## Week 6 Section
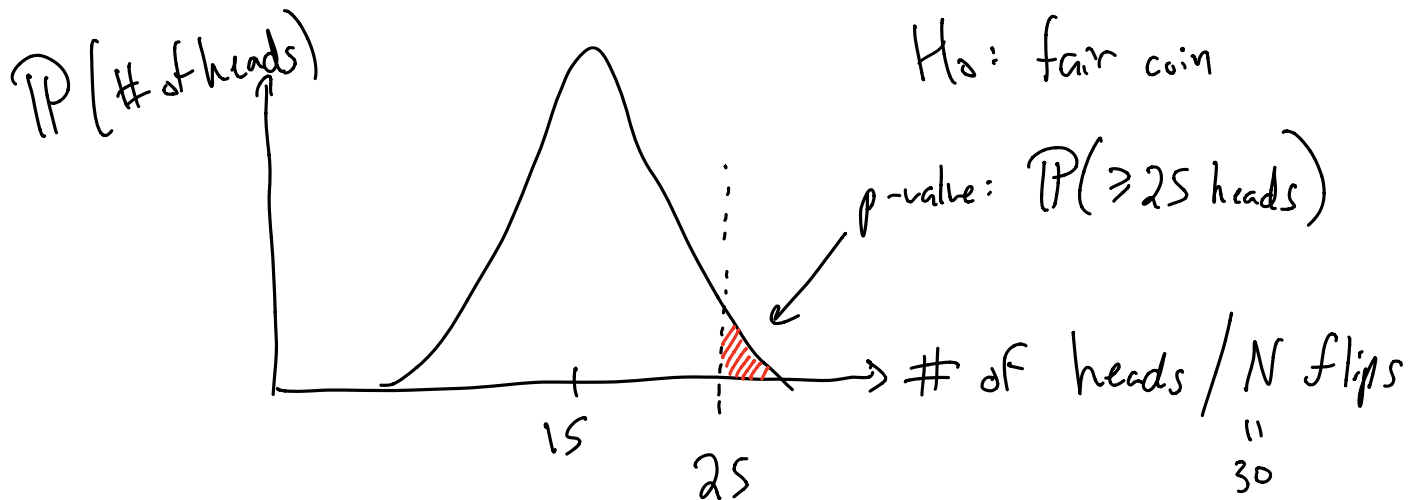
Some of what we discussed in lecture:

→ given some null hypothesis on how data is generated, how surprising is a particular observation of data?

$\mathbb{P}(\text{\# of heads})$

Ho: fair coin

p-value: $\mathbb{P}(\geq 25 \text{ heads})$

→ \# of heads / N flips

15

25

$\overset{||}{30}$

For section today, a mix of things but they're all connected---

→ given data, how do we make a <u>new</u> hypothesis?

→ revisit Bayes Rule
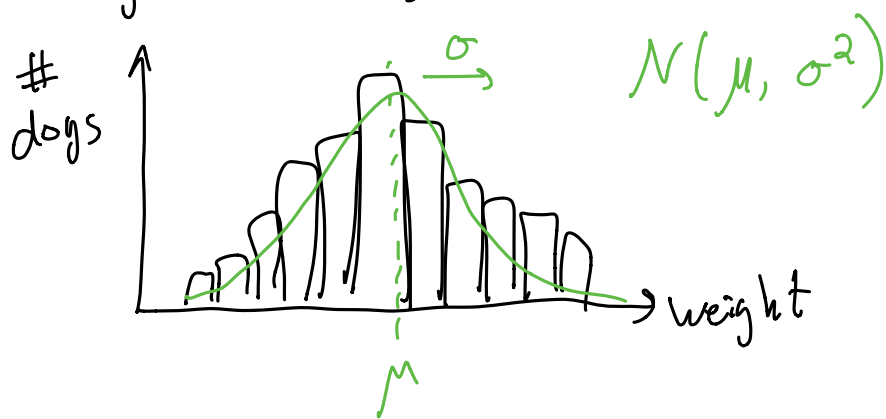
→ estimation of parameters for binomial & normal

→ see how T scores arise from estimation of $\int$

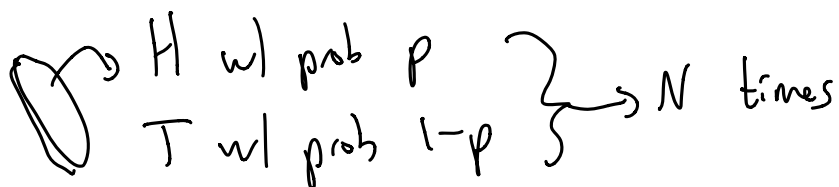→ concepts behind p-set

# Parameter Estimation / Inference

Given some data, $X_1, ..., X_N$, we'd like to describe the process that produced the data.

Ex. 1: Weights of Siberian huskies



$$N(\mu, \sigma^2)$$

$\Rightarrow$ what particular $\mu, \sigma$ describe $X_1, ..., X_n$?

Ex. 2: # of heads out of N flips of a coin?

H w/ prob $p$
T w/ prob $1-p$ $\Rightarrow$ N times

$\Rightarrow$ what particular $p$ describes # of heads / flips?

