



Computer  
Science

# CSC380: Principles of Data Science

**Statistics 5 & Midterm review**

**Xinchen Yu**

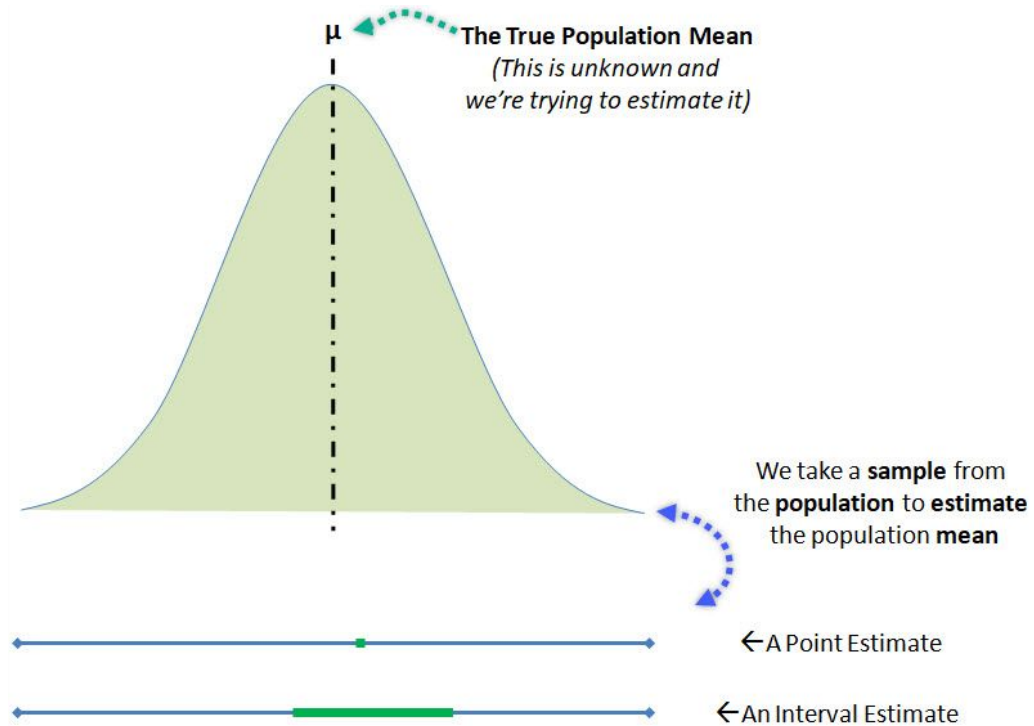
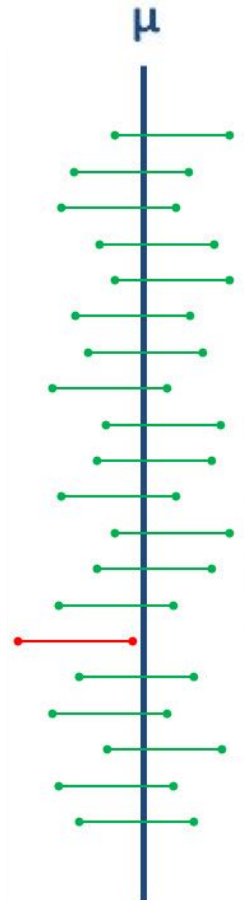
- HW5 has been out.
  - Due Friday, Mar 15
- Practice problems will be out by end of this weekend.
  - Solutions will be out by Mar 10.
- Lecture on Tuesday Mar 12:
  - Another review session
  - revisit solutions of some questions in HW1 - 4
  - Q & A

- Midterm
  - What you can bring?
    - Cheat sheet: letter size, double-sided
    - Scientific calculator
  - Time: Mar 14, Thursday, 3:30-4:45 pm
  - Location: C E Chavez Bldg, Rm 111 (same as lecture room)

$\frac{\sqrt{2}}{3}$  is ok

# Review: Interval estimate

4



# Review: Gaussian (Corrected)

Suppose  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  with unknown  $\mu$  & known  $\sigma^2$ .

