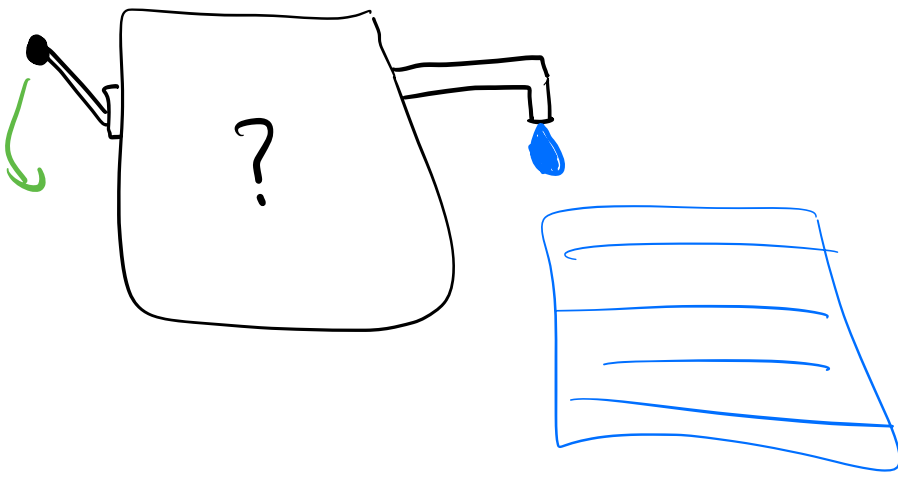


Probability and Statistics : A data scientist's perspective on the basics

Where does data come from?



Probability: how to model the ?

Statistics: Given data only, what can we learn about ?

Probability: basic terms

Coin flip

experiment - process that results in one of a set of outcomes

$\{H, T\}$

sample space - the set of possible outcomes

lands T

event - subset of sample space

the probability of an outcome s is written $P(s)$

$$P(H) = 0.5$$

$$P(T) = 0.5$$

and satisfies:

$$- 0 < P(s) < 1 \quad (\text{between } 0 \text{ and } 1)$$

$$- \sum_{s \in S} P(s) = 1 \quad (\text{they all sum to } 1)$$

$$V(T) = 0$$

$$V(H) = 1$$

A random variable is a function that maps
an outcome to a number.

The expected value of a R.V. V is the sum of V 's
values, each weighted by their probability:

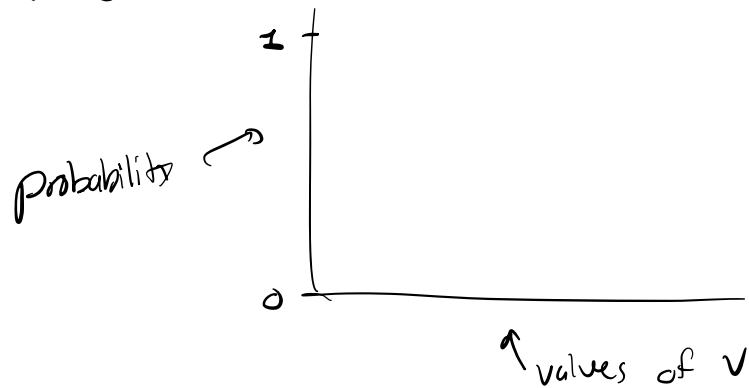
$$E(V) = 0.5 \cdot 0 + 0.5 \cdot 1 \\ = 0.5$$

$$E(V) = \sum_{s \in S} P(s) V(s)$$

Exercises 1-3

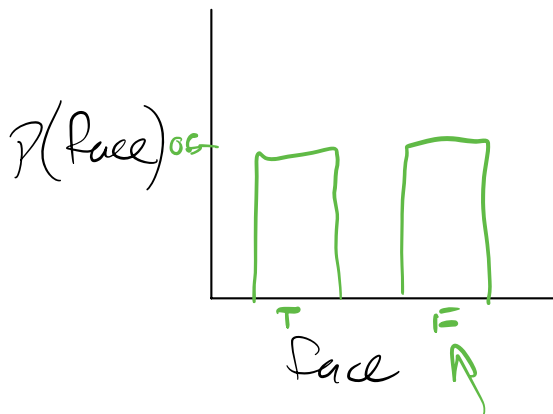
Probability Distributions:

A random variable V 's **probability density function (PDF)** is a function mapping each value of V to its probability.

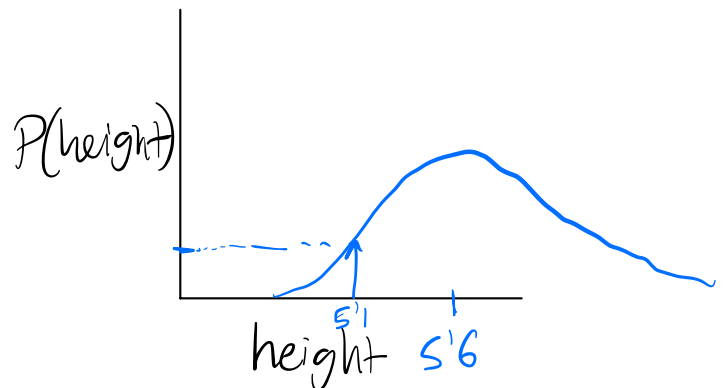


Examples:

Flip a Fair coin



grow a human to adulthood



Statistics

If all you have is data, what can you say about the generating process?

Histograms are essentially an empirical measurement of a PDF.

(see notebook)

Ex. 4

Summary Statistics

distill data to fewer numbers

Central tendency measures

2 2 2 3 2 99



(arithmetic)
- mean $\frac{1}{n} \sum_{i=1}^n x_i$

- median middle value



- quantiles

- geometric mean $\left(\prod_{i=1}^n a_i \right)^{\frac{1}{n}}$



↑ 25th %ile ↑ 75th %ile

↑ use for ratios

Variability measures

- standard deviation

$$\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}} = \sigma$$

- variance σ^2

Ex. #5