Notation: $\theta$ is a _hypothesis_ on how the data was generated
     ↳ particular values of the parameters of the underlying
                                        probability distribution.

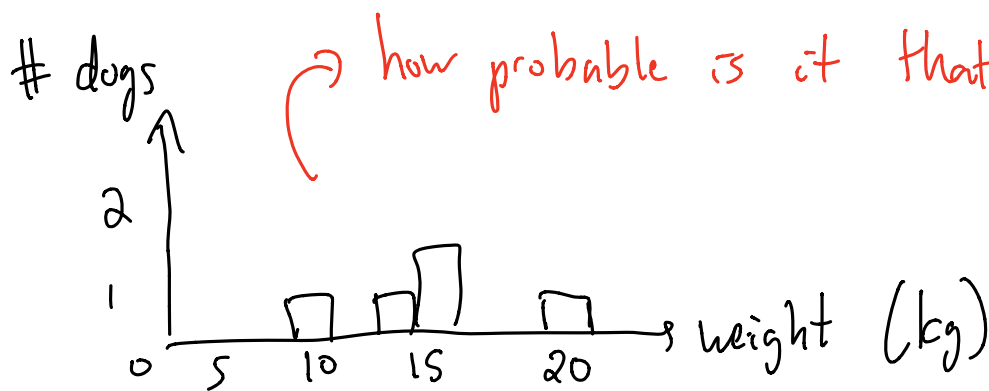Ex. for dog weights: $\theta = \{\mu = 15 \text{ kg}, \sigma = 5 \text{ kg}\}$

Ex. for an unbiased coin: $\theta = \{p = 0.5\}$

The big question we'd like to answer:

> given data $D$, what's the most probable $\theta$?

$(\theta$: a hypothesis on how data $D$ was generated$)$

Ex. I weigh 5 dogs and make a histogram:



→ how probable is it that $\mu = 5$? meh.
$\mu = 15$? probable!

Ex. I get 25 heads in 30 flips
     ↳ "fair" hypothesis: how probable is it
                        that $p = \mathbb{P}(\text{heads}) = 0.5$?

Suppose we had some set of hypotheses, $\theta_1, \ldots, \theta_M$

Given data, how probable is a particular hypothesis $\theta_k$?

**likelihood:**
how probable is the data given our hypothesis is true

**prior:**
how probable is our hypothesis prior to seeing data?

$$\mathbb{P}\left(\theta_k | D\right) = \frac{\mathbb{P}(D | \theta_k)\, \mathbb{P}(\theta_k)}{\mathbb{P}(D)}$$

**posterior:**
how probable is our hypothesis $\theta_k$ post seeing the data?

**Marginal:**
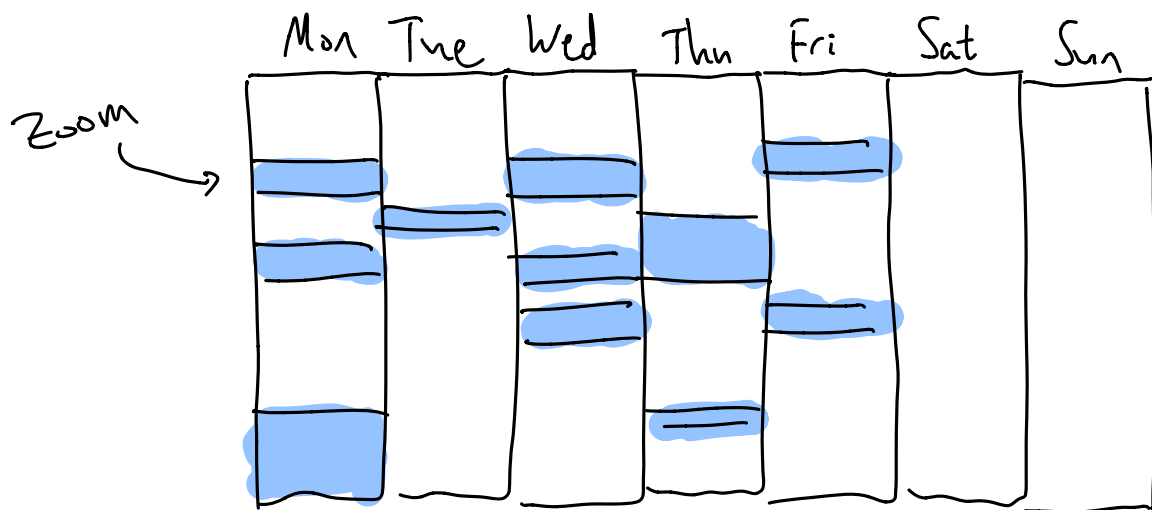how probable is the data under all possible hypotheses?

Marginal recap:

$$\mathbb{P}(D) = \sum_{i=1}^{M} \mathbb{P}(D \ \overline{AND} \ \theta_i \text{ true}) \leftarrow \text{add over all possible hypotheses}$$

$$= \sum_{i=1}^{M} \mathbb{P}(D | \theta_i)\, \mathbb{P}(\theta_i) \leftarrow \begin{array}{l} \mathbb{P}(X \text{ AND } Y) \\ = \mathbb{P}(X|Y)\,\mathbb{P}(Y) \end{array}$$

$$= \int \mathbb{P}(D | \theta)\, \mathbb{P}(\theta)\, d\theta$$

if an $\infty$ number of $\theta$'s

# An example: My Weekly Schedule



Zoom →

$\Theta$: What day it is.     Monday, Thursday, who knows...