**(Fact 1)**  $\hat{\mu} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$   $\sqrt{n} \frac{\hat{\mu} - \mu}{\sigma} \sim \mathcal{N}(0, 1)$

**(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

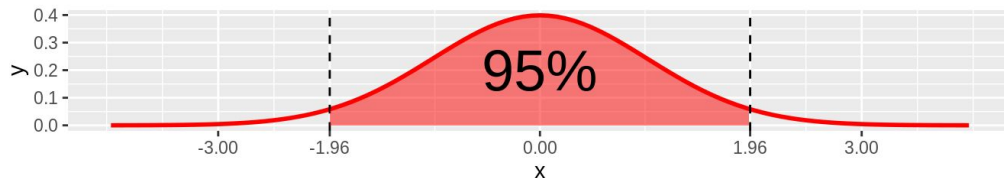
$z = 2.58$ : RHS  $\approx .99$ ,

**Let:**  $Z \rightarrow \sqrt{n} \frac{\hat{\mu} - \mu}{\sigma}$

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

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

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



$$\begin{aligned}\phi(1.96) &= P(Z \in [-\infty, 1.96]) = 0.975 \\ P(Z \in [-1.96, 1.96]) &= 1 - 2 \cdot (1 - 0.975) = 0.95\end{aligned}$$

# Review: Gaussian (Corrected)

Suppose  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  with unknown  $\mu$  & known  $\sigma^2$ .

**(Fact 1)**  $\hat{\mu} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$   $\sqrt{n} \frac{\hat{\mu} - \mu}{\sigma} \sim N(0,1)$  → T-dist

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

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

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

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

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

**Let:**  $Z \rightarrow \sqrt{n} \frac{\hat{\mu} - \mu}{\sigma}$

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

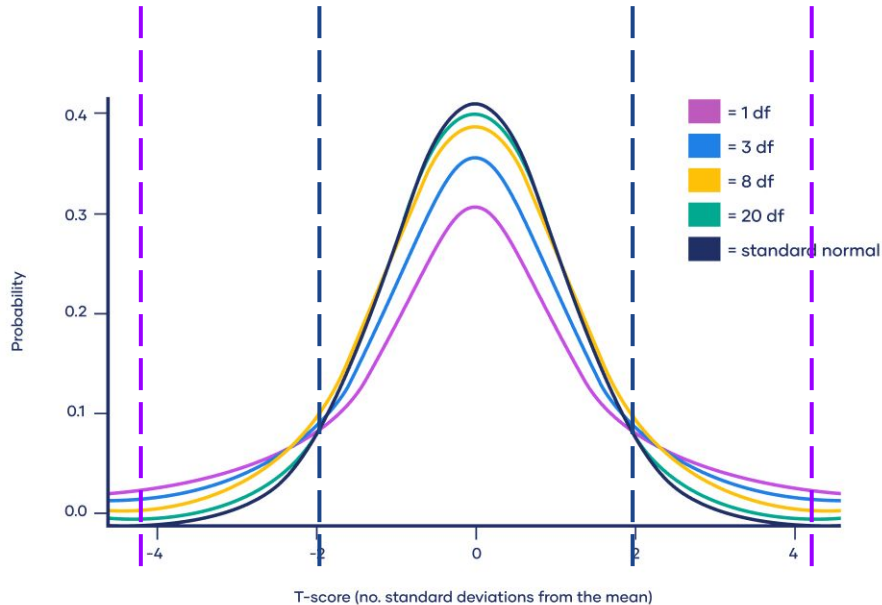
$$P\left(\hat{\mu} \in \left[\mu - \frac{2.58\sigma}{\sqrt{n}}, \mu + \frac{2.58\sigma}{\sqrt{n}}\right]\right) \geq 0.99$$

Q: what if  $X$  from an arbitrary distribution (e.g. uniform)?

Q: what if  $\sigma^2$  is unknown and sample size is small ( $< 30$ )?

