



Computer  
Science

# **CSC380: Principles of Data Science**

## **Statistics 1**

**Chicheng Zhang**

- Expectations of HW code submission
  - For questions that ask 'paste your code', we ask you to paste your code in your solution PDF file, *in addition to* submitting to 'HW code'
- My (Chicheng's) office hours will be hybrid up until Feb 21; see Piazza
- HW2 due Feb 8 (Wed)
- HW3 out Feb 9 (Thu)

- Two players, 7 rounds. Fair game.
- $W$  := the event of you winning the game
- $S=(i,j)$  := you won  $i$  times and the opponent won  $j$  times.
- Goal: Compute  $a_{i,j} = P(W|S = (i,j))$  for every  $i$  and  $j$ .

- 1) If  $i = 4$  and  $j < 4$ ,  $a_{i,j} = 1$ . OTOH, if  $i < 4$  and  $j = 4$ ,  $a_{i,j} = 0$ . Can you see why?
- 2) Let  $R_i$  be a random variable where  $R_i = 1$  if you win round  $i$  and  $R_i = 0$  if you lose that round. Note that  $P(R_i = 0) = P(R_i = 1) = \frac{1}{2}$ .
- 3) Recall that by the law of total probability  $P(W | S = (i,j)) = P(W, R_{i+j+1} = 1 | S = (i,j)) + P(W, R_{i+j+1} = 0 | S = (i,j))$ .
- 4) By the probability chain rule  $P(W, R_{i+j+1} | S = (i,j)) = P(W | R_{i+j+1}, S = (i,j))P(R_{i+j+1} | S = (i,j))$ .
- 5) Although it requires rigorous argument, for this problem, you can take it as given that  $P(W | R_{i+j+1} = 1, S = (i,j)) = P(W | S = (i+1,j))$  and  $P(W | R_{i+j+1} = 0, S = (i,j)) = P(W | S = (i,j+1))$ . (Can you see why, intuitively?)
- 6) Find a way to write down  $a_{i,j}$  as a function of  $a_{i+1,j}$  and  $a_{i,j+1}$ . This will help you compute the answers in a recursive manner.

$$P(W \mid S = (i, j)) = P(W, R_{i+j+1} = 1 \mid S = (i, j)) + P(W, R_{i+j+1} = 0 \mid S = (i, j)) \quad [\text{Hint 3}]$$

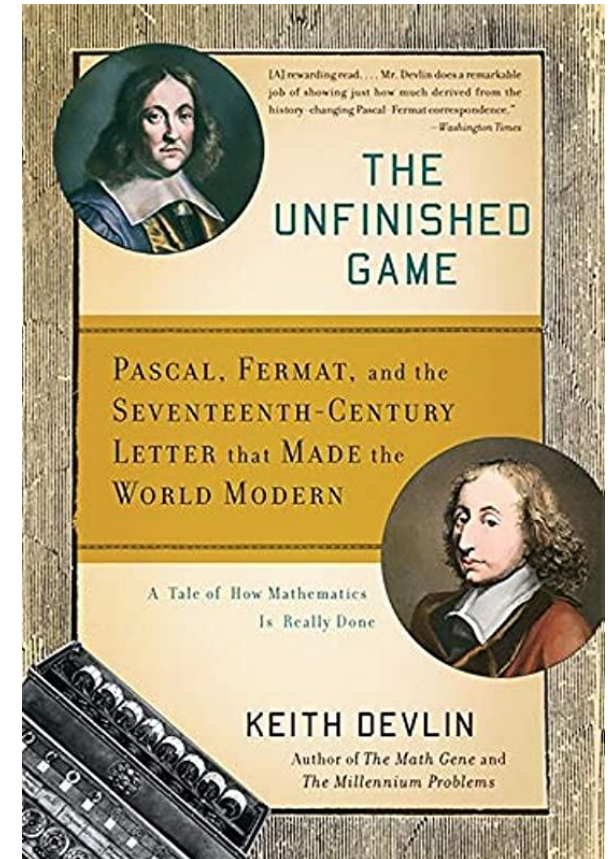
$$\begin{aligned} &= P(W \mid R_{i+j+1} = 1, S = (i, j))P(R_{i+j+1} = 1 \mid S = (i, j)) \\ &\quad + P(W \mid R_{i+j+1} = 0, S = (i, j))P(R_{i+j+1} = 0 \mid S = (i, j)) \quad [\text{Hint 4}] \end{aligned}$$

$$= P(W \mid S = (i+1, j)) \times \frac{1}{2} + P(W \mid S = (i, j+1)) \times \frac{1}{2} \quad [\text{Hint 5}]$$

$$= \frac{1}{2} \times \left( P(W \mid S = (i+1, j)) + P(W \mid S = (i, j+1)) \right)$$

- Another approach for calculating  $a_{i,j} = P(W|S = (i,j))$
- Let's say we play all 7 games anyways
- Conditioned on  $S = (i,j)$ ,  
You win  $\Leftrightarrow$  You win at least  $4 - i$  out of the remaining  $7 - i - j$  rounds
- E.g.  $a_{3,2} = P(X \geq 1)$  for  $X \sim \text{Bin}(2, 0.5)$

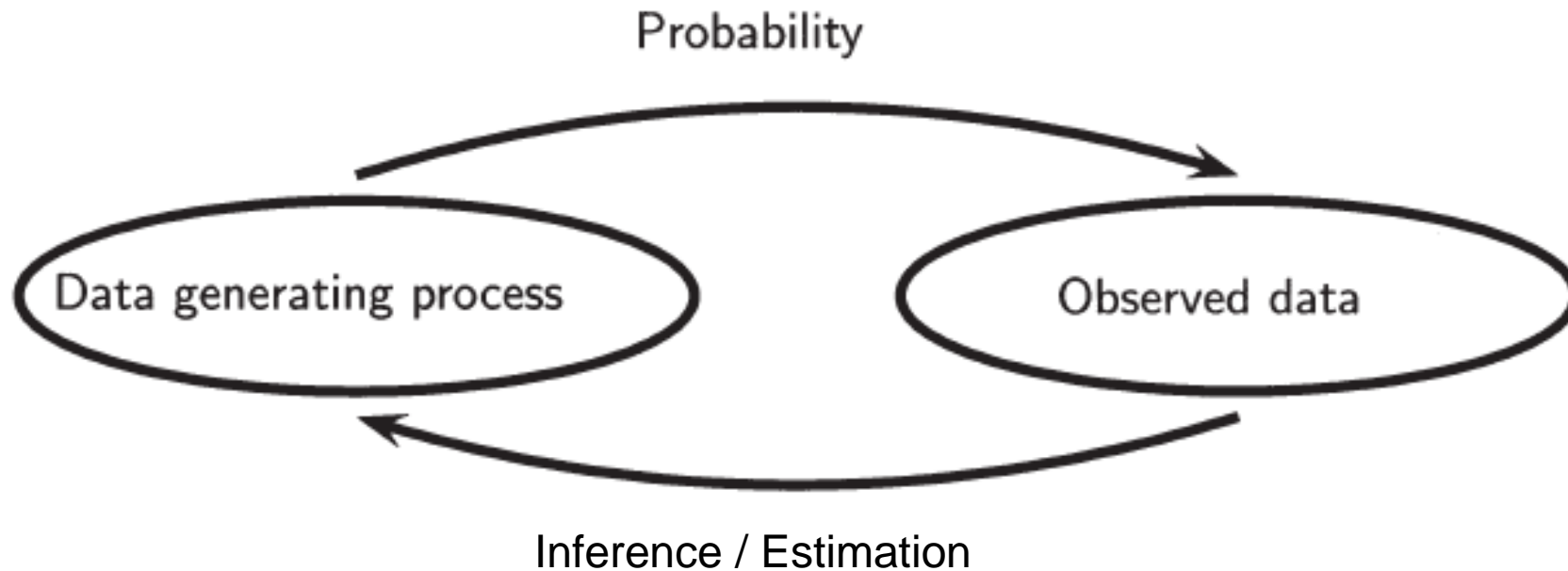
- “Problem of points”
- [https://en.wikipedia.org/wiki/Problem\\_of\\_points](https://en.wikipedia.org/wiki/Problem_of_points)
- Motivates modern probability theory



- Probability provides a mathematical formalism to reason about **random events**
  - Knowing the distribution, how can we compute probability of the event of interest? (e.g., two fair dice,  $P(\text{sum} = 3 \mid X_1 = 1)$  )
- Statistics is centered on **data**
  - Fitting models to data (estimation)
  - **E.g.**, I don't know the distribution, but I have samples drawn from it. Let's estimate what the distribution is!  $\Rightarrow$  **reverse engineering!**
  - Answering questions from data (statistical inference, hypothesis testing)
  - Interpretation of data
- Statistics *uses* probability to address these tasks

*Probability: **Given a distribution**, compute probabilities of data/events.*

E.g., If  $X_1, \dots, X_{10} \sim \text{Bernoulli}(p=.1)$ , what is the probability of  $\sum_{i=1}^{10} X_i \geq 3$ ? e.g., data = outcome of coin flip



E.g., We observed  $X_1, \dots, X_{10} \in \{0,1\}$ . What is the head probability?

*Statistics: **Given data**, compute/infer the distribution or its properties.*



*Suppose that we toss a coin 100 times. We don't know if the coin is fair or biased...*

**Question 1** Suppose that we observe 52 heads and 48 tails. Is the coin fair? Why or why not?

**Question 2** Now suppose that out of 100 tosses we observed 73 heads and 27 tails. Is the coin fair? Why or why not?

**Question 3** How might we estimate the bias of the coin with 73 heads and 27 tails?



We can model each coin toss as a Bernoulli random variable,

$$X \sim \text{Bernoulli}(\pi) \Rightarrow p(X = x) = \pi^x (1 - \pi)^{1-x}$$

Recall that  $\pi$  is the coin bias (probability of heads) and that,

$$\mathbf{E}[X] = \pi$$

Suppose we observe N coin flips  $x_1, \dots, x_N$ , estimate  $\pi$  as,

$$\hat{\pi} = \frac{1}{N} \sum_{n=1}^N x_n$$

*This is called empirical mean or sample mean*

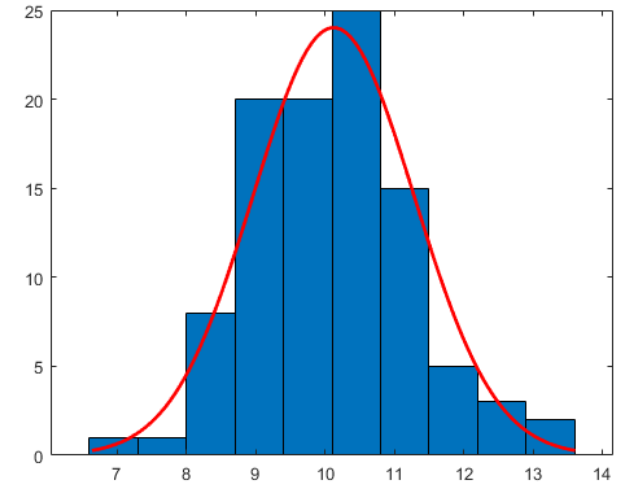
More generally, can use sample average of  $f(x_i)$ 's to estimate  $\mathbf{E}[f(X)]$  (plug-in principle)

# Estimating Gaussian Parameters

12

Suppose we observe the heights of  $N$  students at UA, and we model them as Gaussian:

$$\{x_i\}_i^N \sim \mathcal{N}(\mu, \sigma^2) \quad (\text{A.K.A. Normal})$$



How can we estimate  $\mu$ ?

$$\mu = \mathbb{E}[X] \approx \frac{1}{N} \sum_i x_i$$

Estimate  $\mu$  using sample mean

$$\hat{\mu} = \frac{1}{N} \sum_i x_i \quad (\text{abbrev. } \bar{x})$$

How can we estimate  $\sigma$ ?

$$\sigma^2 = \text{var}(X) = \mathbb{E}[(X - \mu)^2] \approx \frac{1}{N} \sum_i (x_i - \mu)^2 \approx \frac{1}{N} \sum_i (x_i - \hat{\mu})^2$$

Estimate  $\sigma$  using

$$\hat{\sigma} = \sqrt{\frac{1}{N} \sum_i (x_i - \hat{\mu})^2}$$

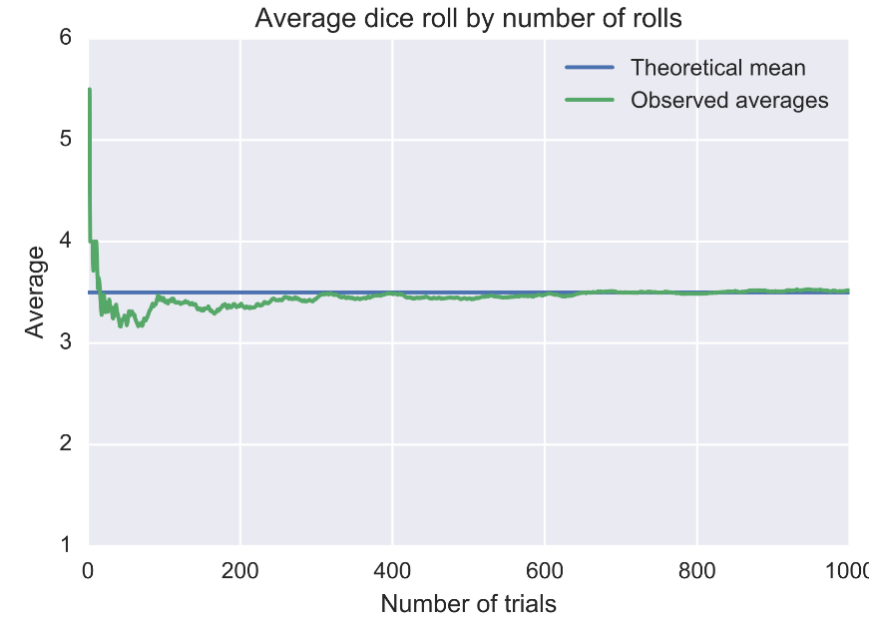
*Claim: sample mean converges to the true mean.*

(Theorem) Let  $X_1, \dots, X_N, \dots$  be drawn iid from a distribution with mean  $\mu$ . Let  $\hat{\mu}_N = \frac{1}{N} \sum_{i=1}^N X_i$  be the sample mean. Then

$$\lim_{N \rightarrow \infty} \hat{\mu}_N = \mu$$

This is the **law of large numbers**

- Weak Law: Converges to mean with high probability
- Strong Law: Stronger notion of convergence; will converge at all times! (if variance is finite)

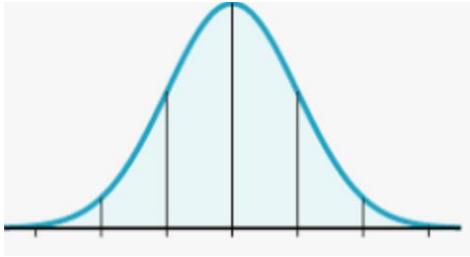


Limitation: it does not say how fast it will converge!

