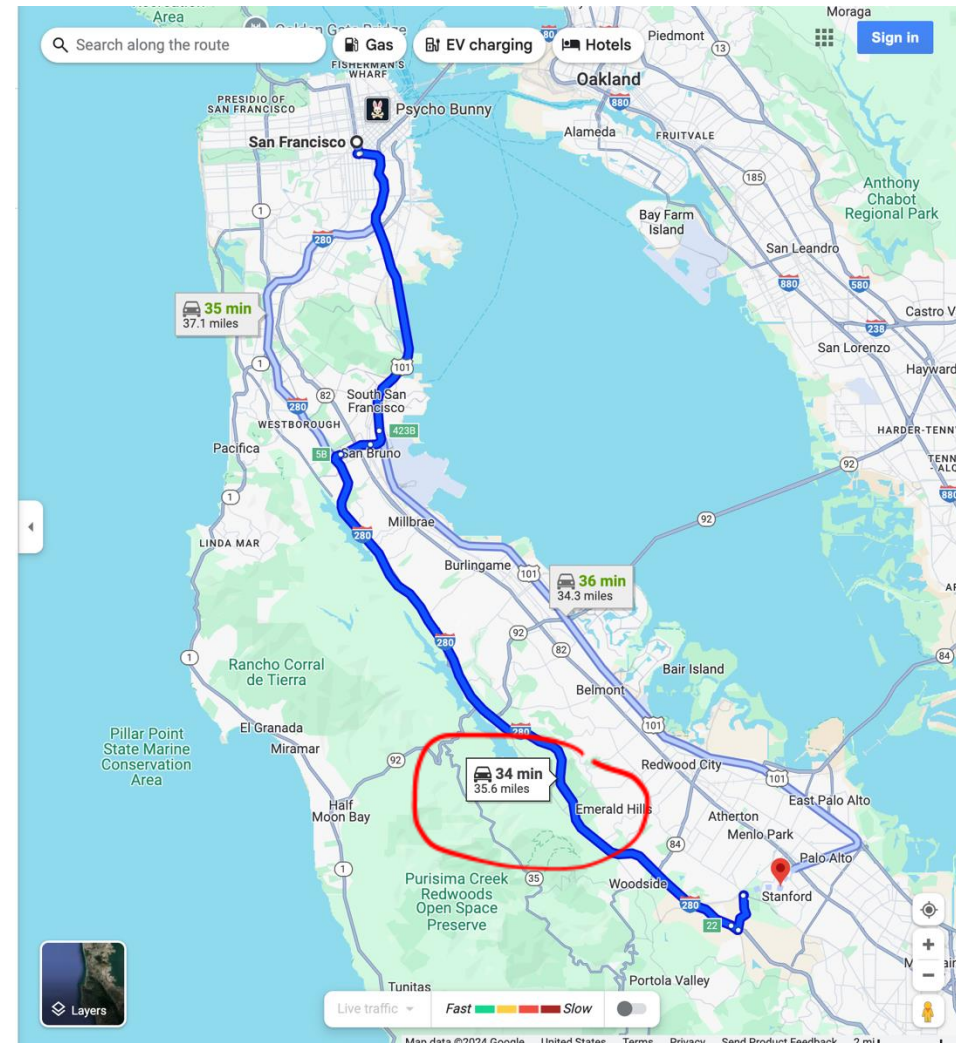
Traffic Predictions Over Time

Should You Wait to Plan Your Route?