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

# Z score versus T score



When alpha is 0.05:

- For standard normal distribution:

$$P(X \in [-1.96, 1.96]) = 0.95$$

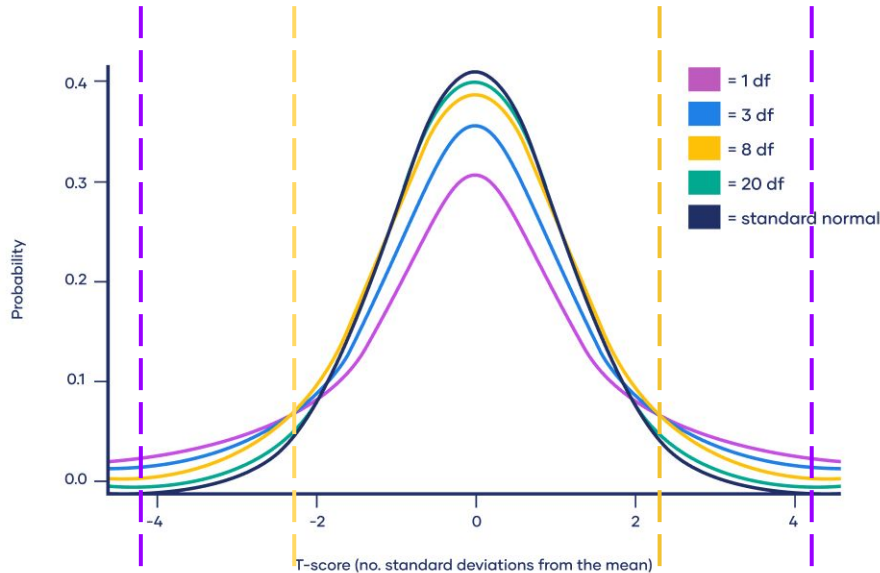
- For T distribution when  $n = 2$ :

$$P(X \in [-4.30, 4.30]) = 0.95$$

# Review: T scores for different df

Let's compare t scores when we only have 2 and 6 observations in the sample:

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



(recall: 1.96 for gaussian)

```
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
```



# Method 2: Bootstrap

Suppose  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  with unknown  $\mu$  & known  $\sigma^2$ .

**(Fact 1)**  $\hat{\mu} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$   $\sqrt{n} \frac{\hat{\mu} - \mu}{\sigma} \sim \mathcal{N}(0, 1)$

**(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  $\approx .99$ ,

**Let:**  $Z \longrightarrow \sqrt{n} \frac{\hat{\mu} - \mu}{\sigma}$

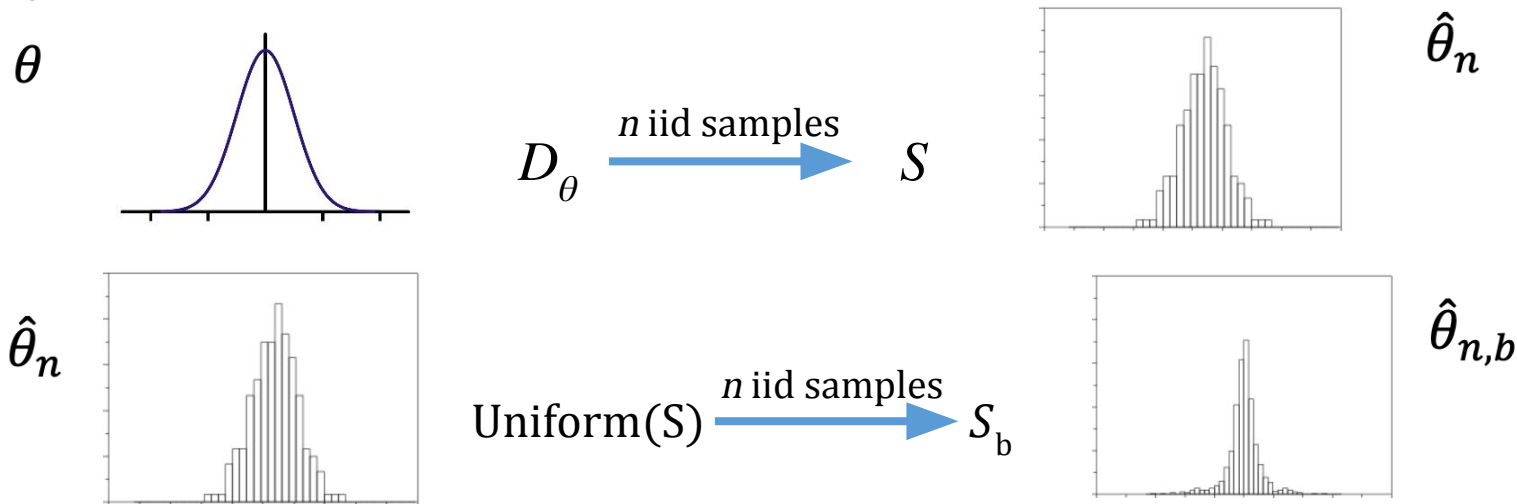
$$P\left(\hat{\mu} \in \left[\mu - \frac{1.96\sigma}{\sqrt{n}}, \mu + \frac{1.96\sigma}{\sqrt{n}}\right]\right) \geq 0.95$$