Data: how many Zoom meetings I had today.

$\Theta \xrightarrow{\text{generates}} $ data, but also data $\xrightarrow{\text{infer}} \Theta$

Ex. Marginal:

$\mathbb{P}(3 \text{ meetings}) = \mathbb{P}(3 \text{ meetings AND it's Monday})$
$$+ \cdots +$$
$$\mathbb{P}(3 \text{ meetings AND it's Sunday})$$

$$= \sum_{i=1}^{7} \mathbb{P}(3 \text{ meetings AND it's } i^{th} \text{ day of week})$$

$$= \sum_{i=1}^{7} \mathbb{P}(3 \text{ meetings} \mid i^{th} \text{ day}) \, \mathbb{P}(i^{th} \text{ day})$$

## Comparing two posteriors / hypotheses

$$\mathbb{P}(\text{Monday} \mid 3 \text{ meetings}) = \frac{\mathbb{P}(3 \text{ meetings} \mid \text{Mon}) \, \mathbb{P}(\text{Mon})}{\mathbb{P}(3 \text{ meetings})}$$

$$\mathbb{P}(\text{Thursday} \mid 3 \text{ meetings}) = \frac{\mathbb{P}(3 \text{ meetings} \mid \text{Thu}) \, \mathbb{P}(\text{Thu})}{\mathbb{P}(3 \text{ meetings})}$$

To compare these two, I can take a ratio.
The denominator cancels out!

## Back to our q:

given data $D$, what's the most probable $\theta$?

We can scan over a lot of $\theta$'s to look
for a particular $\theta$ w/ the highest $\mathbb{P}(\theta \mid D)$,

$$\mathbb{P}(\theta \mid D) = \frac{\mathbb{P}(D \mid \theta) \, \mathbb{P}(\theta)}{\cancel{\mathbb{P}(D)}}$$

(can ignore denom when comparing posteriors generally)

If we have some prior beliefs $\mathbb{P}(\theta)$,

(it just feels like a Thursday...)

the $\theta$ that maximizes $\mathbb{P}(D|\theta)\,\mathbb{P}(\theta)$ is

the <u>maximum a posteriori</u> (MAP) estimate.

If we further assume uniform priors on $\theta$...

$$\mathbb{P}(\theta|D) = \frac{\mathbb{P}(D|\theta)\,\cancel{\mathbb{P}(\theta)}}{\cancel{\mathbb{P}(D)}}$$
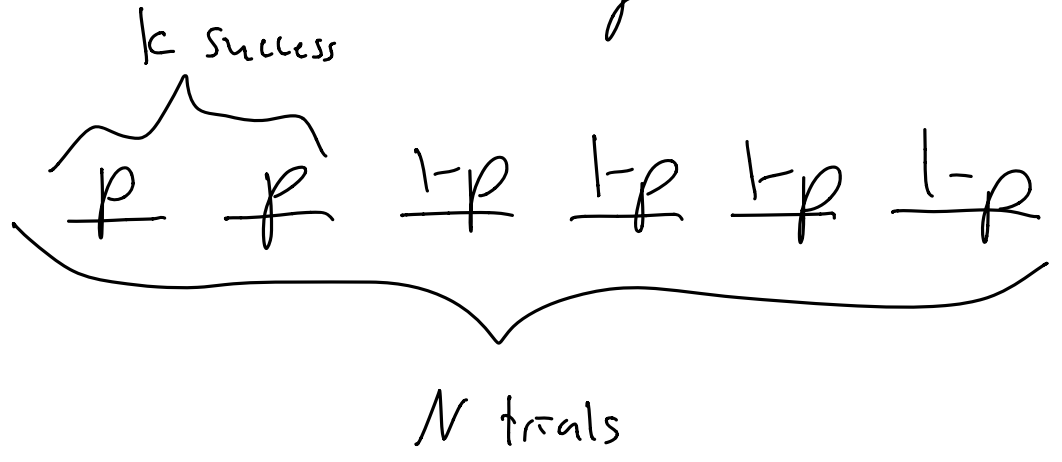
then... the most probable $\theta$ given the data
is the $\theta$ that maximizes the likelihood