Prediction 1 day in advance

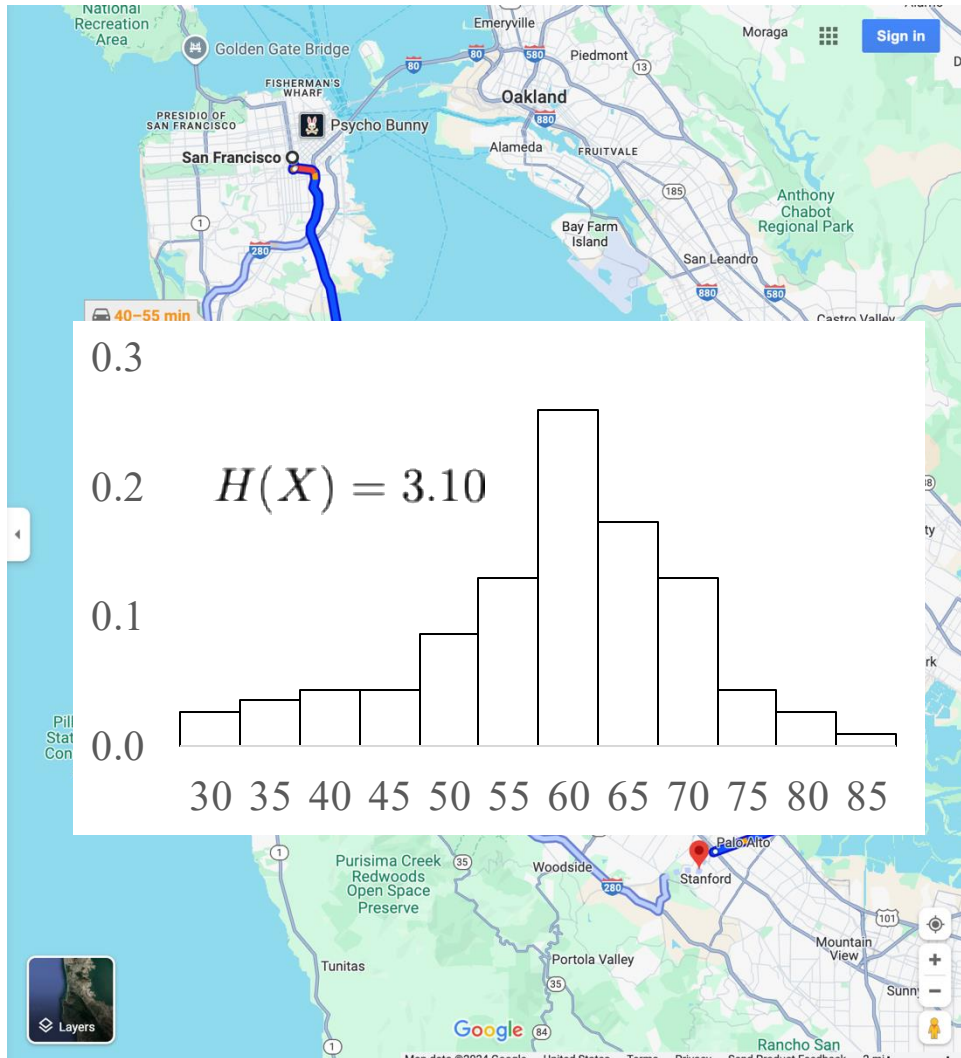


Prediction 30 mins in advance