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

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

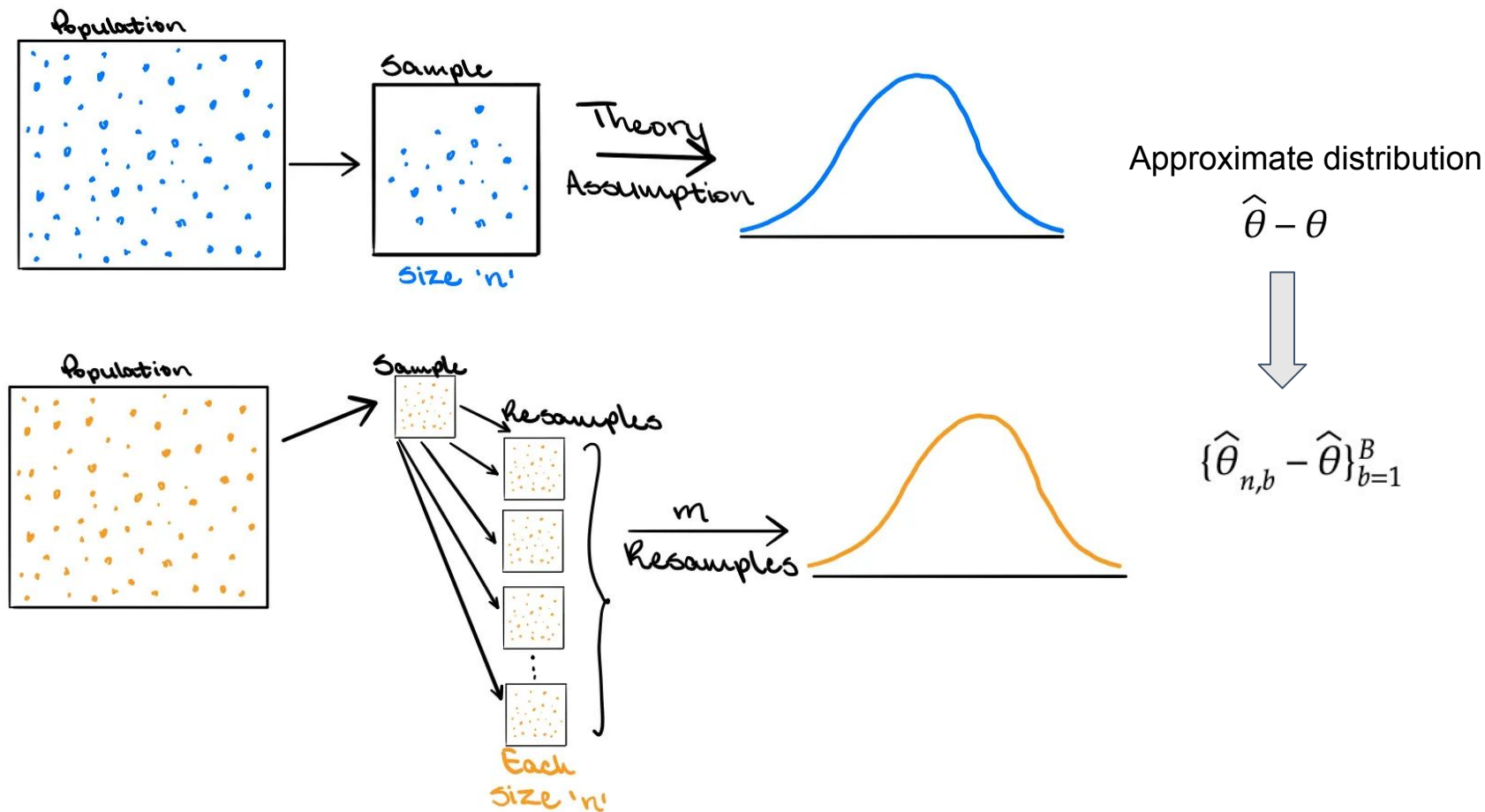
Directly approximate distributions of  $\hat{\mu} - \mu$