$$\mathbb{P}(\theta|D) \propto \mathbb{P}(D|\theta)$$
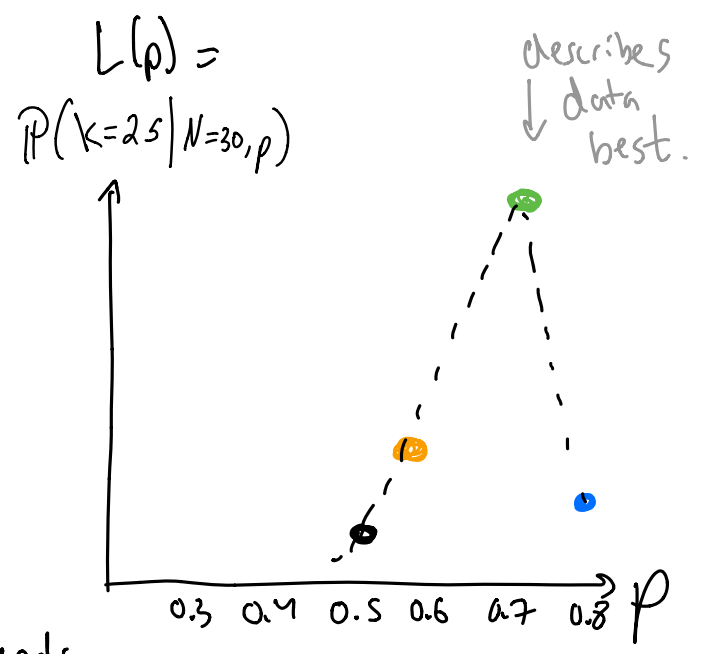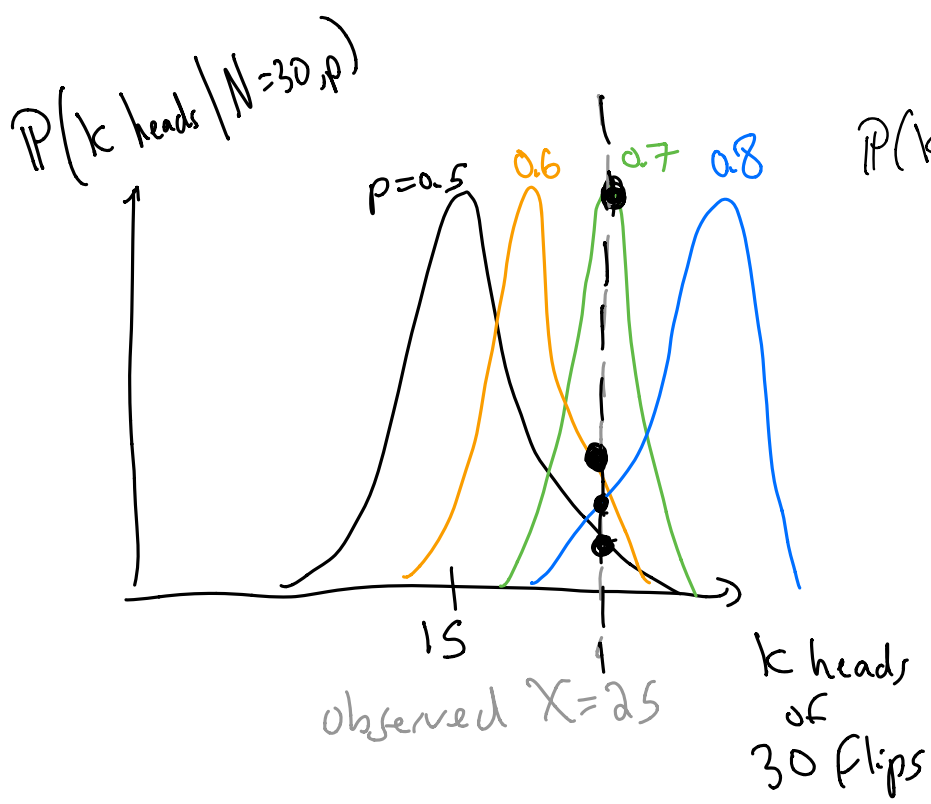
This is Maximum Likelihood Estimation!

# Concept behind likelihood $\mathbb{P}(D|\theta)$ w/ binomial

binomial process:   $n$ trials

each trial has success prob. $p$

how many successes out of $N$ trials?

$k$ success

$$\underbrace{\underbrace{p \quad p}_{} \quad 1-p \quad 1-p \quad 1-p \quad 1-p}_{N \text{ trials}}$$

$$\mathbb{P}\left(X=k \text{ successes} \mid N, p\right) = \binom{N}{k} p^k (1-p)^{N-k}$$



$\mathbb{P}(k \text{ heads} \mid N=30, p)$

$p=0.5$   0.6   0.7   0.8

15

observed $X=25$

$k$ heads of 30 flips

$L(p) =$
$\mathbb{P}(k=25 \mid N=30, p)$

describes data best.

0.3  0.4  0.5  0.6  0.7  0.8   $p$

Key point: the likelihood measures overlap of data w/ a data generation/probability process parameterized by $\theta$.

## Max Likelihood of Binomial

Observed data: 25 heads / 30 flips.

What's the max. likelihood estimate of $p$?

The $p$ that maximizes the likelihood.

$$\mathbb{P}(k \mid N, p) = L(p) = \binom{N}{k} p^{k} (1-p)^{N-k}$$

This is a function of $p$. We can maximize it by finding $p$ s.t. $\frac{\partial}{\partial p} L(p) = 0$.

$$\frac{\partial}{\partial p} L(p) = \frac{\partial}{\partial p} \left[ \binom{N}{k} p^{k} (1-p)^{N-k} \right]$$

gross... chain rule --- take the log!