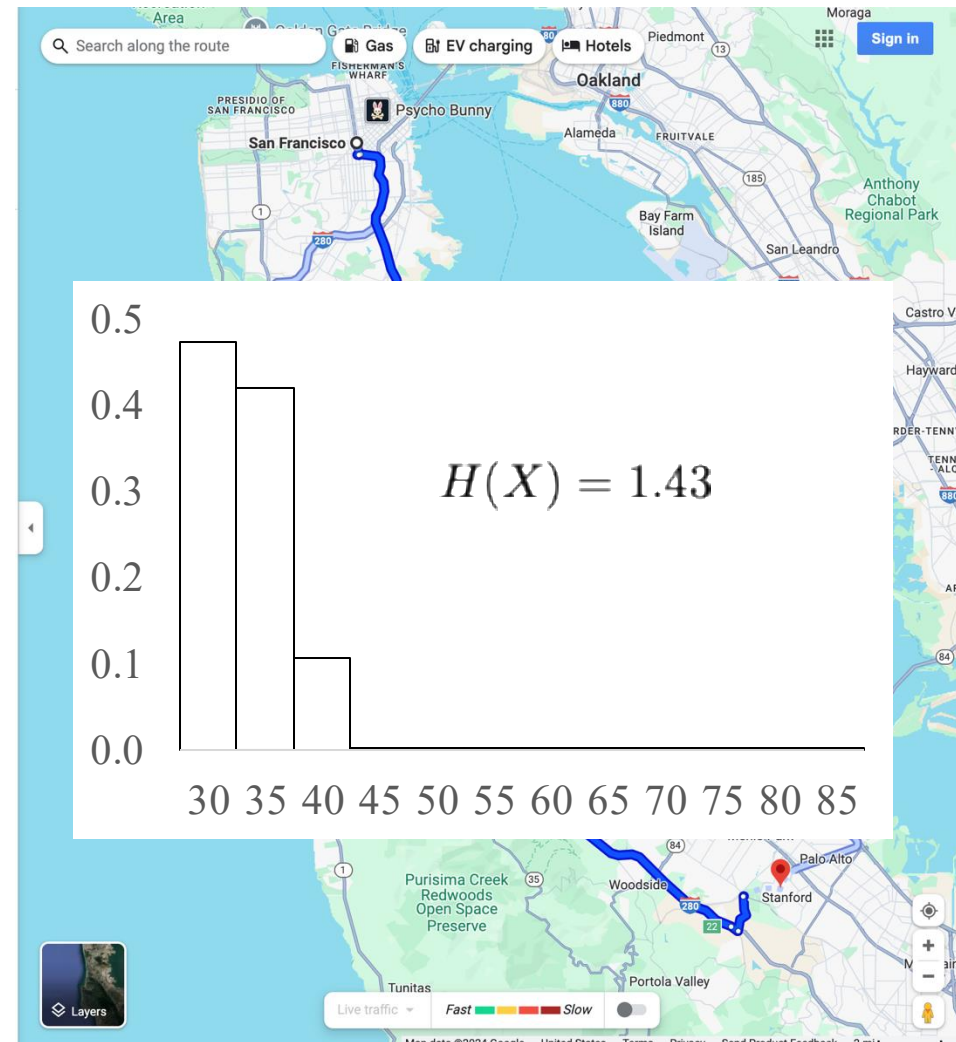


Should You Wait to Plan Your Route?