- Key idea: approximate  $\nu$ , the distribution of  $\hat{\theta}_n - \theta$
- Insight:



- Use empirical distribution of  $\hat{\theta}_{n,b} - \hat{\theta}_n$ 's to approximate  $\nu$ , obtaining approximations of  $v_{\alpha/2}$  and  $v_{1-\alpha/2}$
- This empirical distribution can be obtained by drawing multiple  $S_b$ 's (bootstrap subsample)

# Method 2: Bootstrap



# Method 2: Bootstrap example

Sample data: 30, 37, 36, 43, 42, 43, 43, 46, 41, 42

Sample mean:  $\bar{x} = 40.3$

We want to know the distribution of:  $\delta = \bar{x} - \mu$

Can approximate the distribution:  $\delta^* = \bar{x}^* - \bar{x}$

Let's resample data with same size and generate 20 bootstrap samples:

43	36	46	30	43	43	43	37	42	42	43	37	36	42	43	43	42	43	42	43
43	41	37	37	43	43	46	36	41	43	43	42	41	43	46	36	43	43	43	42
42	43	37	43	46	37	36	41	36	43	41	36	37	30	46	46	42	36	36	43
37	42	43	41	41	42	36	42	42	43	42	43	41	43	36	43	43	41	42	46
42	36	43	43	42	37	42	42	42	46	30	43	36	43	43	42	37	36	42	30
36	36	42	42	36	36	43	41	30	42	37	43	41	41	43	43	42	46	43	37
43	37	41	43	41	42	43	46	46	36	43	42	43	30	41	46	43	46	30	43
41	42	30	42	37	43	43	42	43	43	46	43	30	42	30	42	30	43	43	42
46	42	42	43	41	42	30	37	30	42	43	42	43	37	37	37	42	43	43	46
42	43	43	41	42	36	43	30	37	43	42	43	41	36	37	41	43	42	43	43

## Method 2: Bootstrap example

43	36	46	30	43	43	43	37	42	42	43	37	36	42	43	43	42	43	42	43
43	41	37	37	43	43	46	36	41	43	43	42	41	43	46	36	43	43	43	42
42	43	37	43	46	37	36	41	36	43	41	36	37	30	46	46	42	36	36	43
37	42	43	41	41	42	36	42	42	43	42	43	41	43	36	43	43	41	42	46
42	36	43	43	42	37	42	42	42	46	30	43	36	43	43	42	37	36	42	30
36	36	42	42	36	36	43	41	30	42	37	43	41	41	43	43	42	46	43	37
43	37	41	43	41	42	43	46	46	36	43	42	43	30	41	46	43	46	30	43
41	42	30	42	37	43	43	42	43	43	46	43	30	42	30	42	30	43	43	42
46	42	42	43	41	42	30	37	30	42	43	42	43	37	37	37	42	43	43	46
42	43	43	41	42	36	43	30	37	43	42	43	41	36	37	41	43	42	43	43