# Probability tool: Central Limit Theorem (CLT)

14

Let  $X_1, \dots, X_N$  be iid with mean  $\mu$  and variance  $\sigma^2$  then the sample mean  $\bar{X}_N$  approaches a Normal distribution



$$\lim_{N \rightarrow \infty} \bar{X}_N \rightarrow \mathcal{N} \left( \mu, \frac{\sigma^2}{N} \right)$$

=> the convergence rate is  $\frac{\sigma}{\sqrt{N}}$ !!

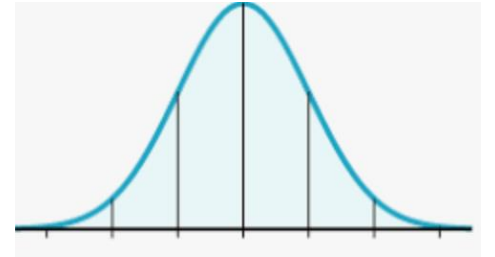
Actually, a mathematically rigorous version is

$$\lim_{N \rightarrow \infty} \frac{\sqrt{N}}{\sigma} (\bar{X}_N - \mu) \rightarrow \mathcal{N}(0, 1)$$

## Comments

- LLN says estimates  $\bar{X}_N$  “pile up” near true mean, CLT says *how* they pile up
- Very remarkable since we make **no assumption about how  $X_i$  are distributed**
- Variance of  $X_i$  **must be finite**, i.e.  $\sigma^2 < \infty$  (e.g., Cauchy distribution has  $\sigma^2 = \infty$ )

- Let  $X_1, \dots, X_N$  be drawn iid from  $\mathcal{N}(\mu, \sigma^2)$
- What's the distribution of  $\bar{X}_N$ ?



$$\Rightarrow \sum_{i=1}^N X_i \sim \mathcal{N}(N\mu, N\sigma^2)$$

$$\Rightarrow \bar{X}_N \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{N}\right)$$

$$\Leftrightarrow \frac{\sqrt{N}(\bar{X}_N - \mu)}{\sigma} \sim \mathcal{N}(0, 1)$$

## Recall: for normal distributions

- Closed under independent addition:

$$X \sim \mathcal{N}(\mu_x, \sigma_x^2) \quad Y \sim \mathcal{N}(\mu_y, \sigma_y^2) \quad , \quad X \perp Y$$

$$X + Y \sim \mathcal{N}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

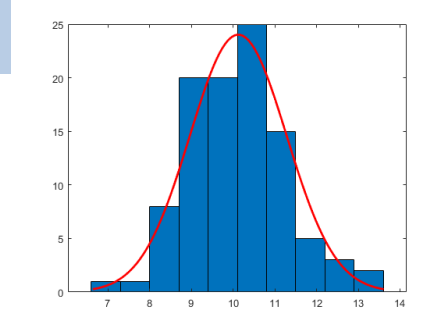
- Closed under affine transformation (a and b constant):

$$aX + b \sim \mathcal{N}(a\mu_x + b, a^2\sigma_x^2)$$

# Parameter Estimation

16

We *pose* a model in the form of a probability distribution, with unknown **parameters of interest**  $\theta$ ,



$\mathcal{D}_\theta$

e.g., assume Gaussian:  $\theta = (\mu, \sigma^2)$

Observe data, typically *independent identically distributed (iid)*,

$$p(X_1 = x_1, \dots, X_N = x_N) = p(X_1 = x_1) \cdots p(X_N = x_N)$$

$$x_1, \dots, x_N \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_\theta,$$

Compute an **estimator** to estimate parameters of interest,

$$\hat{\theta}(\{x_i\}_i^N) \approx \theta$$

*Many different types of estimators, each with different properties*

- A and B are independent but **non**identically distributed
  - E.g., two coin flips A and B with  $P(A=H) = \frac{1}{2}$  and  $P(B=H) = \frac{1}{4}$



## Dependent identical distribution

- First coin ( $X_1$ ): fair coin
- Second coin ( $X_2$ ):
  - if  $X_1=H$ , throw an unfair coin  $P(H) = \frac{1}{4}$ ,  $P(T) = \frac{3}{4}$
  - If  $X_1=T$ , throw an unfair coin  $P(H) = \frac{3}{4}$ ,  $P(T) = \frac{1}{4}$

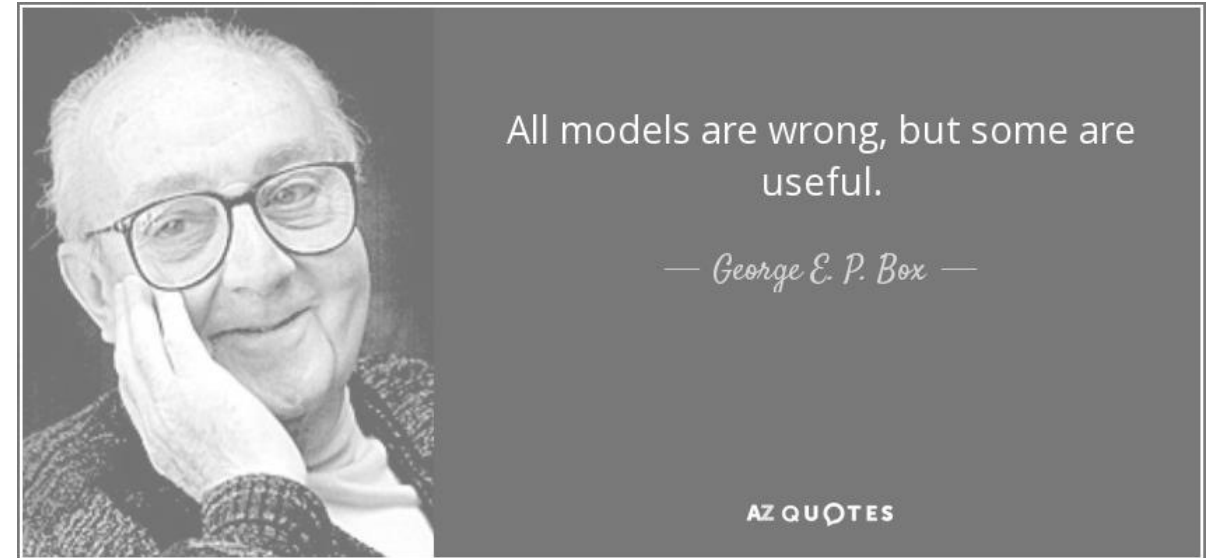
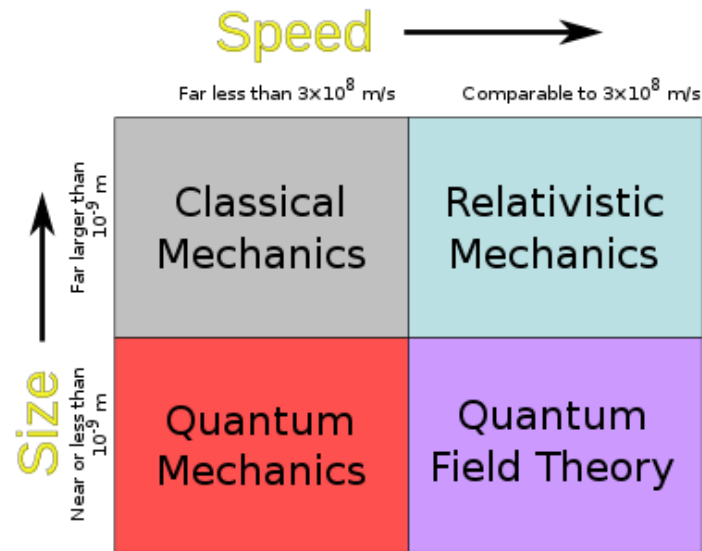
	B=H	B=T	
A=H	1/8	3/8	1/2
A=T	3/8	1/8	1/2
	1/2	1/2	

(joint probability table)

- $P(A=H)=P(B=H)$  but A and B are not independent (prove it!)

In general, i.i.d. is necessary to have estimators close to the true parameter

- In the previous example, we assumed that the heights follow a normal distribution.
- Does it?
- Another example: Physics



- There are ways to check if one model is better than the other (will be covered much later)

A **statistic** is a function of the data that does not depend on any unknown parameter.

## Examples

- Sample mean  $\bar{x}$
- Sample variance  $s^2$  or  $\hat{\sigma}^2$       Note:  $\sigma^2$
- Sample STDEV  $s$  or  $\hat{\sigma}$
- Standardized scores  $(x_i - \bar{x})/s$
- Order statistics  $x_{(1)}, x_{(2)}, \dots, x_{(n)}$
- Sample (noncentral) moments  $\bar{x}^m = \frac{1}{n} \sum_{i=1}^n x_i^m$

An **estimator**  $\hat{\theta}(x)$  is a **statistic** used to infer the unknown parameters of a statistical model.

Q: Gaussian distribution with unknown mean and variance.  
Which of these are estimators?

*Suppose that we toss a coin 100 times. We observe 52 heads and 48 tails...*

Question 1 I define an estimator that is *always*  $\hat{\theta} = 0$ , regardless of the observation. Is this an estimator? Why or why not?

Question 2 Is the estimator above a **good** estimator? Why or why not?

Question 3 What are some properties that could define a **good** estimator?



- **Consistency (asymptotic notion)** Given enough data, the estimator *converges* to the true parameter value

$$\lim_{n \rightarrow \infty} \hat{\theta}(x_1, \dots, x_n) \rightarrow \theta$$

This convergence can be measured in a number of ways: in probability, in distribution, absolutely

A bare minimum requirement!

Otherwise, you may collect more data that will give us a worse estimator!



- **Efficiency (nonasymptotic notion)** It should have low error with finite n, e.g.

$$\text{MSE}(\hat{\theta}_n) := E[(\hat{\theta}_n - \theta)^2]$$

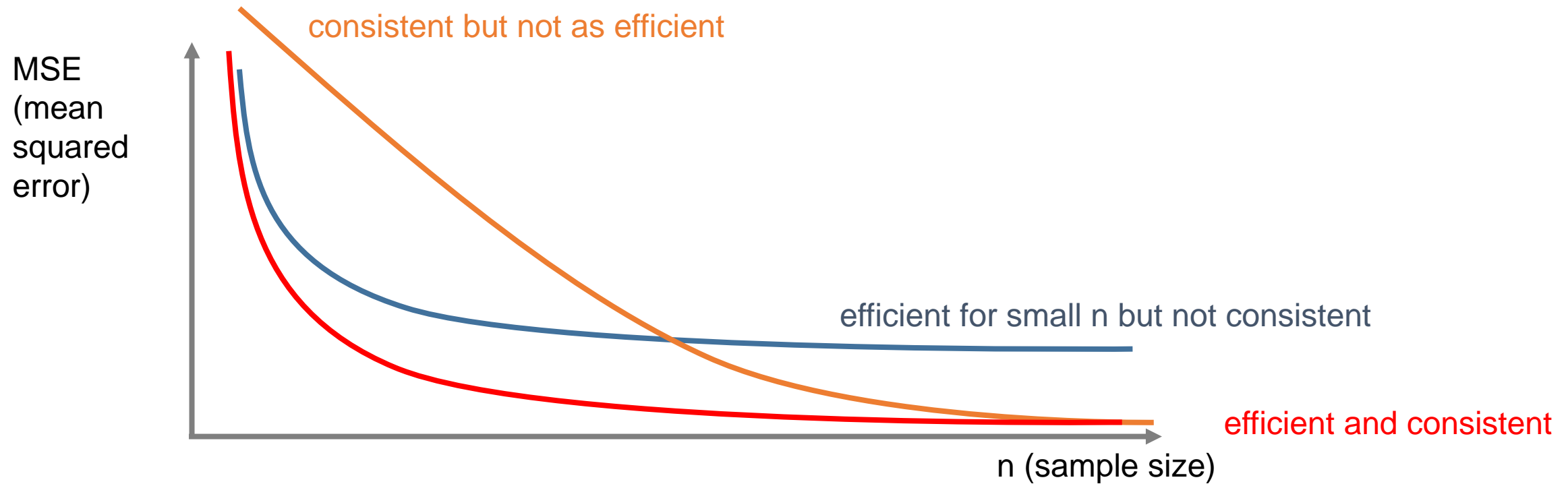
Mean squared error should be small

looks like variance but it's different!

Q: spot the difference from  $\text{Var}(\hat{\theta}_n)$ ?

# Two Desirable Estimator Properties

23



- **Unbiasedness**: For any  $n$ ,  $\mathbf{E}[\hat{\theta}(X_1, \dots, X_n)] = \theta$ 
  - E.g., sample mean is unbiased. If  $X_1, \dots, X_n \sim D$  with  $\mathbf{E}_{X \sim D}[X] = \mu$

$$\mathbf{E}[\bar{X}_N] = \frac{1}{N} \sum_i \mathbf{E}[X_i] = \mu$$

- Traditionally, considered to be a good property.
- In modern statistics, **not a necessary condition** to be a good estimator.
  - An unbiased estimator may be **less efficient** compared to some other **biased** estimator.
- Biased estimators can still be **consistent**.
- Consistency  $\approx$  asymptotically unbiased.

E.g., for some estimator  
 $E[\hat{\theta}(X_1, \dots, X_n)]$  can be  $\mu + \frac{1}{n}$

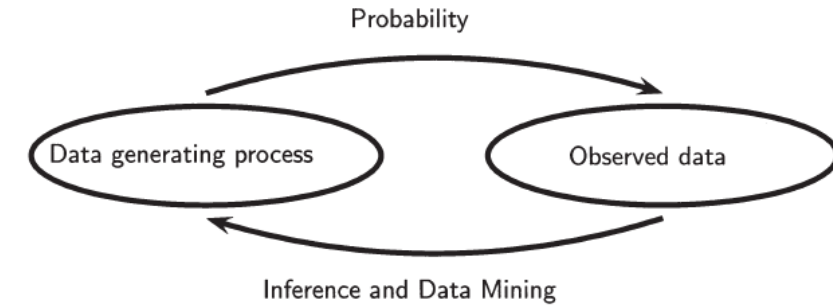


Computer  
Science

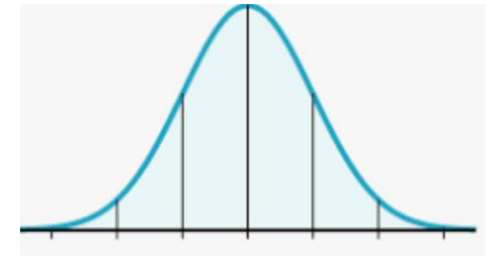
# **CSC380: Principles of Data Science**

## **Statistics 2**





- Statistics: **Given data**, compute/infer the distribution or its properties.
- Probability tools: Law of Large Numbers (LLN), Central Limit Theorem (CLT)
  - Justifies the “Plug-in principle” for estimation
- Basic setup of *estimation*:
  - Data  $x_1, \dots, x_N \sim \mathcal{D}_\theta$ ,  $\{\mathcal{D}_{\theta'}: \theta' \in \Theta\}$ : class of models parameterized by  $\theta' \in \Theta$
  - Estimator:  $\hat{\theta}(x_1, \dots, x_N)$
- Desirable properties of estimators: consistency, efficiency (MSE), unbiasedness



- Maximum likelihood estimation (MLE)
- Basic statistical properties of simple estimators: sample mean and sample variance
- The bias-variance tradeoff of statistical estimation

*Suppose that we toss a coin 100 times. We observe 73 heads and 27 tails...*

Question Let  $\theta$  be the coin bias (probability of heads). What is a more likely estimate? What is your reasoning?