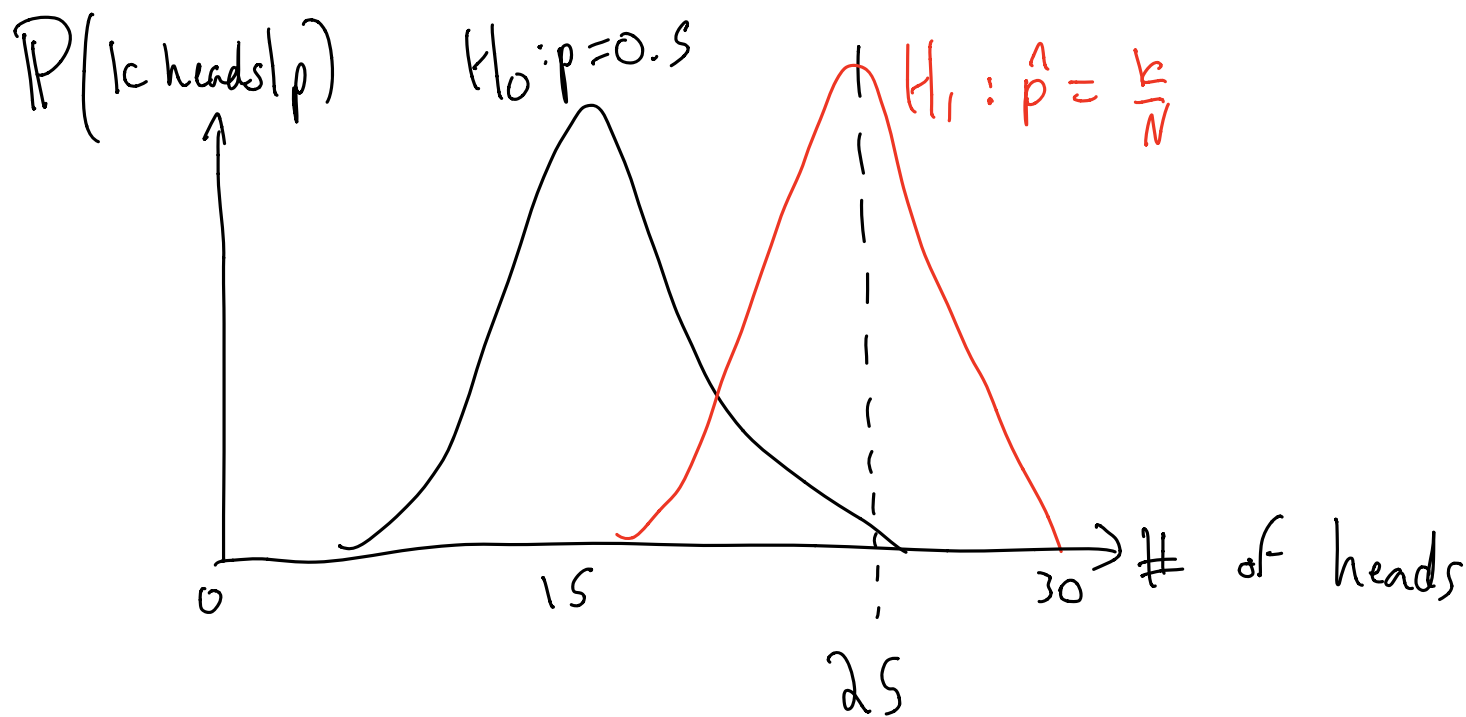$$\log L(p) = \log \binom{N}{k} + k \log p + (N-k) \log(1-p)$$

$$\frac{\partial}{\partial p} \log L(p) = \frac{k}{p} - \frac{N-k}{1-p} = 0$$

$$\Rightarrow \frac{k}{p} = \frac{N-k}{1-p}$$

$$\Rightarrow k - kp = Np - kp \Longrightarrow \hat{p} = \frac{k}{N}$$

So, given we observe k heads out of N flips, the max. likelihood estimate of $p$ is $\frac{k}{N}$.

$\mathbb{P}\left(k \text{ heads} | p\right)$    $H_0 : p = 0.5$    $H_1 : \hat{p} = \frac{k}{N}$

# of heads

0          15          30

25

Things to consider:
- is $H_0$ invalidated?
- does this mean $H_1$ has to be the correct model?
- what if we had prior beliefs on
  $\mathbb{P}(\theta) = \mathbb{P}(\text{prob. of heads})$?
  (maybe we really trust the coin is fair)

## MLE for a normal: $\quad X_1, \ldots, X_n \sim N(\mu, \sigma^2)$

For one data point,

$$\mathbb{P}(X = x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

For $n$ data points, write the joint likelihood:

$$L(\mu, \sigma^2) = \mathbb{P}(X_1 = x_1, \ldots, X_n = x_n \mid \mu, \sigma^2)$$

$$= \mathbb{P}(X_1 = x_1 \mid \mu, \sigma^2) \cdot \ldots \cdot \mathbb{P}(X_n = x_n \mid \mu, \sigma^2)$$

$$= \prod_{i=1}^{n} \mathbb{P}(X_i = x_i \mid \mu, \sigma^2)$$

$$= \prod_{i=1}^{n} (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$= (2\pi\sigma^2)^{-\frac{n}{2}} \prod_{i=1}^{n} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\sum_{i=1}^{n} \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

Take the log:

$$\log L(\mu, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \sum_{i=1}^{n} \frac{(x_i - \mu)^2}{2\sigma^2}$$

$\hookrightarrow$ we'll find MLE again by taking derivatives

## MLE for $\mu$:

$$\frac{\partial}{\partial \mu} \log L(\mu, \sigma^2) = \sum_{i=1}^{n} \frac{2(X_i - \mu)}{2\sigma^2} = 0$$

$$\implies \hat{\mu} = \sum_{i=1}^{n} X_i \Big/ n$$

The MLE for $\mu$ is just the sample mean.