Calculate sample mean for each column (bootstrap sample), compute:  $\delta^* = \bar{x}^* - \bar{x}$

Sort the 20 differences:

-1.6, -1.4, -1.4, -0.9, -0.5, -0.2, -0.1, 0.1, 0.2, 0.2, 0.4, 0.4, 0.7, 0.9, 1.1, 1.2, 1.2, 1.6, 1.6, 2.0

If confidence level is 80%, find out top 10% and bottom 10%:

-1.6, -1.4, -1.4, -0.9, -0.5, -0.2, -0.1, 0.1, 0.2, 0.2, 0.4, 0.4, 0.7, 0.9, 1.1, 1.2, 1.2, 1.6, 1.6, 2.0

The bootstrap confidence interval is:

$$[\bar{x} - \delta_{.1}^*, \bar{x} - \delta_{.9}^*] = [40.3 - 1.6, 40.3 + 1.4] = [38.7, 41.7]$$

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**

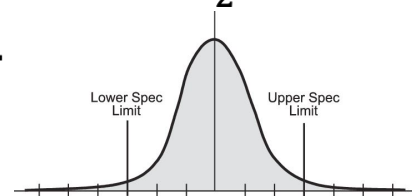
2. Compute estimate  $\hat{\theta}_n(X_1^*, \dots, X_n^*)$

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. Why?

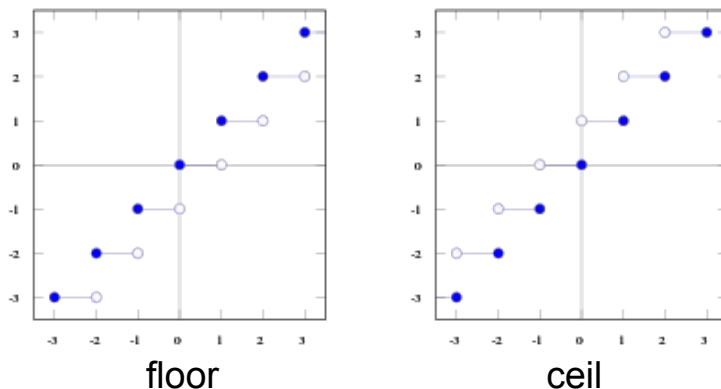


$$P\left(v_{\frac{\alpha}{2}} \leq \hat{\theta}_n - \theta \leq v_{1-\frac{\alpha}{2}}\right) \geq 1 - \alpha$$

## 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  $[\bar{X}_n - |Q_U|, \bar{X}_n + |Q_L|]$



## Midterm Review



- Prioritize reviewing basic concepts & ideas
- Understand the motivations and links between concepts
- “Memorization with understanding”
  
- Try to solve these on your own, then discuss with classmates
  - examples in the slides
  - HW questions (esp. if you did not get them right the first time)
  - practice problems

- What will not included in the midterm?
  - Code related questions
  - Pure proof questions
    - But may need you to provide justifications

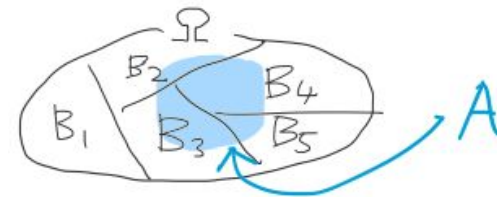
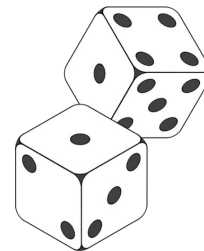
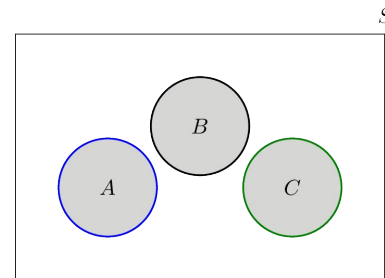
# Probability

- Basic definitions: outcome space, events
- Probability  $P$ : maps events to  $[0, 1]$  values
  - Three axioms
  - Axiom 3: additivity
- Special case of  $P$ : each outcomes is equally likely