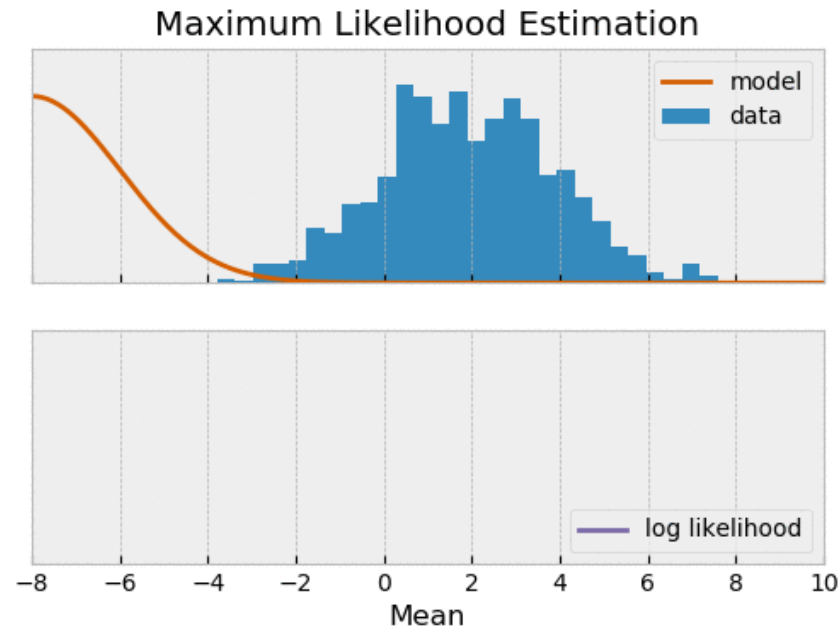
**A:**  $\hat{\theta} = 0.73$ , strong preference for heads

**B:**  $\hat{\theta} = 0.50$ , fair coin (we observed unlucky outcomes)

**Likelihood (informally)** Probability/density of the observed outcomes from a particular model



*Suppose we observe  $N$  data points from a Gaussian model  $\mathcal{N}(\mu, \sigma^2)$  and wish to estimate its mean parameter  $\mu$*

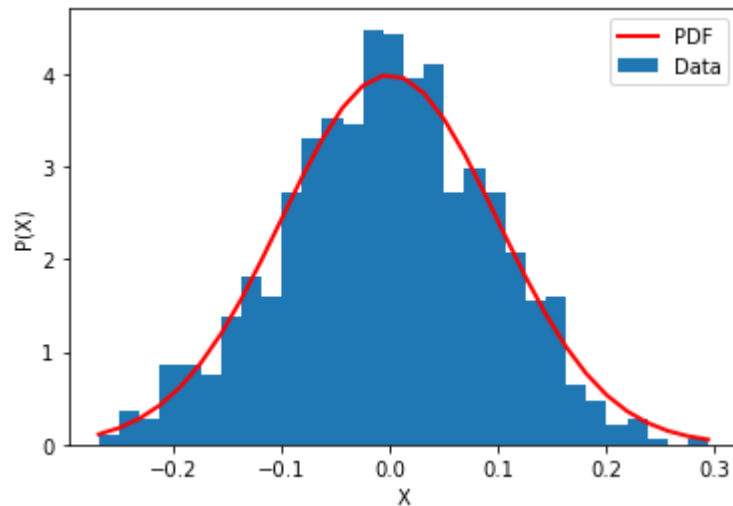


***Likelihood Principle:*** *Given a statistical model, the likelihood function describes all evidence of a parameter that is contained in the data.*

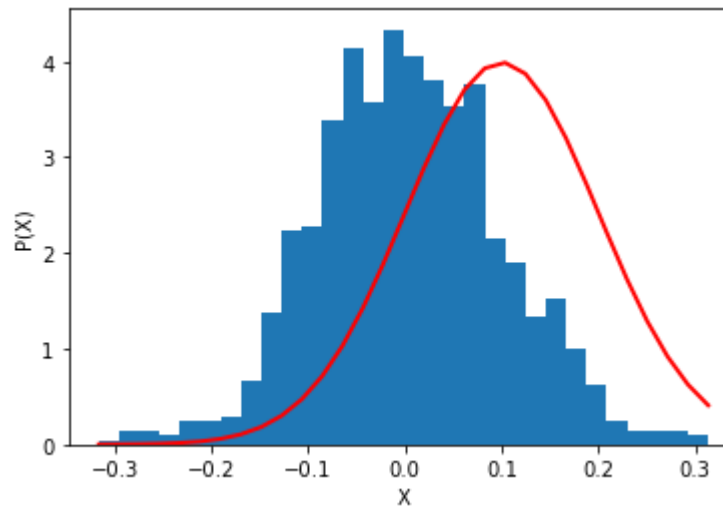
Suppose we observe  $N$  data points from a Gaussian model  $\mathcal{N}(\mu, \sigma^2)$  and wish to estimate **both**  $\mu$  and  $\sigma$

Say we only need to choose from the following three Gaussians...

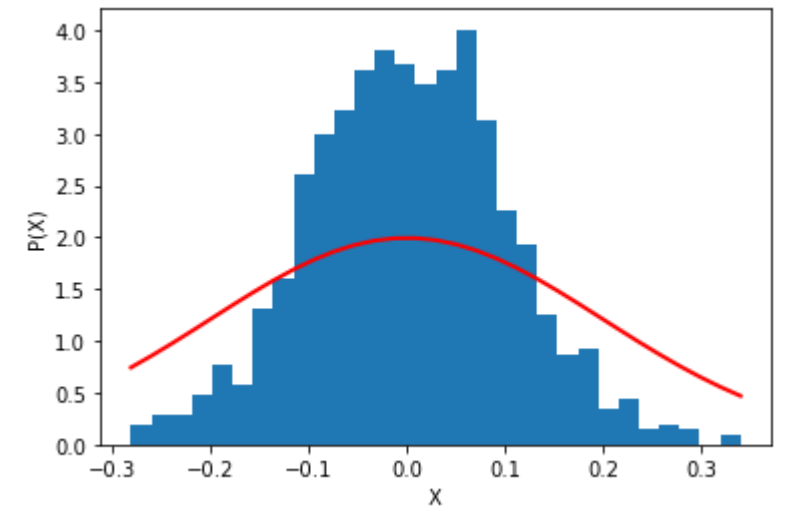
High  
Likelihood



Low  
Likelihood (mean)



Low  
Likelihood (variance)



Suppose  $x_i \sim p(x; \theta)$ , then what is the **joint probability** over  $N$  *independent identically distributed* (iid) observations  $x_1, \dots, x_N$ ?

$$p(x_1, \dots, x_N; \theta) = \prod_{i=1}^N p(x_i; \theta)$$

what appears after ‘;’ are parameters, not random variables.

- We call this the **likelihood function**, often denoted  $\mathcal{L}_N(\theta)$
- It is a function of the parameter  $\theta$ , the data are fixed
- Describes how well parameter  $\theta$  describes data (goodness of fit)

*How could we use this to estimate a parameter  $\theta$  ?*

**Maximum Likelihood Estimator (MLE)** as the name suggests, finds the parameter  $\theta$  that maximizes the likelihood function.

$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \mathcal{L}_N(\theta) = \prod_{i=1}^N p(x_i; \theta)$$

**Question** How do we find the MLE?

1. closed-form
2. iterative methods

# Finding the maximum/maximizer of a function

34

Example: Suppose  $f(\theta) = -a\theta^2 + b\theta + c$  with  $a > 0$

It is a quadratic function.

=> finding the 'flat' point suffices

Compute the gradient and set it equal to 0 (stationary points)

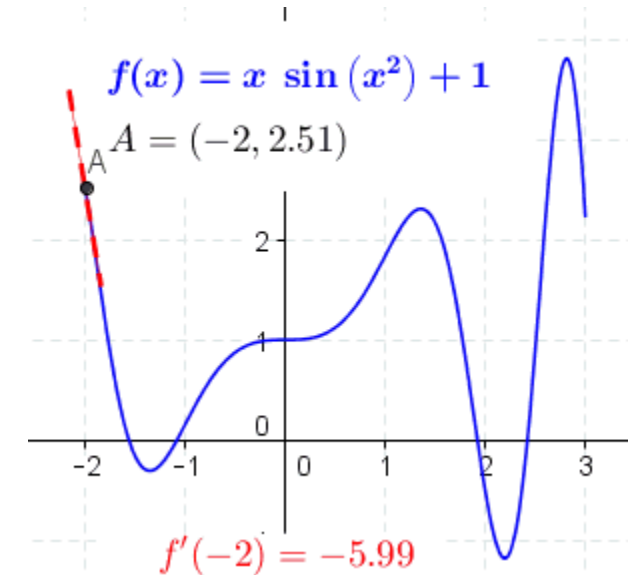
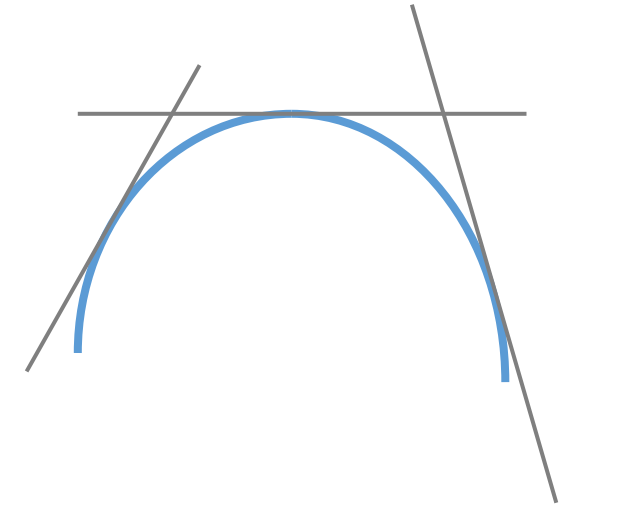
$$f'(\theta) = -2a\theta + b \Rightarrow \theta = \frac{b}{2a}$$

Closed form!

Q: Does this trick of grad=0 work for other functions?

=> Yes for **concave** functions!

=> Roughly speaking, functions that curves down only, never upwards



(gradient illustration)



# Finding the maximum/maximizer of a function

35

What if there is no closed form solution?

Example:  $f(\theta) = \frac{1}{2}x(ax - 2\log(x) + 2)$

$$f'(\theta) = ax - \log(x)$$

No known closed form for  $ax = \log(x)$

Iterative methods:

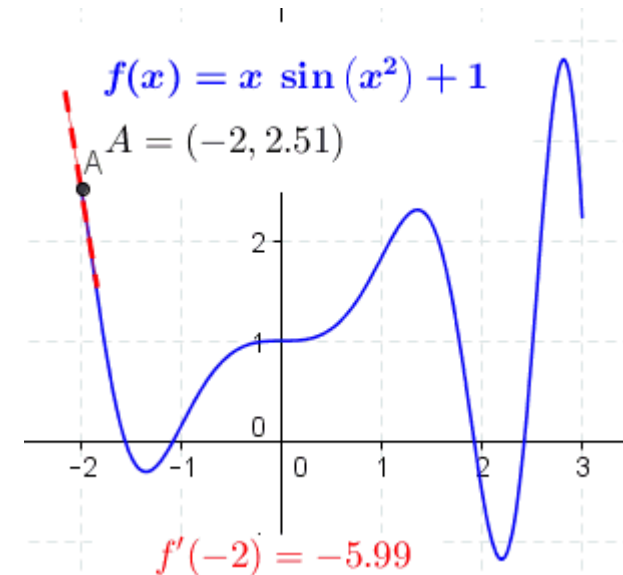
- Hillclimbing - gradient ascent (or *descent* if you are minimizing)
- Newton's method
- Etc. (beyond the scope of our class)

Iterative methods for optimization

=> Will find the global maximum

for **concave** functions (convex optimization)

=> More generally, finds a local maximum but could also get stuck at *stationary point*.



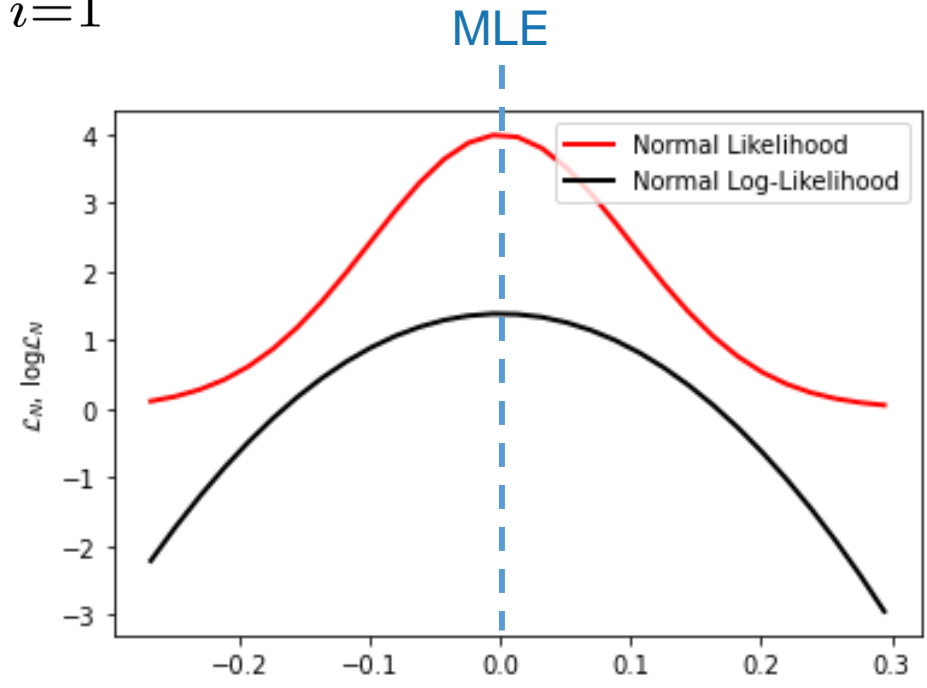
Q: find the local maxima and global maximum

Maximizing **log**-likelihood makes the math easier (as we will see) and doesn't change the answer (logarithm is an increasing function)

$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \log \mathcal{L}_N(\theta) = \sum_{i=1}^N \log p(x_i; \theta)$$

Derivative is a linear operator so,

$$\frac{d}{d\theta} \log \mathcal{L}_N(\theta) = \sum_{i=1}^N \underbrace{\frac{d}{d\theta} \log p(x_i; \theta)}_{\substack{\text{One term per data point} \\ \text{Can be computed in parallel} \\ \text{(big data)}}}$$

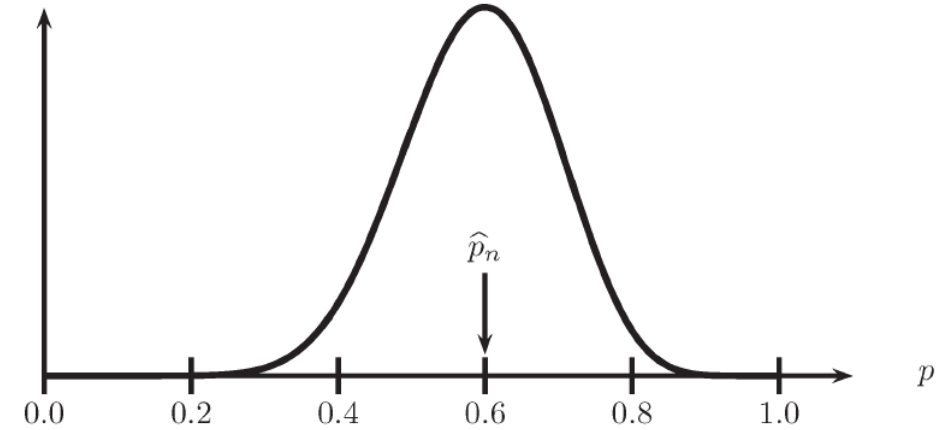


Taking log, it becomes a quadratic function!