Prediction 1 day in advance

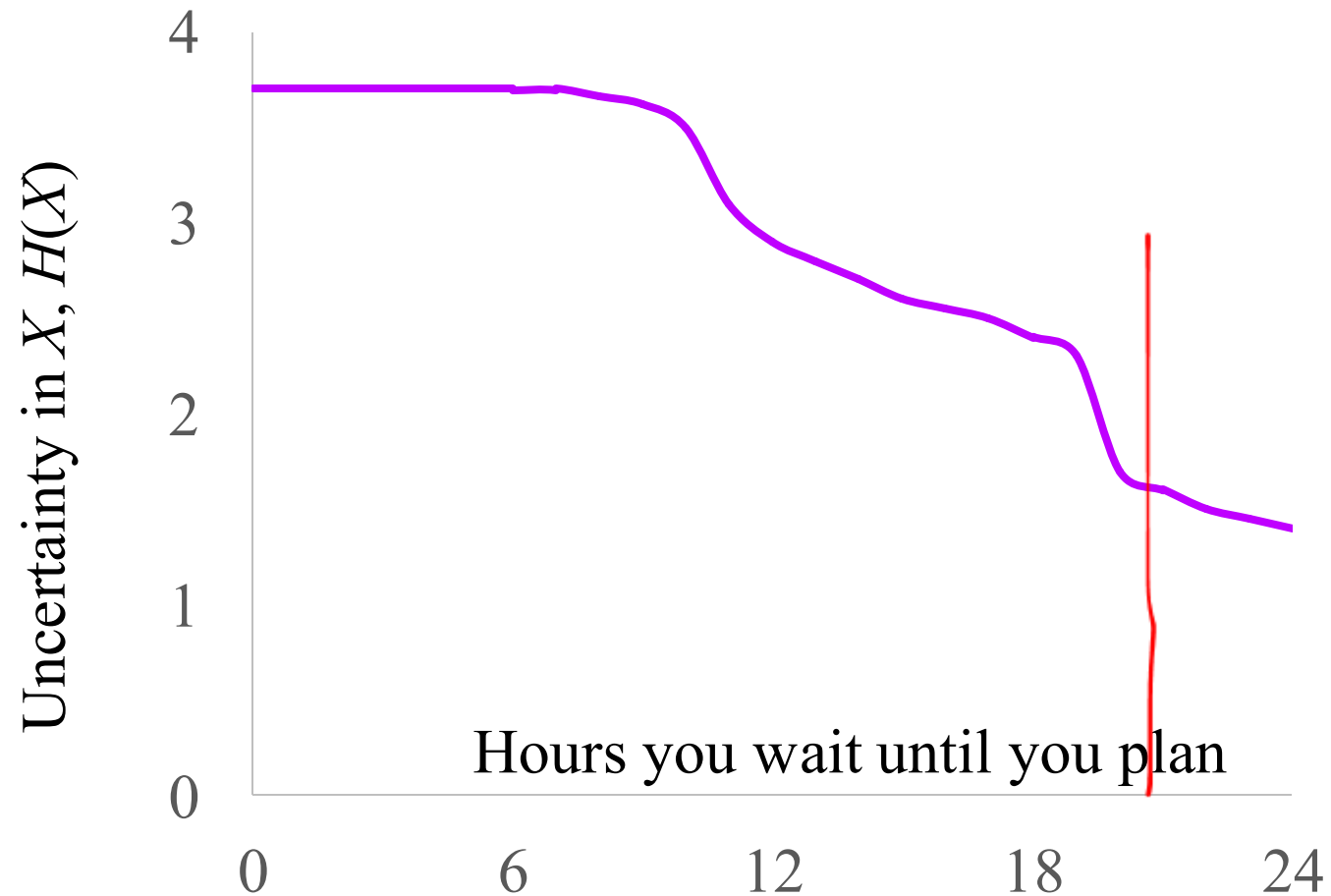


Prediction 30 mins in advance



Should You Wait to Plan Your Route?

Let X be the amount of time to drive from SF to Stanford



If time:

Indus Valley Script. A language?

The Indus civilization - one of the world's earliest urban societies - emerged 5,300 years ago. The script they used is one of the last remaining undeciphered alphabets.



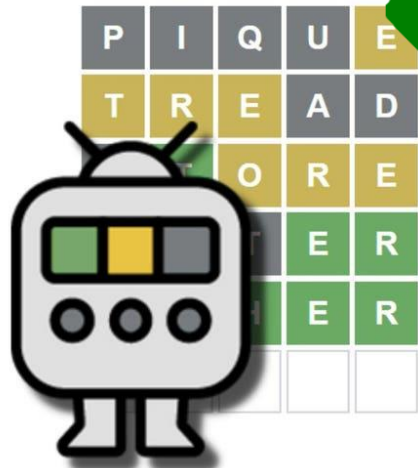
Let X be the first char
Let Y be the second char