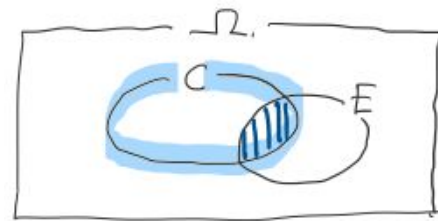
$$P(E) = \frac{|E|}{|\Omega|}$$

Number of elements  
in event set
Number of possible  
outcomes (36)

- distributive law, inclusion-exclusion rule; law of total probability



- $P(E \cap C) = P(E|C)P(C) = P(C|E)P(E)$

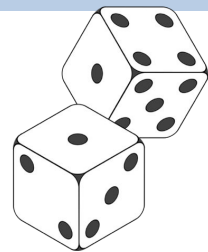


- Conditional probability
  - Chain rule, chain rule + law of total probability, bayes rule
  - Important application: medical diagnosis
  - Approach: write down the joint probability table

- Independence of events:

$$P(A, B) = P(A)P(B)$$

- Conditional / joint / marginal probability



- Discrete random variable  $X$  (e.g., sum of two dice)
- Representation of its distribution: probability mass function (PMF)
  - Tabular representation of joint distribution of 2 RVs  $(X,Y)$
- RVs: law of total probability, conditional probability, chain rule, bayes rule, independence, conditional independence
- Useful discrete distributions
  - Uniform
  - Bernoulli
  - Binominal

- Continuous random variable  $X$ :  $P(X = x) = 0$  for any  $x$

- Probability density function (PDF)

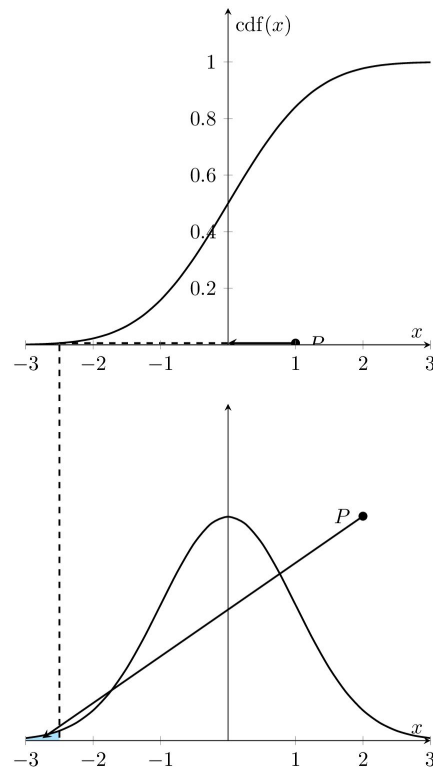
$$P(a < X \leq b) = \int_a^b p(x)dx \quad p(x) = \frac{dF(x)}{dx}$$

- Cumulative distribution function (CDF)

$$P(a < X \leq b) = F(b) - F(a)$$

- Useful continuous distributions

- Uniform
- Gaussian (important properties)



- Moments of random variables: expectation, variance, covariance
- Calculate mean (expectation) and variance of RVs
  - Linearity of expectation:  $E[X + cY] = E[X] + cE[Y]$  for constant  $c$
  - $E[X^2]$
  - $E[XY]$ 
    - If independent:  $E[X]E[Y]$
    - If not independent:  $E[XY] = \sum_{(x,y)} xy \cdot p(x, y)$
  - $E[X \mid Y = y]$
  - $\text{Var}[c] = 0$
  - $\text{Var}[cX]$
  - $\text{Var}[X + c] = \text{Var}[X]$
  - $\text{Var}[X+Y]$  when independent

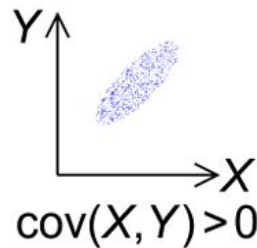
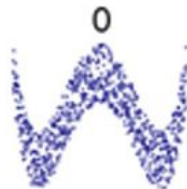
$$\text{Var}[X] = E[(X - E[X])^2]$$

- Expectation and variance of useful distributions (esp. Bernoulli, Gaussian)



- Measures *linear relationship* between  $X, Y$

$$\text{Cov}(X, Y) = 0 \not\Rightarrow X \perp Y$$



- Pearson correlation:  $\rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$ , where  $\sigma_X = \sqrt{\text{Var}(X)}$

- Important property:  $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$ 
  - What if  $X, Y$  are independent?