[ Source: Wasserman, L. 2004 ]

**Example** N coin tosses with  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ . We don't know the coin bias  $p$ . The likelihood function is,

$$\mathcal{L}_n(p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^S (1-p)^{n-S}$$



*Likelihood function for Bernoulli with  $n=20$  and  $\sum_i x_i = 12$  heads*

where  $S = \sum_i x_i$ . The log-likelihood is,

$$\log \mathcal{L}_n(p) = S \log p + (n - S) \log(1 - p)$$

Set the derivative of  $\log \mathcal{L}_n(p)$  to zero and solve,

$$\hat{p}^{\text{MLE}} = S/n = \frac{1}{n} \sum_{i=1}^n x_i$$

Maximum likelihood is equivalent to sample mean in Bernoulli

⇒ this showcases how MLE is aligned to our intuition!

**Example** Let  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  with parameters  $\theta = (\mu, \sigma^2)$  and the likelihood function (ignoring some constants) is:

$$\begin{aligned}\mathcal{L}_n(\mu, \sigma) &= \prod_i \frac{1}{\sigma} \exp \left\{ -\frac{1}{2\sigma^2} (X_i - \mu)^2 \right\} \\ &= \sigma^{-n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_i (X_i - \mu)^2 \right\} \\ &= \sigma^{-n} \exp \left\{ -\frac{nS^2}{2\sigma^2} \right\} \exp \left\{ -\frac{n(\bar{X} - \mu)^2}{2\sigma^2} \right\}\end{aligned}$$

Exercise: Show that

$$\sum_i (X_i - \mu)^2 = nS^2 + n(\bar{X} - \mu)^2$$

Where  $\bar{X} = \frac{1}{n} \sum_i X_i$  and  $S^2 = \frac{1}{n} \sum_i (X_i - \bar{X})^2$  are sample mean and sample variance, respectively.

Continuing, write log-likelihood as:

$$\ell(\mu, \sigma) = -n \log \sigma - \frac{nS^2}{2\sigma^2} - \frac{n(\bar{X} - \mu)^2}{2\sigma^2}.$$

Solve zero-gradient conditions:

$$\frac{\partial \ell(\mu, \sigma)}{\partial \mu} = 0 \quad \text{and} \quad \frac{\partial \ell(\mu, \sigma)}{\partial \sigma} = 0,$$

To obtain maximum likelihood estimates of mean / variance:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i (X_i - \hat{\mu})^2$$

- The probability/density of data given parameter is mathematically the same object as likelihood of a parameter given data
- The difference is the point of view!
  - From the probabilistic perspective, the parameter is fixed and PMF/PDF is viewed as a function of the possible data
  - From the statistical perspective, the data is given (thus fixed) and we view likelihood as a function of the parameter.
- Statistics is inherently about reverse engineering.

- MLE is a very important tool.
- Usually, you write a function that computes log likelihood and then you will use existing libraries (e.g., cvxpy) to find the maximizer (next slide)
- There are efforts to develop ‘probabilistic programming’ (next slide)

```
# Import packages.  
import cvxpy as cp  
import numpy as np
```

```
# Generate data.
```

```
m = 20  
n = 15  
np.random.seed(1)  
A = np.random.randn(m, n)  
b = np.random.randn(m)
```

```
# Define and solve the CVXPY problem.
```

```
x = cp.Variable(n)  
cost = cp.sum_squares(A @ x - b)  
prob = cp.Problem(cp.Minimize(cost))  
prob.solve()
```

```
# Print result.
```

```
print("\nThe optimal value is", prob.value)  
print("The optimal x is")  
print(x.value)
```

Not the usual variable, but it is an object specifically designed to work with optimization algorithms.

Cost is not an actual value; it is an object of cvxpy that encodes the operation of  $\sum_i (\langle A_{i,:}, x \rangle - b_i)^2$  as a tree.

Alternative: `cp.Maximize`

cvxpy.Problem has many numerical methods to find the optimal solution



## Turing.jl

Bayesian inference with probabilistic programming.

aims to do 'declarative' programming for probabilistic models, just like SQL in databases!

```
using Turing
using Optim
```

```
@model function gdemo(x)
     $\sigma^2 \sim \text{InverseGamma}(2, 3)$ 
     $m \sim \text{Normal}(0, \sqrt{\sigma^2})$ 

    for i in eachindex(x)
         $x[i] \sim \text{Normal}(m, \sqrt{\sigma^2})$ 
    end
end
```

```
# Create some data to pass to the model.
```

```
data = [1.5, 2.0]
```

```
# Instantiate the gdemo model with our data.
```

```
model = gdemo(data)
```

```
# Generate a MLE estimate.
```

```
mle_estimate = optimize(model, MLE())
```

```
DynamicPPL.Model{typeof(gdemo), (:x,), (), (), Tuple{Vector{Float64}}, Tuple{},
DynamicPPL.DefaultContext}(gdemo, (x = [1.5, 2.0],), NamedTuple(),
DynamicPPL.DefaultContext())
```

ModeResult with maximized lp of -0.07  
2-element Named Vector{Float64}

A		
<hr/>		
	+	
$:\sigma^2$		0.0625
$:m$		1.75

```
# to access the value
mle_estimate.values[: $\sigma^2$ ]
```

Under some mild assumptions on the model class  $\{\mathcal{D}_{\theta'}: \theta' \in \Theta\}$ :

1) The MLE is a **consistent** estimator:

$$\lim_{n \rightarrow \infty} \hat{\theta}_n^{\text{MLE}} \xrightarrow{P} \theta_*$$

Roughly, converges to the true value.

2) The MLE is an **asymptotically efficient**: roughly, has the lowest mean squared error among all consistent estimators.

3) The MLE is an **asymptotically Normal**: roughly, the estimator (which is a random variable) approaches a Normal distribution (more later).

4) The MLE is **functionally invariant**: if  $\hat{\theta}^{\text{MLE}}$  is the MLE of  $\theta$  then  $g(\hat{\theta}^{\text{MLE}})$  is the MLE of  $g(\theta)$ .

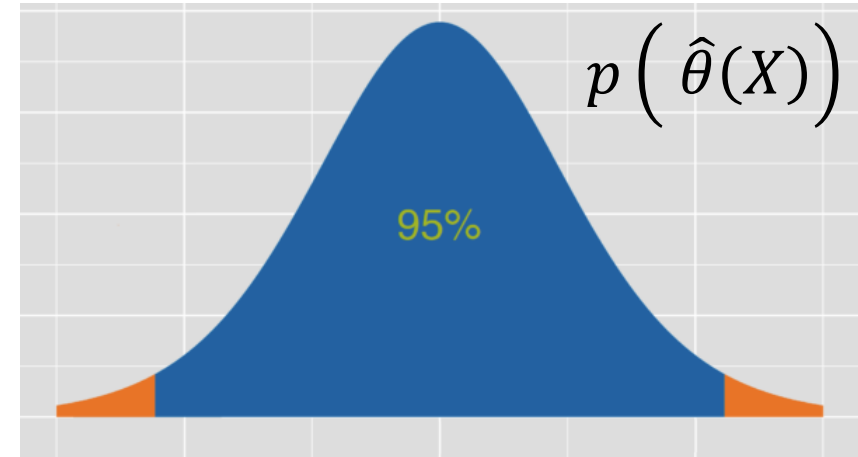
*Recall: An estimator  $\hat{\theta}$  is a RV (Random Variable).*

**Example** Let  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$   
and estimate  $\hat{p}$  be the *sample mean*,

$$\hat{p} = \frac{1}{N} \sum_i X_i$$

**Question** Is  $\hat{p}$  an unbiased estimator of  $p$ ?

Notation:  $X := (X_1, \dots, X_N)$



$$\mathbf{E}[\hat{p}(X)] = \mathbf{E} \left[ \frac{1}{N} \sum_i X_i \right] \stackrel{(a)}{=} \frac{1}{N} \sum_i \mathbf{E}[X_i] \stackrel{(b)}{=} \frac{1}{N} Np = p$$

(a) Linearity of Expectation Operator

(b) Mean of Bernoulli RV =  $p$

**Conclusion** On average  $\hat{p} = p$  (it is *unbiased*);

**In general** sample mean unbiasedly estimates mean (not nec. Bernoulli)

**Example** Let  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$  and estimate  $\hat{p}$  be the *sample mean*. Calculate the variance of  $\hat{p}$ :

quiz candidate

$$\begin{aligned}\text{Var}(\hat{p}) &= \text{Var}\left(\frac{1}{N} \sum_i X_i\right) \stackrel{(a)}{=} \frac{1}{N^2} \text{Var}\left(\sum_i X_i\right) \stackrel{(b)}{=} \frac{1}{N^2} \sum_i \text{Var}(X_i) \\ &\stackrel{(c)}{=} \frac{1}{N^2} \sum_i p(1-p) = \frac{1}{N} p(1-p) = \frac{1}{N} \text{Var}(X)\end{aligned}$$

(a)  $\text{Var}(cX) = c^2 \text{Var}(X)$

(b) Independent RVs

(c)  $\text{Var}(X) = p(1-p)$  for Bernoulli

**In General** Variance of sample mean  $\bar{X}$  for RV with variance  $\sigma^2$ ,

STDEV of sample mean  
decreases as  $1/\sqrt{N}$

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{N}$$

Decreases linearly with  
number of samples  $N$

# Unbiasedness of the Sample Variance?

Recall: Sample mean is an unbiased estimator for the true mean.

***How about the sample variance?***

**Ex.** Let  $X_1, \dots, X_N$  be drawn (iid) from any distribution with  $\text{Var}(X) = \sigma^2$  and,

$$\hat{\sigma}^2 = \frac{1}{N} \sum_i (X_i - \hat{\mu})^2$$

Then the sample variance is a **biased estimator**,

Source of bias:  
plug-in mean estimate

$$\mathbf{E}[\hat{\sigma}^2] = \frac{1}{N} \sum_i \mathbf{E}[(X_i - \hat{\mu})^2] = \text{boring algebra} = \frac{N-1}{N} \sigma^2 \quad \text{tends to underestimate}$$

Correcting bias yields unbiased variance estimator:

Q: is this estimator consistent or not?  
Consistent! (needs further justifications)

$$\hat{\sigma}_{\text{unbiased}}^2 = \frac{N}{N-1} \hat{\sigma}^2 = \frac{1}{N-1} \sum_i (X_i - \hat{\mu})^2$$

quiz candidate:

show that  $\mathbf{E} \left[ \frac{1}{N} \sum_{i=1}^N (X_i - \mu)^2 \right]$  is unbiased  
(note  $\mu = E[X_1] = \dots = E[X_N]$ )



Computer  
Science

# **CSC380: Principles of Data Science**

## **Statistics 3**

# Numpy Background

49

- Often, you have a matrix of data: e.g., movie review score

User \ Movie	Inception	Jurassic park	Batman
A	5	2	3
B	1	4	2
C	4	3	3
D	1	2	3

Numpy arrays can be 2d

```
A = np.array([[1,2,3],[4,5,6]])
```

```
A[0,1]
```

```
⇒ 2
```

```
mean(A,0)
```

```
⇒ array([2.5, 3.5, 4.5])
```

```
mean(A,1)
```

```
⇒ array([2., 5.])
```

means  $\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$

access A[0,1] means 1<sup>st</sup> row, 2<sup>nd</sup> column

computes average for each column

computes average for each row

var(A,0), var(A,1) works the same way!

**Task**: Compare the **MSE** (mean squared error) of the two variance estimators for  $N=5$ .

```
import numpy as np
```

```
import numpy.random as ra
```

```
X = ra.randn(10_000,5) # 10k by 5 matrix of  $N(0,1)$  => 10k random trials
```

```
np.mean((var(X,1,ddof=0) - 1)**2)
```

```
=> 0.36310526687176103
```

X					
X[0,:]	0.35	1.45	-0.22	-2.95	-3.09
	.....				
X[9999,:]	-1.78	-2.31	0.43	0.77	0.16

This estimates  $E[(\hat{\sigma}^2 - 1)^2]$ , the MSE of estimator  $\hat{\sigma}^2$  with respect to the target  $\sigma^2 = 1$

ddof=0 uses  $1/N$

ddof=1 uses  $1/(N-1)$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_i (X_i - \hat{\mu})^2$$

$$\hat{\sigma}_{\text{unbiased}}^2 = \frac{N}{N-1} \hat{\sigma}^2 = \frac{1}{N-1} \sum_i (X_i - \hat{\mu})^2$$



```
np.mean((var(X,1,ddof=1) - 1)**2)
```

```
⇒ 0.5071783438808787
```

```
Recall: np.mean((var(X,1,ddof=0) - 1)**2)
```

```
⇒ 0.36310526687176103
```

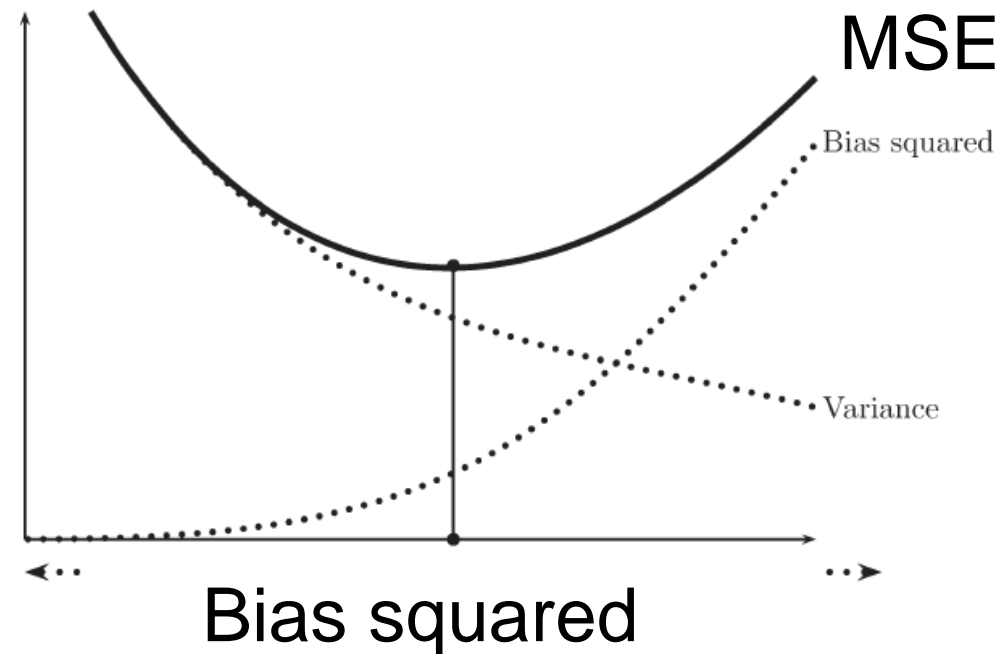
- In this case,  $\widehat{\sigma}^2$  (biased version) is more accurate than  $\hat{\sigma}_{\text{unbiased}}^2$ ! (but recall that it will underestimate)
- *There is a tradeoff between bias and variance!!*