## MLE for $\sigma^2$:

$$\frac{\partial}{\partial \sigma^2} \log L(\mu, \sigma^2) = -\frac{n}{2\sigma^2} + \sum_{i=1}^{n} \frac{(X_i - \mu)^2}{2(\sigma^2)^2} = 0$$

Skipping some steps---
$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \mu)^2$$

But wait--- if we only observe $X_1, ..., X_n$, we don't know $\mu$!

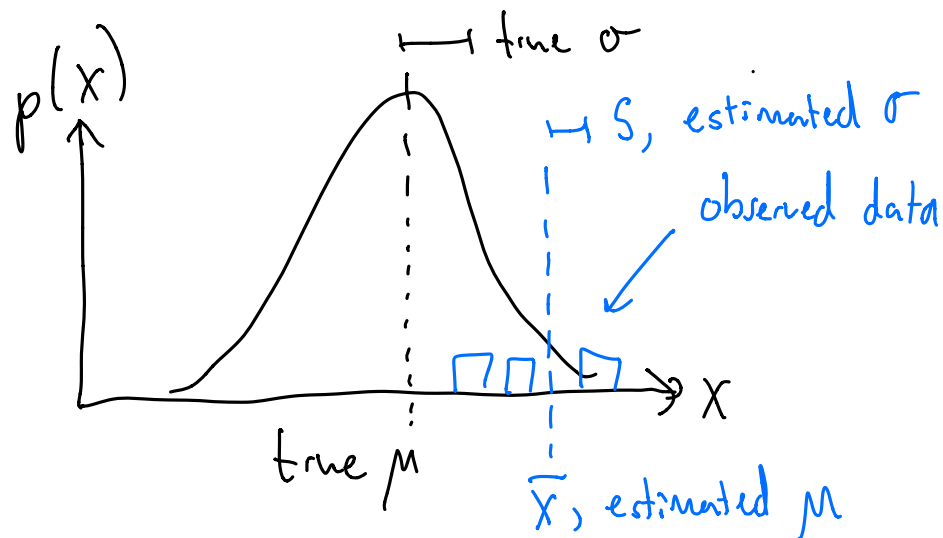Replacing $\mu$ w/ $\bar{X}$:  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X})^2$

It turns out this is a biased estimator of the population $\sigma^2$. (see Bessel's correction)

Unbiased estimate of population variance :
$$S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$$

$X_1, \ldots, X_n \overset{iid}{\sim} N(M, \sigma)$, $n$ small



true $\sigma$

$\vdash S$, estimated $\sigma$

observed data

true $M$

$\overline{X}$, estimated $M$

Q: How likely is it that $M$ is the population mean?

Compute distance b/w $M$ and $\overline{X}$, scaled by

typical fluctuation in $\overline{X}$, which is $\frac{\sigma}{\sqrt{n}}$. (standard error'' of the mean)
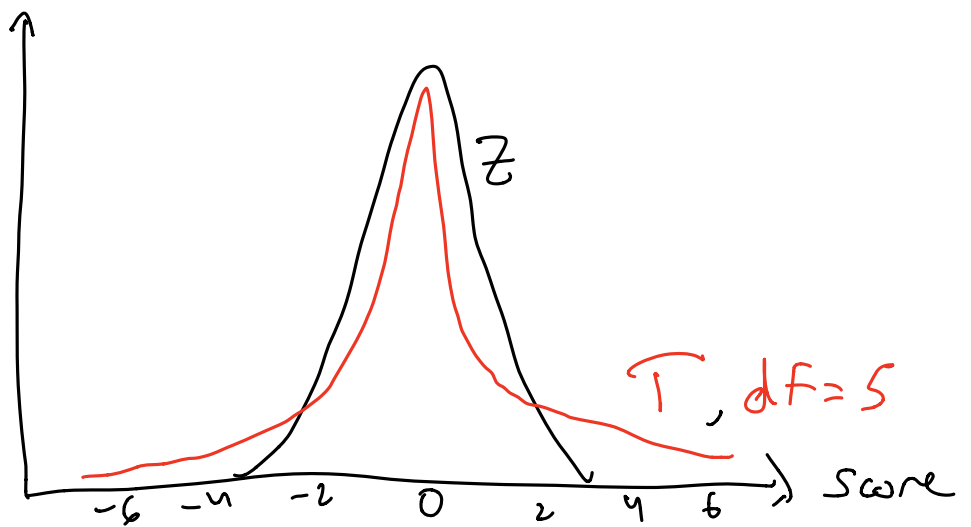
$\Rightarrow$ Standard score: $\dfrac{\overline{X} - M}{\sigma/\sqrt{n}}$

If $\underline{\sigma \text{ is known}}$, this is a $Z$ score, $Z = \dfrac{\overline{X} - M}{\sigma/\sqrt{n}}$

If $\underline{\sigma \text{ is unknown}}$, we estimate it with $S$,

this is a $T$ score: $T = \dfrac{\overline{X} - M}{S/\sqrt{n}}$, w/ $-1$ degrees of freedom

Distribution of
T, Z scores
(when $\mu$ really is
population mean)



$\Longrightarrow$ T distribution has fatter tails.
It allows for $\overline{X}$ to be "far away" from $\mu$
because we had to estimate $\sigma^2$ from the data.