- a) What is $P(Y = y \mid X = x)$?
- b) $E[H(Y \mid X = x)]$

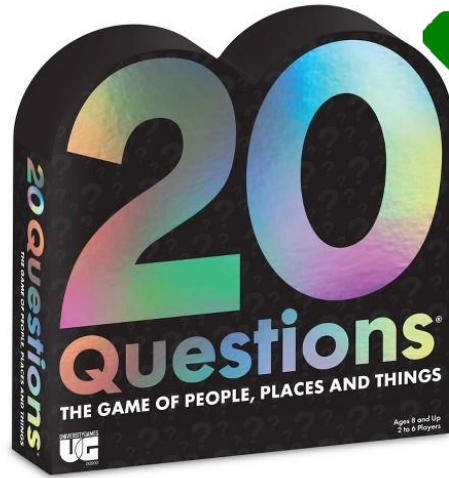
We are using entropy analysis to show that the script represents a true linguistic system. Let **all_examples** be the list of all recorded examples of the script, where each item in the list is one example string:

```
all_examples = ['𑀩 𑀲 𑀲 𑀲',  
               '𑀲 𑀩 𑀲 𑀲 𑀲',  
               '𑀩 𑀲',  
               ...]
```

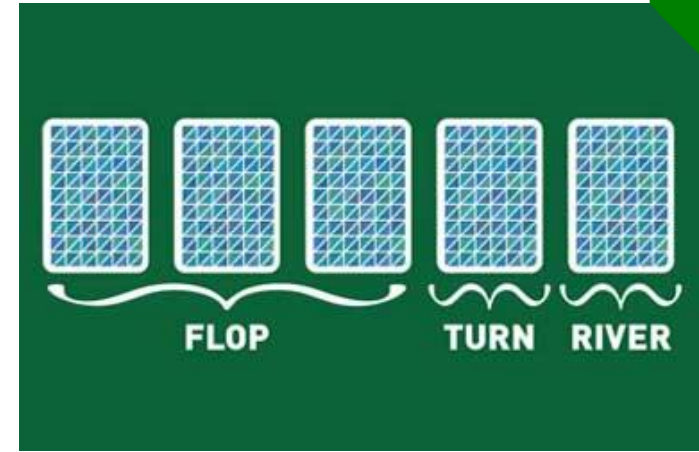
WorldeBot



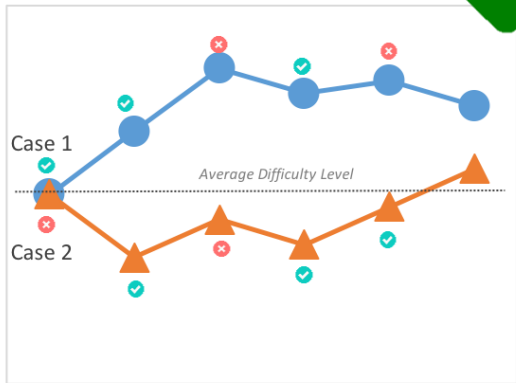
Decision Trees



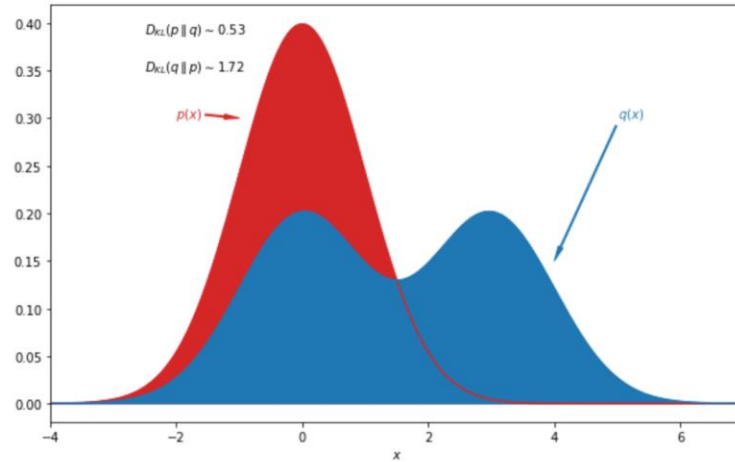
Value of Info in Poker



Adaptive Tests



Comparing Distributions



Compression of Data





DALL·E 3

Distance Between Two Distributions?

Recall this