*Is an unbiased estimator “better” than a biased one? It depends...*

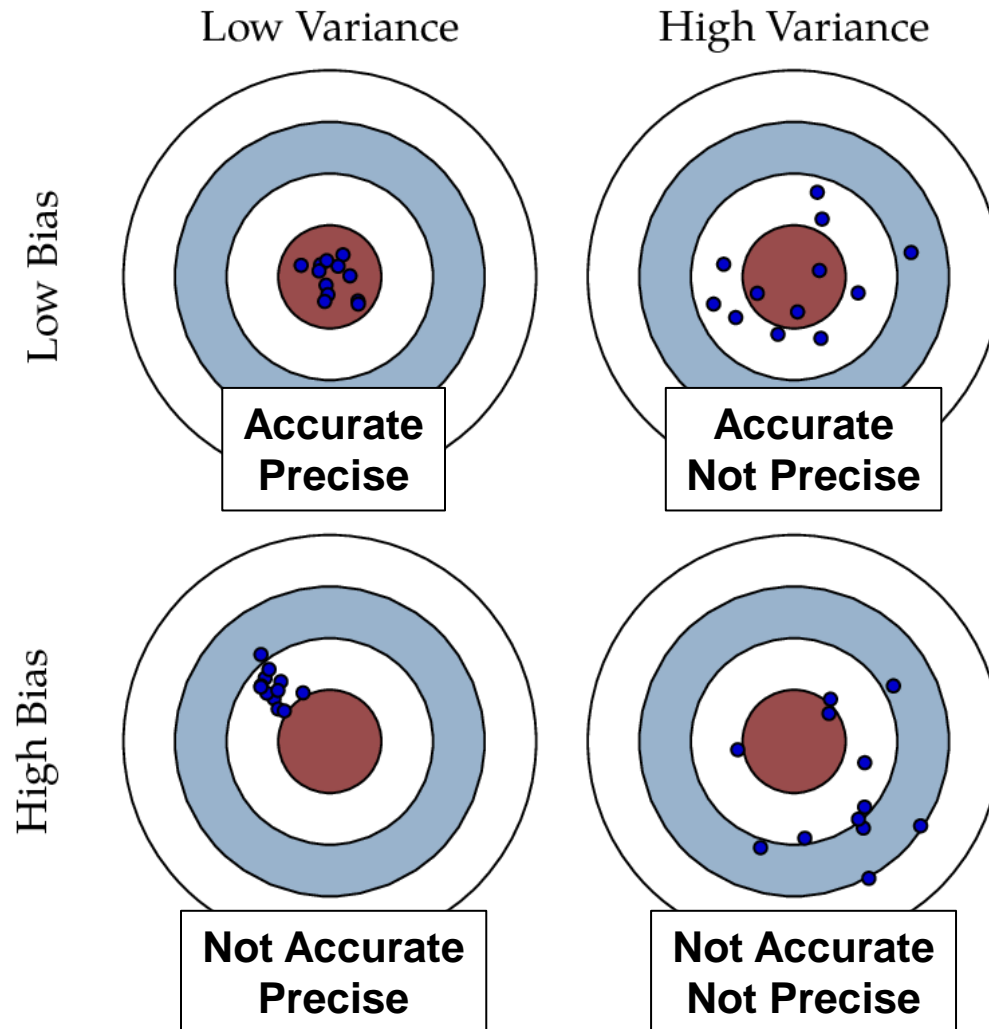
Evaluate the quality of estimate  $\hat{\theta}$  using **mean squared error**,

$$\text{MSE}(\hat{\theta}) = \mathbf{E} \left[ (\hat{\theta} - \theta)^2 \right] = \text{bias}^2(\hat{\theta}) + \mathbf{Var}(\hat{\theta})$$

- $\text{bias}(\hat{\theta}) = \mathbf{E}[\hat{\theta}] - \theta$
- MSE for unbiased estimators is just,  
$$\text{MSE}(\hat{\theta}) = \mathbf{Var}(\hat{\theta})$$
- Bias-variance is a fundamental tradeoff in statistical estimation
- MSE increases as **square** of bias
- Biased estimator can be more accurate than an unbiased one.



*Suppose an archer takes multiple shots at a target...*



- **Target** =  $\theta$
- **Each shot** = an estimate  $\hat{\theta}$
- Bias  $\approx$  systematic error
- Variance  $\approx$  random error

quiz candidate

$$\begin{aligned}\text{MSE}(\hat{\theta}) &= \mathbf{E} \left[ (\hat{\theta}(X) - \theta)^2 \right] \\ &= \mathbf{E} \left[ \left( \hat{\theta} - \mathbf{E}[\hat{\theta}] + \mathbf{E}[\hat{\theta}] - \theta \right)^2 \right] \\ &= \mathbf{E}[(\hat{\theta} - \mathbf{E}[\hat{\theta}])^2] + 2(\mathbf{E}[\hat{\theta}] - \theta)\mathbf{E}[\hat{\theta} - \mathbf{E}[\hat{\theta}]] + \mathbf{E}[(\mathbf{E}[\hat{\theta}] - \theta)^2] \\ &= \left( \mathbf{E}[\hat{\theta}] - \theta \right)^2 + \mathbf{E}[(\hat{\theta} - \mathbf{E}[\hat{\theta}])^2] \\ &= \text{bias}^2(\hat{\theta}) + \text{Var}(\hat{\theta})\end{aligned}$$

*Compare the results of two coin flip experiments...*

Experiment 1 Flip 100 times and observe 73 heads, 27 tails

Experiment 2 Flip 1,000 times and observe 730 heads, 270 tails

Question The MLE estimate of coin bias for both experiments is equivalent  $\hat{\theta} = 0.73$ . Which should we trust more? Why?

**Remark** The estimate  $\hat{\theta}(X)$  is a function of random data. So, it is a random variable. It has a distribution.

**Next lecture:** confidence intervals



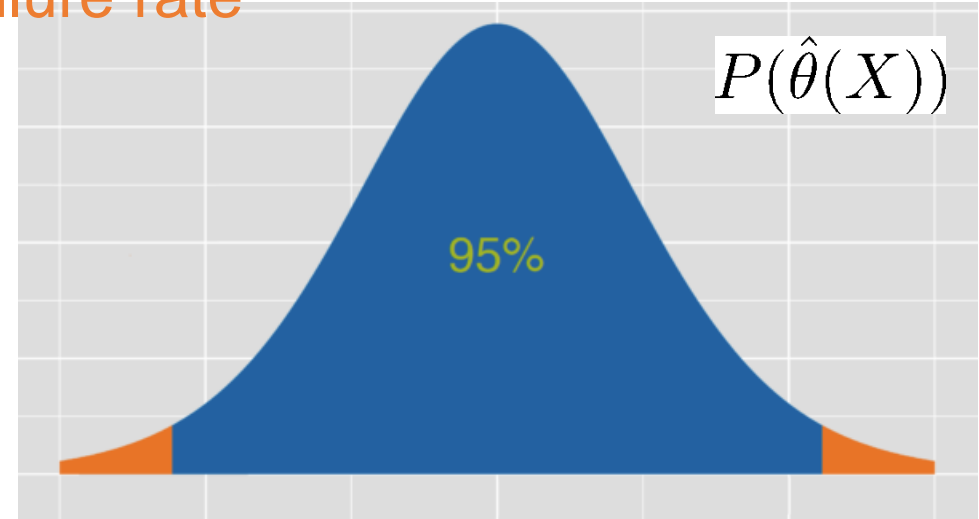
**Informally**, find an interval such that we are *pretty sure* it encompasses the true parameter value.

significance level = failure rate

Given data  $X_1, \dots, X_n$  and ~~confidence~~  $\alpha \in (0, 1)$  find interval  $(a, b)$  such that,

$$P(\theta \in (a, b)) \geq 1 - \alpha$$

**In English** the interval  $(a, b)$  contains the true parameter value  $\theta$  with probability **at least**  $1 - \alpha$



- Intervals must be computed from data: i.e.,  $a(X_1, \dots, X_n)$  and  $b(X_1, \dots, X_n)$
- Interval  $(a, b)$  is **random**, parameter  $\theta$  is **not random** (it is fixed)
- Usually, you compute an estimator  $\hat{\theta}$  and then set  $a = \hat{\theta} - \epsilon_a$  and  $b = \hat{\theta} - \epsilon_b$  for a carefully chosen  $\epsilon_a, \epsilon_b > 0$



Suppose  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  with unknown  $\mu$  & known  $\sigma^2$ . Let  $\hat{\mu} := \frac{1}{n} \sum_i X_i$ .

**(Fact 2)** If  $Z \sim \mathcal{N}(0,1)$ ,

$$P(Z \in [-z, z]) = 1 - 2(1 - \Phi(z))$$

where  $\Phi(z) := P(Z \leq z)$  is the CDF of  $Z$ .

$z = 1.96$ : RHS  $\approx .95$ , 95% confident

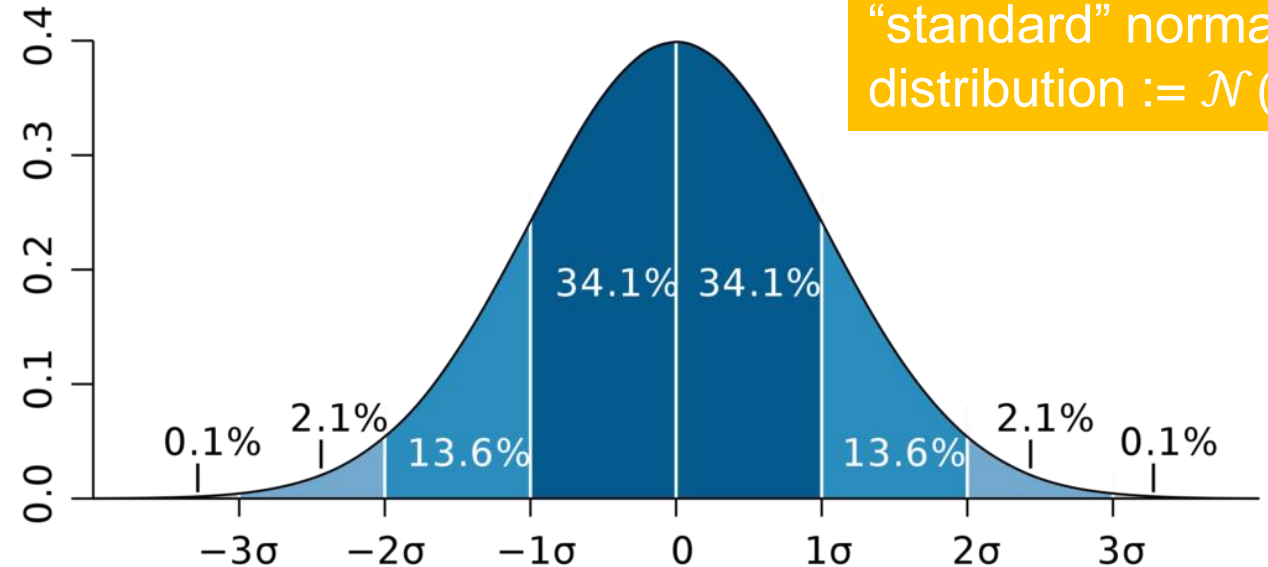
$z = 2.58$ : RHS  $\approx .99$ ,

**(Corollary)**

$$P\left(\hat{\mu} \in \left[\mu - \frac{z\sigma}{\sqrt{n}}, \mu + \frac{z\sigma}{\sqrt{n}}\right]\right) = 1 - 2(1 - \Phi(z))$$

hints: use the fact  $\sqrt{n} \frac{\hat{\mu} - \mu}{\sigma} \sim \mathcal{N}(0,1)$ . Set  $Z :=$

$\sqrt{n} \frac{\hat{\mu} - \mu}{\sigma}$  and use Fact 2.



Terminology:  
“standard” normal  
distribution  $:= \mathcal{N}(0,1)$

*Gaussians almost do not have tails!*

remember: ‘normal algebra’ is very useful  
(and will appear in exams)



Suppose  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  with unknown  $\mu$  & known  $\sigma^2$ . Let  $\hat{\mu} := \frac{1}{n} \sum_i X_i$ .

Finally, by our corollary,

$$P\left(\hat{\mu} \in \left[\mu - \frac{1.96\sigma}{\sqrt{n}}, \mu + \frac{1.96\sigma}{\sqrt{n}}\right]\right) \geq 0.95$$

$$P\left(\hat{\mu} \in \left[\mu - \frac{2.58\sigma}{\sqrt{n}}, \mu + \frac{2.58\sigma}{\sqrt{n}}\right]\right) \geq 0.99$$

This is a confidence bound for the mean  $\mu$  !!

=> Compute  $\left[\hat{\mu} - \frac{1.96\sigma}{\sqrt{n}}, \hat{\mu} + \frac{1.96\sigma}{\sqrt{n}}\right]$ . Done!

note we can switch  $\hat{\mu}$  and  $\mu$

$$P\left(\mu \in \left[\hat{\mu} - \frac{1.96\sigma}{\sqrt{n}}, \hat{\mu} + \frac{1.96\sigma}{\sqrt{n}}\right]\right) \geq 0.95$$

Q: If  $X_1, \dots, X_n$  from an arbitrary distribution, can we still use the same method?

Almost yes, if  $n$  is large enough! => central limit theorem (see later).

Question How should we interpret a confidence interval (e.g. 95%)?

$$P(\theta \in (a(X), b(X))) \geq 0.95$$

Hint Think about what is random and what is not...

This is NOT about the randomness of  $\theta$

**Wrong** If someone reveals  $\theta$  when we have exactly the same data, then  $\theta$  will be in the interval with probability at least 95%

the moment you compute the interval with the data, whether or not  $\theta$  is in the interval is determined.. you just don't know it!

*This is commonly misinterpreted*

Recommended point of view:

- Assume: Heights of UA students follow a normal distribution  $\mathcal{N}(\mu, 1)$  with unknown  $\mu$
- Fork **m parallel universes**. For each universe  $u \in \{1, 2, \dots, m\}$ ,
  - Subsample  $n$  UA students randomly, take the sample mean  $\hat{\mu}^{(u)}$ .
  - Compute the confidence bound  $\left[ \hat{\mu}^{(u)} - \frac{1.96\sigma}{\sqrt{n}}, \hat{\mu}^{(u)} + \frac{1.96\sigma}{\sqrt{n}} \right]$
- The fraction of parallel universes where the random interval includes  $\mu$  is *approximately* at least 0.95 if  $m$  is large enough.
- As  $m$  goes to infinity, the fraction will become arbitrarily close to a value that is at least 0.95.