$\Longrightarrow$ We could have estimated too small a $\sigma^2$, which
means the T score would be larger than it
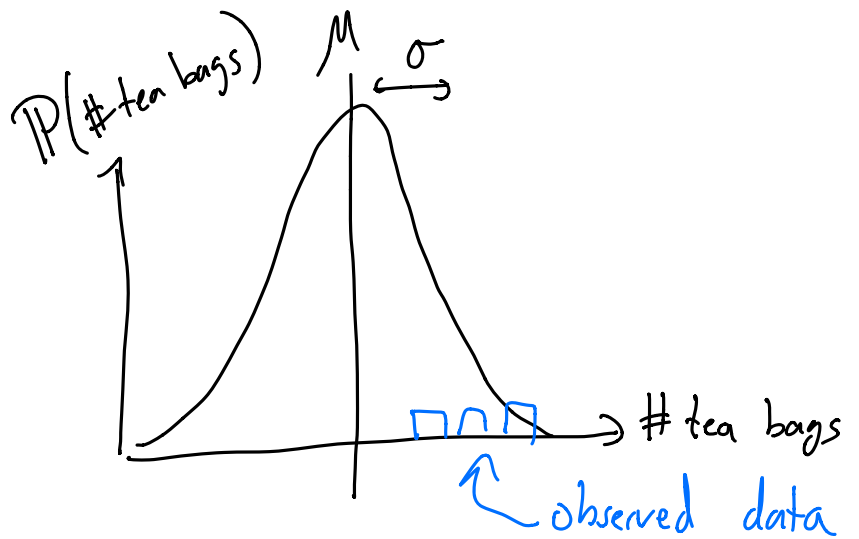should be $\rightarrow$ hence more probability of big T scores

$\Longrightarrow$ At large $n$, we estimate $\sigma^2$ well,
T dist converges to Z dist.

So the T dist can arise from estimation of
an uncertain $\sigma^2$ given data.
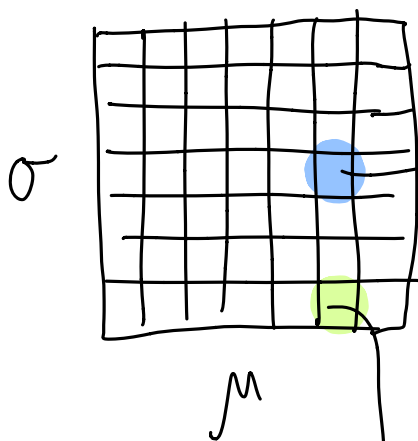On the p-set we will marginalize over many
potential $\sigma^2$'s.

$\mathbb{P}(\#\text{tea bags})$ $\mu$ $\sigma$

There's a true $N(\mu, \sigma^2)$, but we only see a few data points

→ # tea bags

← observed data

$\Longrightarrow$ Goal: <u>bet</u> on candidate $(\mu, \sigma^2)$ given only a few observed datapoints

How? We'll assume a <u>grid</u> of possible $(\mu, \sigma)$:

$\sigma$

$\mu$

one of these $(\mu, \sigma)$ pairs generated the data we observe, $X_1, \ldots, X_n$

↳ for any candidate $(\mu, \sigma^2)$ we can compute how probable they are given our data...

with a posterior, $\mathbb{P}\left(\mu, \sigma^2 \mid X_1, \ldots, X_n\right)$

$$P(M, \sigma^2 \mid X_1, \dots, X_n)$$

$$= \frac{\overset{\text{①}}{P(X_1, \dots, X_n \mid M, \sigma^2)} \ \overset{\text{②}}{P(M, \sigma^2)}}{\underset{\text{③}}{P(X_1, \dots, X_n)}}$$

① likelihood: <span style="color:orange">(independent observations)</span>

$$P(X_1, \dots, X_n \mid M, \sigma^2) = P(X_1 \mid M, \sigma^2) \times \dots \times P(X_n \mid M, \sigma^2)$$

$$= \prod_{i=1}^{n} P(X_i \mid M, \sigma^2)$$

$$= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X_i - M)^2}{2\sigma^2}}$$

② prior: $P(M, \sigma^2) = P(M) P(\sigma^2)$

③ marginal: $P(X_1, \dots, X_n) = \sum_{M} \sum_{\sigma} P(X_1, \dots, X_n \text{ AND } M, \sigma)$

$$= \sum_{M} \sum_{\sigma} P(X_1, \dots, X_n \mid M, \sigma) P(M, \sigma)$$

# Extra: distribution of a sample mean.

The expectation, or average, of $\bar{X}$:

$$E(\bar{X}) = E\left[\frac{1}{N}\sum_{i=1}^{N} X_i\right] = \frac{1}{N}\sum_{i=1}^{N} E[X_i]$$

(mean of sum $=$ sum of means)

$$= \frac{1}{N}(N\mu)$$

$$= \mu.$$

So the average $\bar{X}$ is indeed $\mu$, the population average.