**Recall**: If  $X_1, \dots, X_n$  from an **arbitrary** distribution, can we still use the same method used for Gaussian?

Short answer: YES, if  $n$  is large enough.

Plan

- Law of large numbers
- Central limit theorem
- 3 methods for arbitrary distributions

But first, why do we care?

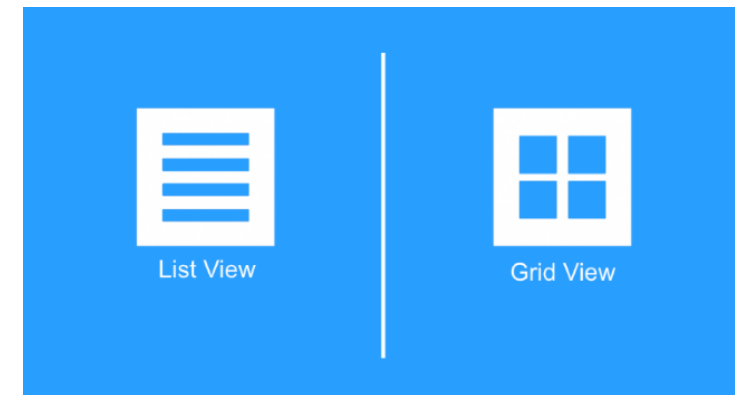
You are an engineer at amazon. You want to see if people buy more items if you change the search result from list view to grid view.

You changed it to grid view for one day. Various metrics: click rate, purchase rate, ...

You compute these, it seems to increase the click rate by 0.05%. You tell the boss about it.

Your boss: How do I know it is not a random fluctuation?

unfortunately, clicks are Bernoulli RVs, not gaussian!



Suppose  $X_1, \dots, X_n \sim \mathcal{D}$ , i.i.d., **but  $\mathcal{D}$  is unknown**. Let  $\hat{\mu} := \frac{1}{n} \sum_i X_i$ .

- In light of CLT, we can pretend that  $\mathcal{D} = \mathcal{N}(\mu, \sigma^2)$  with unknown  $\mu, \sigma^2$

- Q: Can we just use the following?

$$P\left(\hat{\mu} \in \left[\mu - \frac{z\sigma}{\sqrt{n}}, \mu + \frac{z\sigma}{\sqrt{n}}\right]\right) = 1 - 2(1 - \Phi(z))$$

population variance  
= true variance  
population mean  
= true mean

- No, we need to know the **population variance**  $\sigma^2$ .. What would you do?

- Solution: Replace  $\sigma$  with  $\hat{\sigma} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu})^2}$ , the sample variance.

- Q: Why not the 1/n version?

you want to avoid false claims as much as possible!

Suppose  $X_1, \dots, X_n \sim \mathcal{D}$ , i.i.d., but  $\mathcal{D}$  is unknown. Let  $\hat{\mu} := \frac{1}{n} \sum_i X_i$ .

**Summary**: Let  $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu})^2$

For 95% confidence:

$$\left[ \hat{\mu} - \frac{1.96\hat{\sigma}}{\sqrt{n}}, \hat{\mu} + \frac{1.96\hat{\sigma}}{\sqrt{n}} \right]$$

For 99% confidence:

$$\left[ \hat{\mu} - \frac{2.57\hat{\sigma}}{\sqrt{n}}, \hat{\mu} + \frac{2.57\hat{\sigma}}{\sqrt{n}} \right]$$

Suppose  $X_1, \dots, X_n \sim \mathcal{D}$ , i.i.d., but  $\mathcal{D}$  is unknown. Let  $\hat{\mu} := \frac{1}{n} \sum_i X_i$ .

**Summary:** Let  $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu})^2$

For 95% confidence:

$$\left[ \hat{\mu} - 1.96 \frac{\hat{\sigma}}{\sqrt{n}}, \hat{\mu} + 1.96 \frac{\hat{\sigma}}{\sqrt{n}} \right]$$

For 99% confidence:

$$\left[ \hat{\mu} - 2.57 \frac{\hat{\sigma}}{\sqrt{n}}, \hat{\mu} + 2.57 \frac{\hat{\sigma}}{\sqrt{n}} \right]$$

**Two assumptions:**

- The data follows normal distribution.
- Sample variance is equal to the actual variance.

turns out, fixable => method 2



## all quiz candidates

- For  $\hat{\sigma}^2$ , why do we use the unbiased estimator rather than the more accurate biased estimator?
- What does CLT imply about the convergence rate of the sample mean  $\bar{X}_n$  to the population mean  $\mu$ ?
- List the pros and cons of the biased variance estimator ( $1/n$ ) vs unbiased variance estimator ( $1/(n-1)$ ).

# Method 2: Gaussian (Corrected)

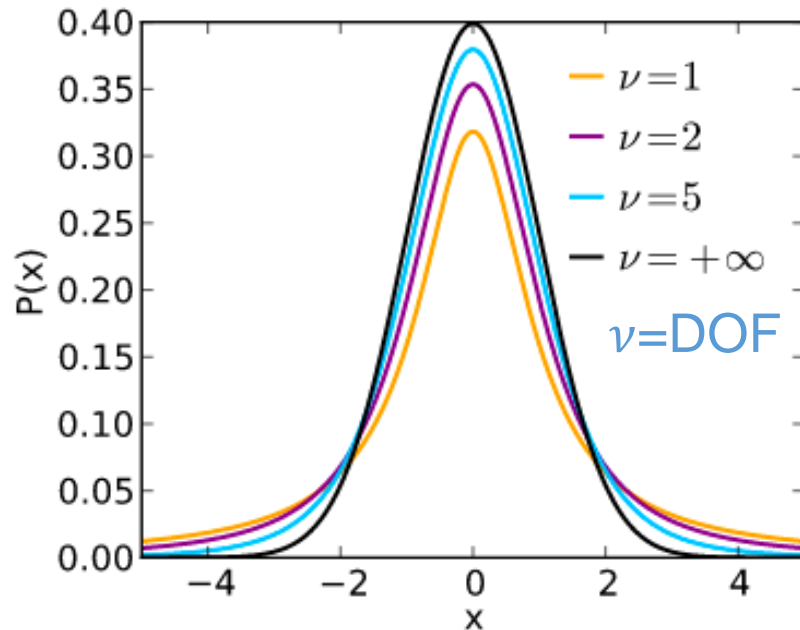
69

**Recall:** Gaussian confidence interval with  $\sqrt{n} \frac{\hat{\mu}_n - \mu}{\sigma} \sim \mathcal{N}(0,1)$ .

What if we use  $\hat{\sigma}^2$  instead of  $\sigma$ ?

(Theorem) Let  $\widehat{UVar}_n := \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu}_n)^2$ . Then,

$$\sqrt{n} \frac{\hat{\mu}_n - \mu}{\sqrt{\widehat{UVar}_n}} \sim \text{student-t}(\text{mean } 0, \text{ scale } 1, \text{ degrees of freedom } = n - 1)$$



as  $\nu \rightarrow \infty$ , it becomes Gaussian.

With a similar derivation we have done before,  
With at least 95% confidence:

$$\left[ \hat{\mu} - t_{\alpha/2, n-1} \frac{\hat{\sigma}}{\sqrt{n}}, \hat{\mu} + t_{\alpha/2, n-1} \frac{\hat{\sigma}}{\sqrt{n}} \right]$$

Where  $t_{\alpha/2, n-1}$  can be computed numerically.

**Key take away:** more conservative!  
=> more likely to be correct.

*Method 2 is strongly preferred over Method 1*

**Common practice:** Apply this method even if we do not know  
whether true distribution is Gaussian.

but often the data is highly nongaussian (e.g., skewed), which leads to..

(recall: 1.96 for gaussian)

much larger number  
compensates for the  
inaccuracy of  $\hat{\sigma}^2$

```
import scipy.stats as st  
alpha = 0.05  
st.t.ppf(1-alpha/2, df=2)  
=> 4.302652729911275
```

```
st.t.ppf(1-alpha/2, df=5)  
=> 2.5705818366147395
```

```
st.t.ppf(1-alpha/2, df=10)  
=> 2.2281388519649385
```

```
st.t.ppf(1-alpha/2, df=30)  
=> 2.0422724563012373
```

```
st.t.ppf(1-alpha/2, df=100)  
=> 1.9839715184496334
```



Computer  
Science

# **CSC380: Principles of Data Science**

## **Statistics 4**

$$\text{Poisson}(x; \lambda) = \frac{1}{x!} \lambda^x e^{-\lambda}.$$

The parameter  $\lambda$  is the *rate* parameter, and represents the expected number of arrivals  $\mathbb{E}[x] = \lambda$ . To fit the model I will need to estimate the rate parameter using some data.

a) During my last three office hours I received  $X_1 = 10, X_2 = 11, X_3 = 8$  students. Write the logarithm of the joint probability distribution  $\log p(X_1, X_2, X_3; \lambda)$ .

- Solution of HW3 a) should look like

$$C + A \log \lambda - B \lambda$$

- Let this function be  $f(\lambda)$ . Compute the derivative  $f'(\lambda)$  and set it equal to 0.
- Find the solution  $\lambda$



derivative C+ A\*log(x) -Bx for x

NATURAL LANGUAGE

MATH INPUT

EXTENDED KEYBOARD

EXAMPLES

UPLOAD

RANDOM

Suppose we observe data  $X_1, X_2, \dots, X_n \sim P(X; \theta)$ :

1. Sample new “dataset”  $X_1^*, \dots, X_n^*$  uniformly from  $X_1, \dots, X_n$  **with replacement**
2. Compute estimate  $\hat{\theta}_n(X_1^*, \dots, X_n^*)$
3. Repeat B times to get the estimators  $\hat{\theta}_{n,1}, \dots, \hat{\theta}_{n,B}$
4. Consider the **empirical distribution** of  $\left\{ \hat{\theta}_{n,b} - \frac{1}{n} \sum_{i=1}^n X_i \right\}_{b=1}^B$  and find its top  $\frac{\alpha}{2}$  quantile and bottom  $\frac{\alpha}{2}$  quantile (denoted by  $Q_U$  and  $Q_L$  respectively).
5.  $(1-\alpha)$  Confidence Interval:  $\left[ \frac{1}{n} \sum_{i=1}^n X_i - |Q_U|, \frac{1}{n} \sum_{i=1}^n X_i + |Q_L| \right]$   
counterintuitively, upper quantile for lower width, lower quantile for upper width

note: there are other variations, but this version is recommended by statisticians.

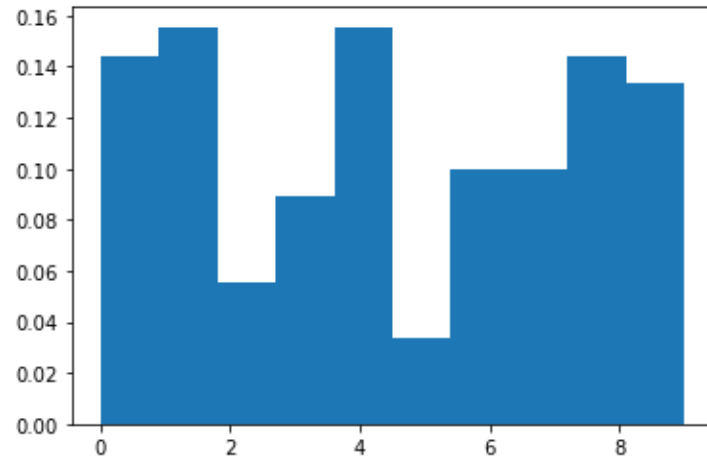
- Suppose we have  $X_1, \dots, X_n$ , i.i.d. sample from a distribution  $\mathcal{D}$
- Let  $V$  be the set of unique values of  $X_1, \dots, X_n$
- Construct a categorical distribution  $\hat{\mathcal{D}}$  where, if  $Y \sim \hat{\mathcal{D}}$ ,

$$\forall v \in V, \quad P(Y = v) = \frac{1}{n} \sum_{i=1}^n I\{X_i = v\} \quad (\text{PMF})$$

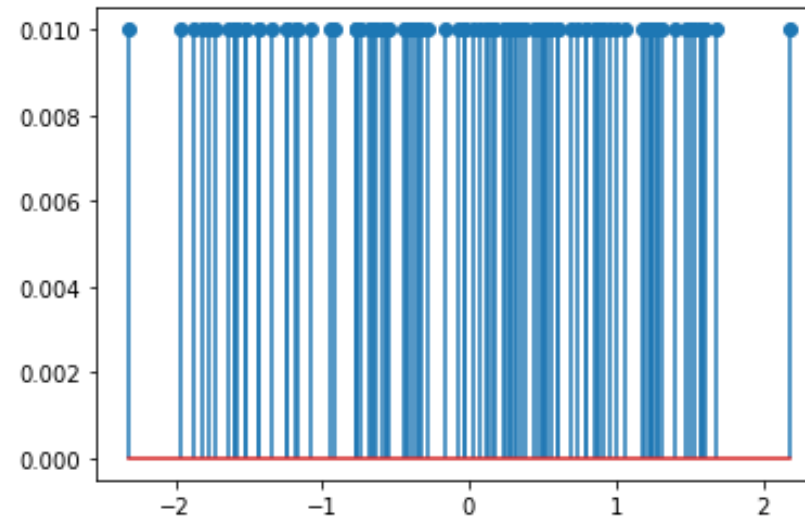
# Empirical Distribution

75

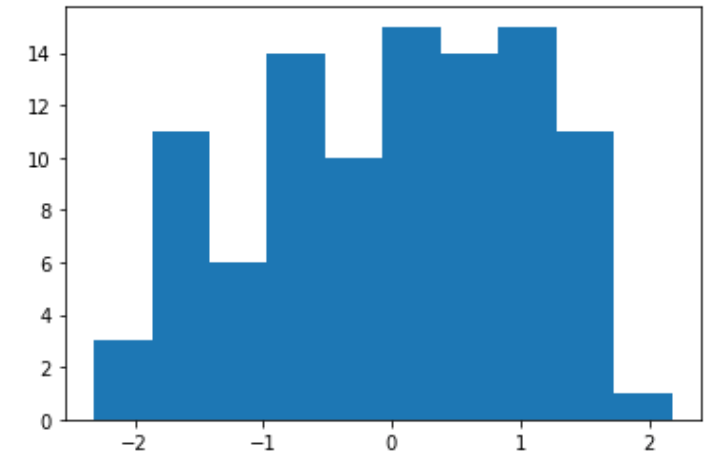
```
import numpy.random as ra
X = ra.randint(10, size=100)
plt.hist(X)
```



```
X = ra.randn(100)
plt.stem(X, 1/100*np.ones(X.shape[0]))
```



```
X = ra.randn(100)
plt.hist(X, 10)
```



not exactly what we call  
“empirical distribution”



## Pseudocode

Input:  $X_1, \dots, X_n, B, \alpha$

- Compute  $\bar{X}_n$
- Bootstrapping B times to obtain  $\{\hat{\theta}_{n,b} - \bar{X}_n\}_{b=1}^B$ ; call this array S
- Sorted S in increasing order.
- $Q_U :=$  the top  $\frac{\alpha}{2}$  quantile; i.e.,  $S[\text{int}(\text{np.ceil}((1-\alpha/2)*(B-1)))]$
- $Q_L :=$  the bottom  $\frac{\alpha}{2}$  quantile; i.e.,  $S[\text{int}(\text{np.floor}((\alpha/2)*(B-1)))]$
- Return  $[\bar{X}_n - |Q_U|, \bar{X}_n + |Q_L|]$

## Example: Generate 300 samples from Bernoulli(0.03)

```
ary = ra.rand(300) < 0.03  
muhat = np.mean(ary)
```

```
LB, UB = calc_ci_bootstrap(ary, 0.05, 10_000)  
(LB, UB)
```

```
# compute lower/upper width  
(muhat-LB, UB - muhat)
```

```
w = calc_confwidth_gaussian(ary, alpha)  
(muhat - w, muhat + w)
```

you will implement them in homework

(0.006, 0.046)

(0.023, 0.016)

(-0.107, 0.167)

asymmetric!! (note muhat=0.03..)

inherently symmetric..  
and nonsensical lower confidence bound

# Confidence Intervals Comparison

78

good = correct  
bad = incorrect

	Gaussian (corrected)	Bootstrap
small n	Bad	Bad
moderate n	Okay / Bad	Okay
large n	Good	Very Good
Computational complexity	Low	High, depends on B

Q: When could it be bad?

When the distribution is far from Gaussian

bad if the estimator  
takes a long time to  
compute

# Hypothesis testing

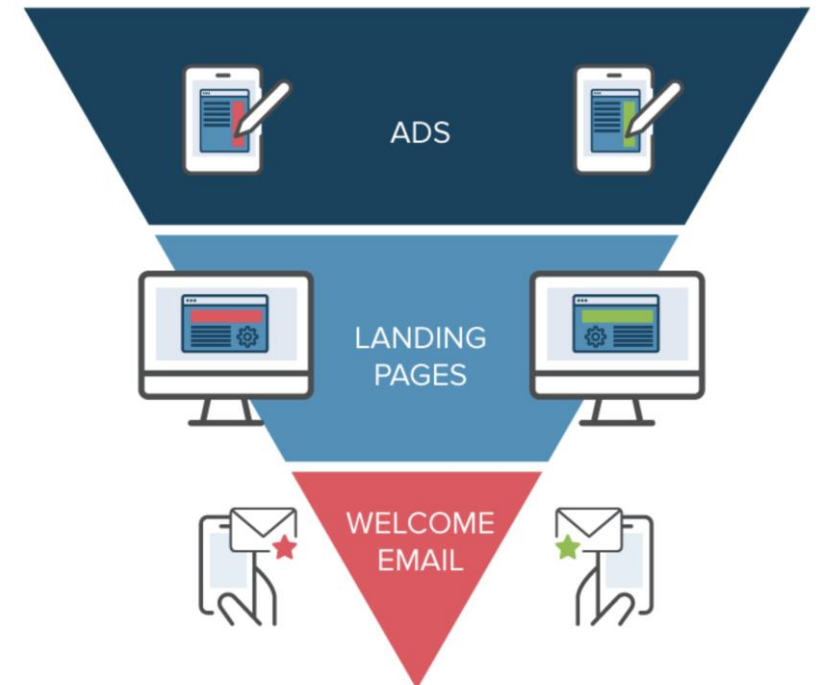
These days, webpages/apps run A/B testing extensively

Try out an alternative of webpage/interface on **randomly chosen subset of users** to gather data (and help guide the decision)

- E.g., recall the **list view** vs grid view to measure the effectiveness
- Can use **2 or more** alternatives

**Review question:** How do we know one is actually better than the other? (i.e., statistically significant)

A/B testing: Compute confidence bounds for alternatives, see if they overlap. If not, you have a clear winner!



(from optimizely.com)

Another example: Search system evaluation (e.g., Google). Compare system A vs B.

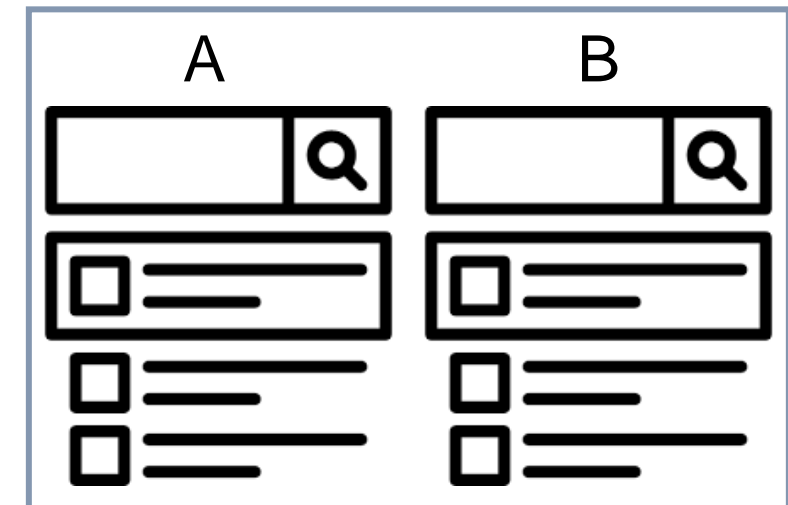
- Each evaluator
  - => a random keyword is picked, and then both systems pick top 10 relevant documents and show them.
  - => the evaluator provides rating (1-5) for both lists.

Evaluator:	1	2	3	4	5	6	...
System A	5	2	2	5	4	2	...
System B	4	1	1	4	3	1	...

Methods to claim which is better.

- 1. The average
- 2. Confidence bound
- 3. Paired t-test

bad  
okay  
great



- Two hypotheses
  - $H_0$ : the system A and B have the same performance (“null” hypothesis)
  - $H_1$ : the system A and B have different performance
- Must define a **trial**
- $n$  trials =  $n$  evaluators’ score
- Let  $\delta_i$  = score of (A) – score of (B) on the data point  $i$  (or, evaluator  $i$ )
- $\widehat{UVar}_n := \frac{1}{n-1} \sum_{i=1}^n (\delta_i - \bar{\delta}_n)^2$  where  $\bar{\delta}_n := \frac{1}{n} \sum_{i=1}^n \delta_i$

We do not know what distribution  $\delta_i$  follows.  $\Rightarrow$  again, assume normality

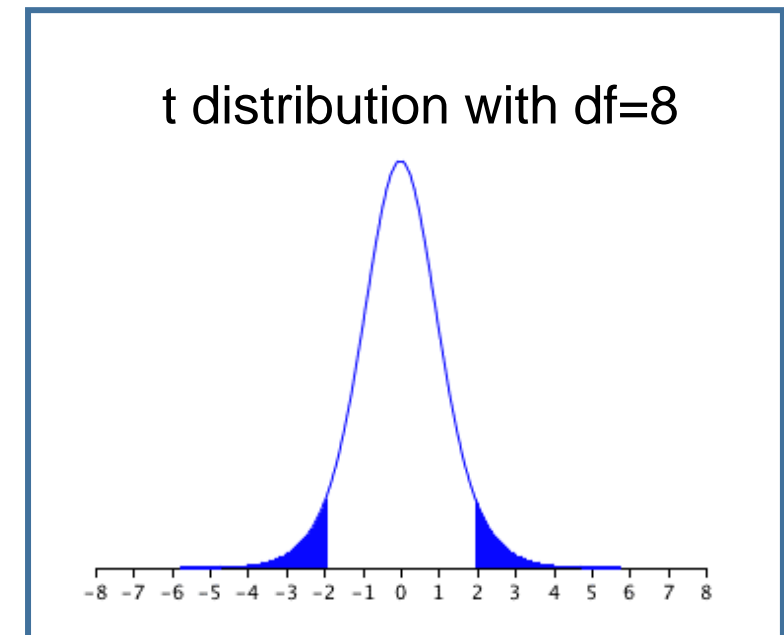
- $H_0$  assumes  $\delta_i \sim N(0, \sigma^2), i = 1, \dots, n$ .

Recall: Let  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ , and  $\hat{\mu}_n := \frac{1}{n} \sum_{i=1}^n X_i$  In our case,  $\mu = 0$

$$\left( \underset{\text{call it t-score}}{T_n} := \sqrt{n} \frac{\hat{\mu}_n - \mu}{\sqrt{\widehat{UVar}_n}} \right) \sim \text{student-t}(\text{mean } 0, \text{ scale } 1, \text{ degrees of freedom } = n - 1)$$

Ask “what is a plausible range of values of  $\hat{\mu}_n$  with significance level  $\alpha = 0.05$ ?”

- Find the *quantile* of student-t distribution! (Recall  $t_{\frac{\alpha}{2}, n-1}$ )
- Let  $X_i \leftarrow \delta_i$  &  $\mu = 0$  and see if  $|T_n|$  crosses the quantile!





Let  $X_i \leftarrow \delta_i$  and see if  $|T_n|$  crosses the quantile!

- Yes: Reject the null hypothesis  $H_0$ .  $\Rightarrow$  claim: the differences are real
- No: Accept the null hypothesis  $H_0$ .  $\Rightarrow$  claim: no statistically significant difference

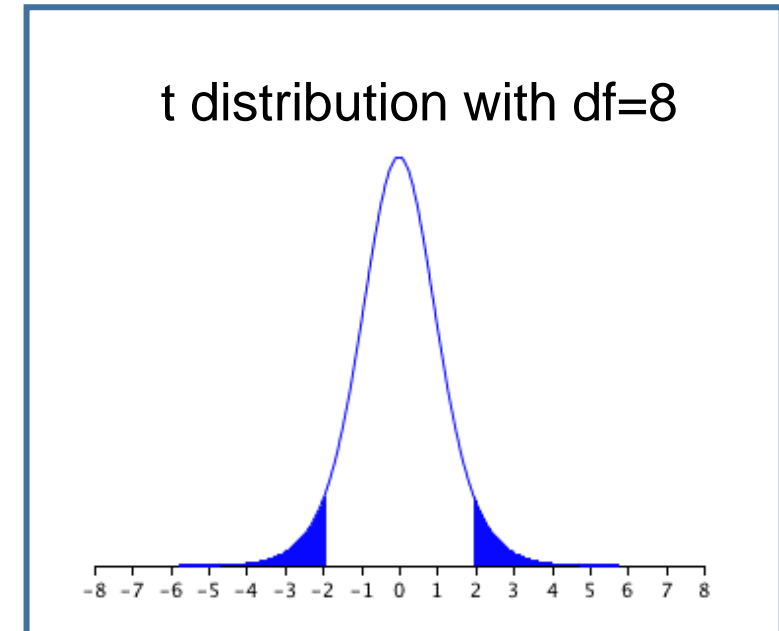
~~shut up and get more data until getting rejecting  $H_0$~~

## p-value:

- Previously, we fixed  $\alpha$  and found a binary answer.
- p-value: the smallest  $\alpha$  with which you can still reject  $H_0$ 
  - Compute: Look at the CDF of t distribution, say  $F(x)$ , and compute  $2(1 - F(|T_n|))$
- Smaller the better. Definitely want it to be below 0.05

“probability that your claim is false”

Instead of saying “our new system passed the paired t-test” people often just report the p-value. Smaller the better.



```
import numpy as np
import scipy.stats as st
```

```
data = [[5,3],
        [3,1],
        [3,1],
        [5,3],
        [4,2],
        [2,1]]
```

```
data = np.array(data)
```

```
n = data.shape[0]
alpha = 0.05
```

```
def calc_width(ary, alpha):
    ... (omitted)
```

Method 2: Gaussian (corrected)

```
for i in range(2):
    muhat = data[:,i].mean()
    width = calc_width(data[:,i],alpha)
    print([muhat-width,muhat+width])

diff = data[:,0] - data[:,1]
uvar = diff.var(ddof=1)
tscore = np.mean(diff)/uvar*np.sqrt(n)
threshold = st.t.ppf(1-alpha/2,n-1) # 'quantile'

# want tscore to be outside [-threshold,threshold]
print('tscore = %f' % tscore)
print('threshold = %f' % threshold)
```

```
output:
[2.3957369914894784, 4.937596341843855]
[0.5624036581561451, 3.1042630085105216]
tscore = 3.061862
threshold = 2.570582
```

- Summary
  - More powerful than confidence bounds!
  - Reason: It focuses on the **difference**, which is what we actually care!
- Other methods
  - Bootstrap
  - Permutation tests: highly recommended for A/B testing when the outcomes are not paired. (A/B testing in the content layout – no guarantee that the same user will see both A and B)

- **Statistical Estimation** infers unknown parameters  $\theta$  of a distribution  $p(X; \theta)$  from observed data  $X_1, \dots, X_n$
- An estimator is a function of the data  $\hat{\theta}(X_1, \dots, X_n)$ , it is a **random variable**, so it has a distribution
- **Confidence Intervals** measure uncertainty of an estimator, e.g.

$$P(\theta \in (a(X), b(X))) \geq 0.95$$

- **Bootstrap** A simple method for estimating confidence intervals

↑ Q: when is this good?

## Caution

- Confidence intervals are often misinterpreted!
- Confidence intervals in practice may not be true for small  $n$

- **Estimator bias** describes systematic error of an estimator
- **Mean squared error (MSE)** measures estimator quality / efficiency,

$$\text{MSE}(\hat{\theta}) = \mathbf{E} \left[ (\hat{\theta} - \theta)^2 \right] = \text{bias}^2(\hat{\theta}) + \mathbf{Var}(\hat{\theta})$$

- **Law of Large Numbers (LLN)** guarantees that sample mean approaches (piles up near) true mean in the limit of infinite data
- **Central Limit Theorem (CLT)** says sample mean approaches a Normal distribution with enough data. Also means  $\frac{1}{\sqrt{n}}$  convergence.
- **LLN** and **CLT** are *asymptotic statements* and do not hold for small/medium data in general

- UNUSED


$$X \sim \text{Bernoulli}(\pi)$$

- What about Variance? Recall  $\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}X)^2]$   
 $\text{Var}(X) = \pi(1 - \pi)$

quiz candidate:  
compute  $\text{Var}(X)$

Q: was there another way to calculate variance?

- Suppose we observe N coin flips  $x_1, \dots, x_N$ . Estimate  $\text{Var}(X)$  as

$$\widehat{\text{Var}} = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\pi})^2$$


we are using the sample mean here!