Now for the spread in $\bar{X}$? Its variance:

$$Var(\bar{X}) = Var\left(\frac{1}{N}\sum_{i=1}^{N} X_i\right)$$

$$= \frac{1}{N^2} Var\left(\sum_{i=1}^{N} X_i\right)$$

$Var(aX) = a^2 Var(X)$

$$= \frac{1}{N^2}\sum_{i=1}^{N} Var(X_i)$$

independent variables

$$= \frac{1}{N^2}\sum_{i=1}^{N}\sigma^2$$

Each $X_i \sim N(\mu, \sigma^2)$

$$= \frac{1}{N^2} N\sigma^2$$

$$= \frac{\sigma^2}{N}$$

$$\Rightarrow SD(\bar{X}) = \sqrt{\frac{\sigma^2}{N}} = \frac{\sigma}{\sqrt{N}}.$$

the more $N$ we observe, the closer $\bar{X}$ is to $\mu$

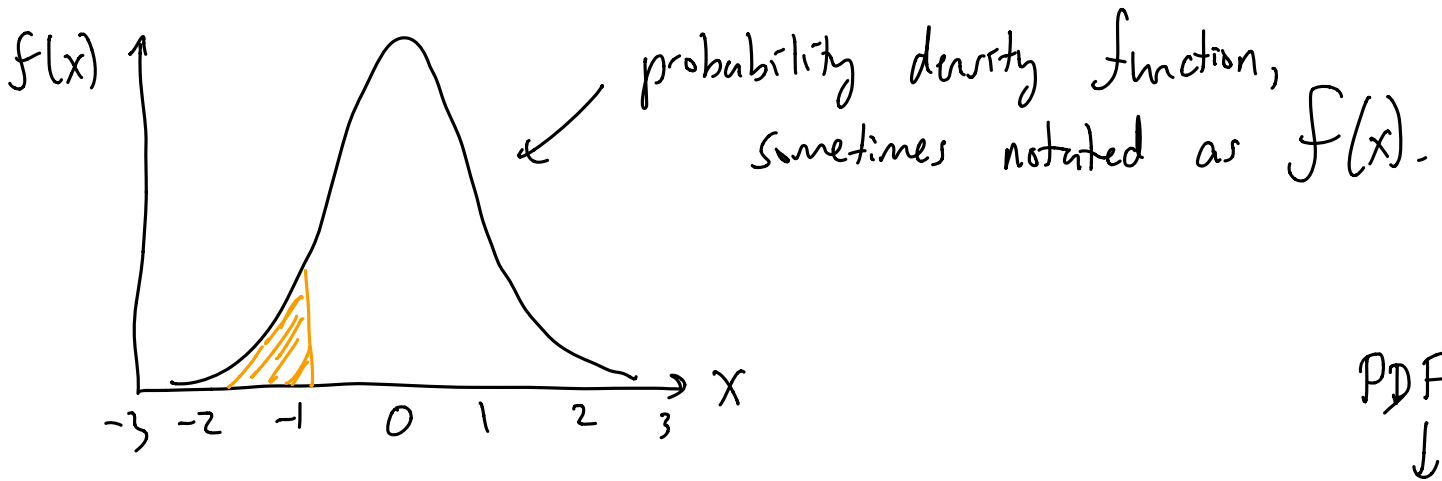# Extra: CDF example

Say we have a normal:
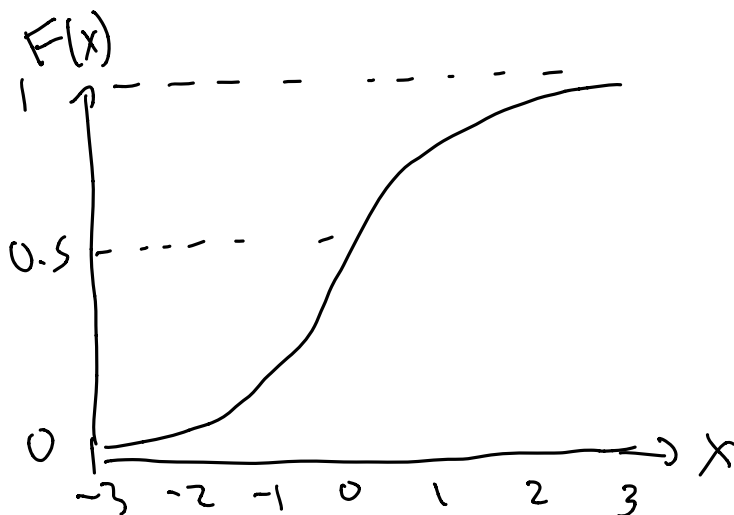


probability density function, sometimes notated as $f(x)$.

PDF
↓

$P(X \leq -1)$ is the integral from $-\infty$ to $-1$ of $f(x)$

A CDF is an integral from $-\infty$ to (some place) of $f(x)$

$$CDF(x) = \int_{-\infty}^{x} f(x') dx', \text{ often notated as } F(x).$$

So $P(X \leq -1) = CDF(-1)$

The further right we integrate to, we can only **add** to the integral, so $CDF(x)$ cannot decrease w/ x.



Q: draw the CDF of a uniform distribution



(Hint: integrate up to x)