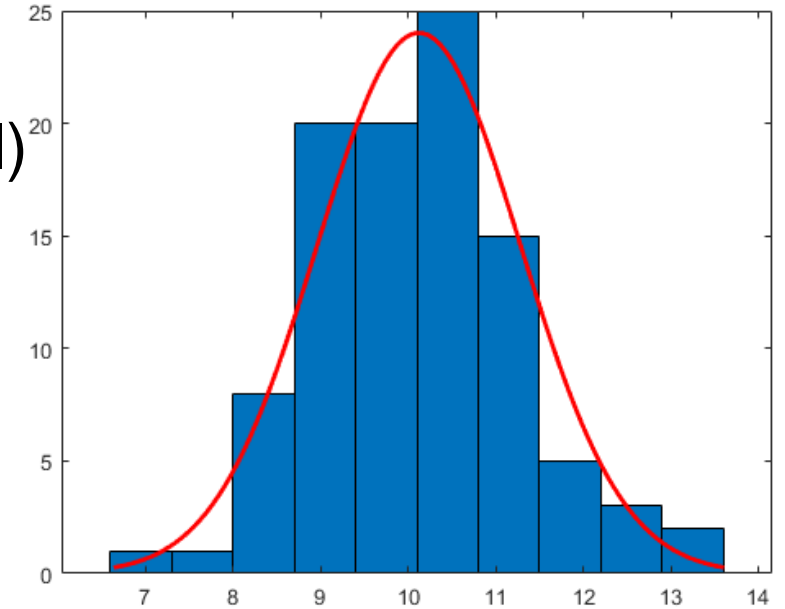
Suppose we observe the heights of  $N$  student at UA, and we model them as Gaussian:  
(A.K.A. Normal)

$$\{x_i\}_i^N \sim \mathcal{N}(\mu, \sigma^2)$$

How can we estimate the **mean**?

$$\hat{\mu} = \frac{1}{N} \sum_i x_i \approx \mu$$

Sample mean  
 $\bar{x}$



How can we estimate the **variance**?

$$\hat{\sigma}^2 = \frac{1}{N} \sum_i (x_i - \hat{\mu})^2 \approx \sigma^2$$





Computer  
Science

# **CSC380: Principles of Data Science**

## **Statistics 2**

Suppose we observe the heights of  $N$  student at UA, and we model them as Gaussian:

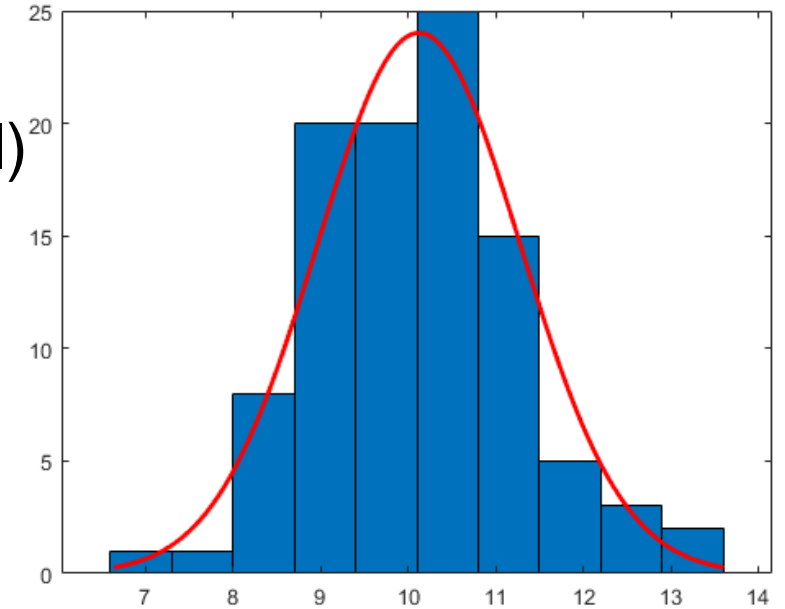
(A.K.A. Normal)

$$\{x_i\}_i^N \sim \mathcal{N}(\mu, \sigma^2)$$

How can we estimate  $\mu$ ?

$$\hat{\mu} = \frac{1}{N} \sum_i x_i \approx \mu$$

Sample mean  
 $\bar{x}$



How can we estimate  $\sigma$ ?

$$\sigma^2 = \text{var}(X) = \text{E}[(X - \mu)^2] \approx \frac{1}{n} \sum_i (x_i - \mu)^2 \approx \frac{1}{n} \sum_i (x_i - \hat{\mu})^2$$

- **Consistency (asymptotic notion)** Given enough data, the estimator *converges* to the true parameter value

$$\lim_{n \rightarrow \infty} \hat{\theta}(x_1, \dots, x_n) \rightarrow \theta$$

This convergence can be measured in a number of ways: in probability, in distribution, absolutely

A bare minimum requirement!

Otherwise, you may collect more data that will give us a worse estimator!

- **Efficiency (nonasymptotic notion)** It should have low error with finite  $n$ , e.g.

$$\text{MSE}(\hat{\theta}_n) = \mathbf{E}[(\hat{\theta}_n - \theta)^2]$$

Mean squared error should be small

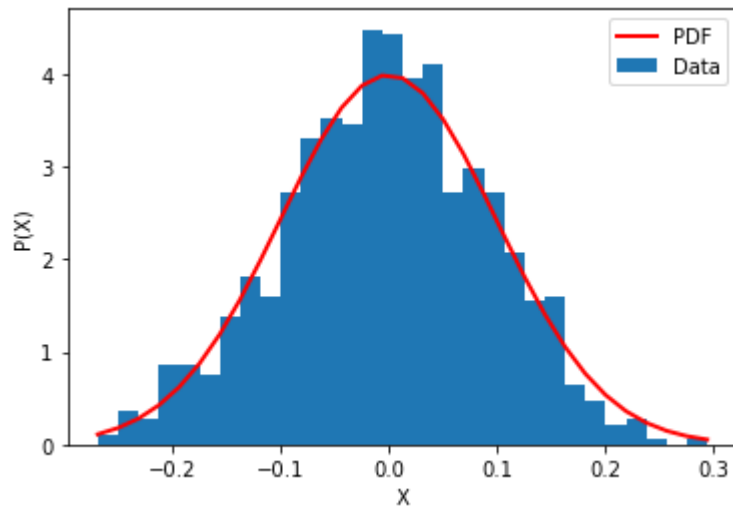
looks like variance but it's different!

Q: spot the difference from  $\text{Var}(\hat{\theta}_n)$ ?

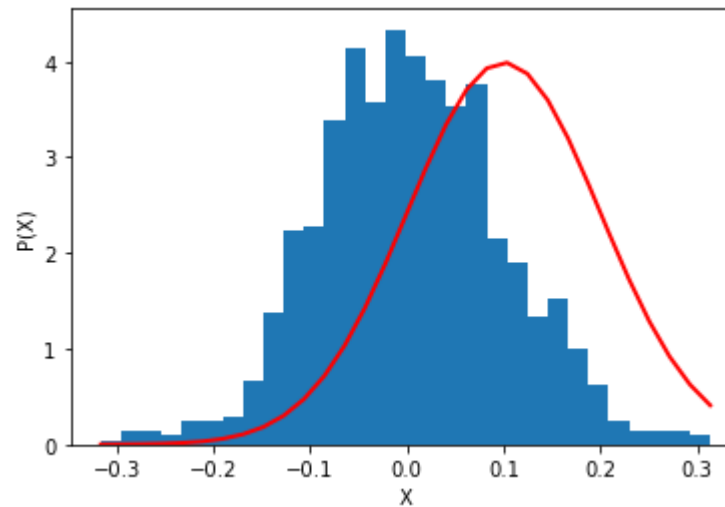
*Suppose we observe  $N$  data points from a Gaussian model and wish to estimate model parameters.*

*Say we only need to choose from the following three Gaussians...*

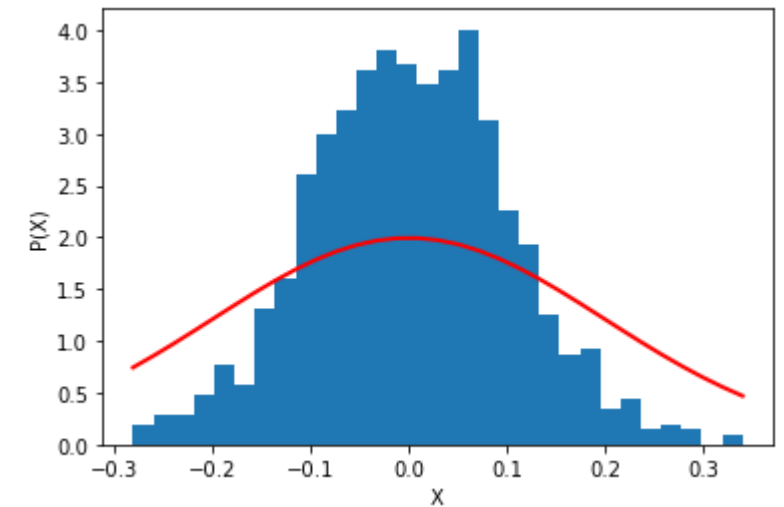
**High  
Likelihood**



**Low  
Likelihood (mean)**



**Low  
Likelihood (variance)**



***Likelihood Principle:*** *Given a statistical model, the likelihood function describes all evidence of a parameter that is contained in the data.*

[ Source: Wasserman, L. 2004 ]

**Example**  $N$  coin tosses with  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ . We don't know the coin bias  $p$ . The likelihood function is,

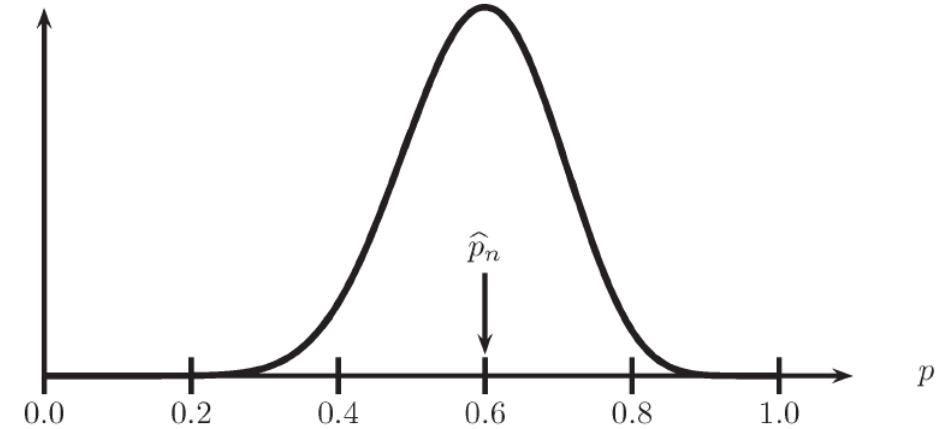
$$X_1, \dots, X_n \sim \text{Bernoulli}(p)$$
$$\mathcal{L}_n(p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^S (1-p)^{n-S}$$

where  $S = \sum_i x_i$ . The log-likelihood is,

$$\log \mathcal{L}_n(p) = S \log p + (n - S) \log(1 - p)$$

Set the derivative of  $\log \mathcal{L}_n(p)$  to zero and solve,

$$\hat{p}^{\text{MLE}} = S/n = \frac{1}{n} \sum_{i=1}^n x_i$$



*Likelihood function for Bernoulli with  $n=20$  and  $\sum_i x_i = 12$  heads*

Maximum likelihood is equivalent to sample mean in Bernoulli

⇒ this showcases how MLE is aligned to our intuition!

- HW2 grades will be out by next Friday.
- HW3 will be released next Tuesday

- Discussion this week: Upload questions by Friday:
  - Your name could appear again if your question is not qualified.
- Questions not on the lecture itself will not count.
  - Do not count: logistics, homework, etc.
  - Next time, you will not be given a second chance.

*Claim: sample mean is a consistent estimator of the true mean.*

We now know the **sample mean** is an unbiased estimator, namely:  $E[\hat{\mu}_n] = \mu, \forall n$

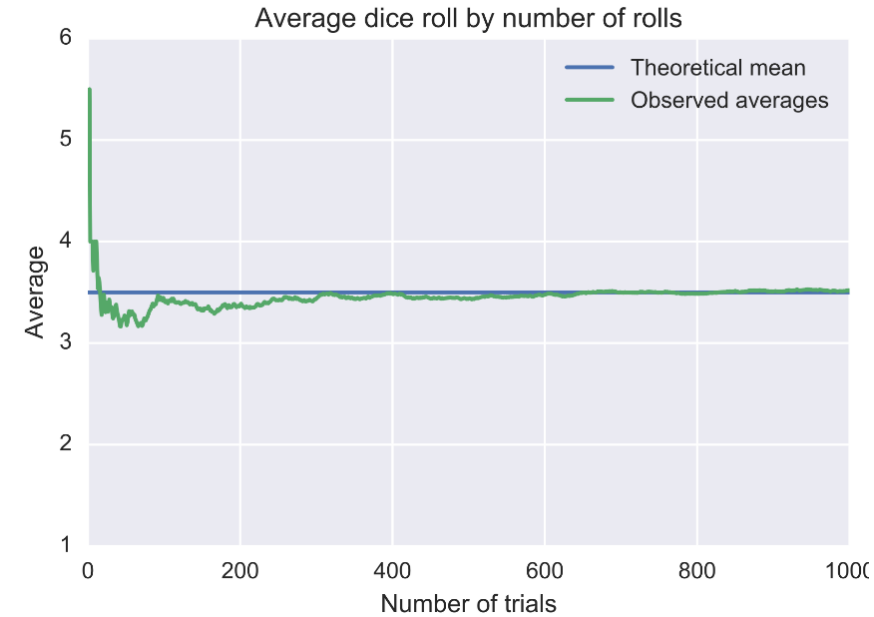
But, expected value is usually different from the actual realization of  $\hat{\mu}_n$ . Will  $\hat{\mu}_n$  become the true mean?

**Yes, in the limit.**

(Theorem)  $\lim_{N \rightarrow \infty} \hat{\mu}_n = \mu$

This is the **law of large numbers**

- Weak Law: Converges to mean with high probability
- Strong Law: Stronger notion of convergence; will converge at all times! (if variance is finite)



Limitation: it does not say how fast it will converge!



Let  $X_1, \dots, X_N$  be iid with mean  $\mu$  and variance  $\sigma^2$  then the sample mean  $\bar{X}_N$  approaches a Normal distribution

$$\lim_{N \rightarrow \infty} \bar{X}_N \rightarrow \mathcal{N} \left( \mu, \frac{\sigma^2}{N} \right)$$

=> the convergence rate is  $\frac{\sigma}{\sqrt{N}}$ !!

Actually, a mathematically rigorous version is

$$\lim_{N \rightarrow \infty} \frac{\sqrt{N}}{\sigma} (\bar{X}_N - \mu) \rightarrow \mathcal{N}(0, 1)$$

## Comments

- LLN says estimates  $\bar{X}_N$  “pile up” near true mean, CLT says *how* they pile up
- Pretty remarkable since we make **no assumption about how  $X_i$  are distributed**
- Variance of  $X_i$  **must be finite**, i.e.  $\sigma^2 < \infty$  (e.g., Cauchy distribution has  $\sigma^2 = \infty$ )

Suppose  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  with unknown  $\mu$  & known  $\sigma^2$ . Let  $\hat{\mu} := \frac{1}{n} \sum_i X_i$ .

**(Fact 1)**  $\hat{\mu} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$

quiz candidate

## Recall:

- Closed under additivity:

$$X \sim \mathcal{N}(\mu_x, \sigma_x^2) \quad Y \sim \mathcal{N}(\mu_y, \sigma_y^2)$$

$$X + Y \sim \mathcal{N}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

- Closed under affine transformation (a and b constant):

$$aX + b \sim \mathcal{N}(a\mu_x + b, a^2\sigma_x^2)$$

(proof)

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

Use this with  $X = \sum_{i=1}^n X_i$ ,  $a = \frac{1}{n}$ ,  $b = 0$ .

- HW3 is out.
- HW2 solution will be posted tonight.
- Midterm: October 11 (Tuesday)
  - homework problems
  - 'quiz candidates'

- Discussion for Week 6
- Ask questions by Thursday 11:59pm