Let random variable  $X$  and  $Y$  independent from each other.

The PMF for  $X$  is:  $P(X = 0) = 0.5, P(X = 1) = 0.5$

The PMF for  $Y$  is:  $P(Y = 0) = 0.25, P(Y = 1) = 0.5, P(Y = 2) = 0.25$

$E[XY^2]$ ?

$$P(XY^2 = 1) = P(X = 1, Y = 1) = 0.5 \cdot 0.5 = 0.25$$

$$P(XY^2 = 4) = P(X = 1, Y = 2) = 0.5 \cdot 0.25 = 0.125$$

$$P(XY^2 = 0) = 1 - 0.25 - 0.125 = 0.625$$

$$\begin{aligned} P(XY^2 = 0) &= P(X = 0, Y = 0) + P(X = 0, Y = 1) + P(X = 0, Y = 2) + P(X = 1, Y = 0) \\ &= 0.125 + 0.25 + 0.125 + 0.125 = 0.625 \end{aligned}$$

$$E[XY^2] = 1 \cdot 0.25 + 4 \cdot 0.125 + 0 \cdot 0.625 = 0.75$$

Let random variable  $X$  and  $Y$  independent from each other.

The PMF for  $X$  is:  $P(X = 0) = 0.5, P(X = 1) = 0.5$

The PMF for  $Y$  is:  $P(Y = 0) = 0.25, P(Y = 1) = 0.5, P(Y = 2) = 0.25$

$E[XY^2]$ ?

$$E[XY^2] = E[X] \cdot E[Y^2]$$

$$E[X] = 0.5$$

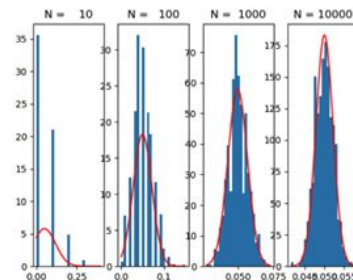
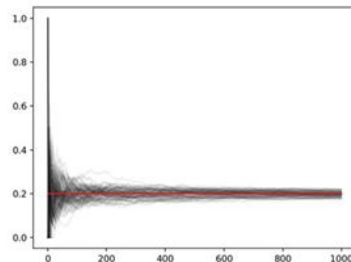
$$E[Y^2] = 0.25 \cdot 0^2 + 0.5 \cdot 1^2 + 0.25 \cdot 2^2 = 1.5$$

$$E[XY^2] = E[X] \cdot E[Y^2] = 0.5 \cdot 1.5 = 0.75$$

# Statistics

# Statistics

- Statistics: make statements about data generation process based on data seen; reverse engineering
- Point estimation
  - Given iid samples  $X_1, \dots, X_n \sim \mathcal{D}_\theta$ , estimate  $\theta$  by constructing *statistics*  $\hat{\theta}_n$
  - Basic estimators: sample mean, sample variance
  - Performance measures: unbiasedness, consistency, MSE (efficiency)
  - Bias-variance decomposition:
    - $$\text{MSE}(\hat{\theta}) = \text{bias}(\hat{\theta})^2 + \text{var}(\hat{\theta})$$
- Useful probability tools:
  - Law of Large Numbers
  - Central Limit Theorem



# Statistics

- Sample mean, sample variance
- Sample variance
  - biased version
  - unbiased version
  - how to determine an estimator is biased or unbiased?
- MSE, Bias, Variance
  - how to calculate expectation and variance if there are more than 1 random variable -- use what we learned in probability lecture 5 & 6

# Statistics

- Calculate bias and variance

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

## Important properties of Gaussian

- Closed under additivity:

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

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

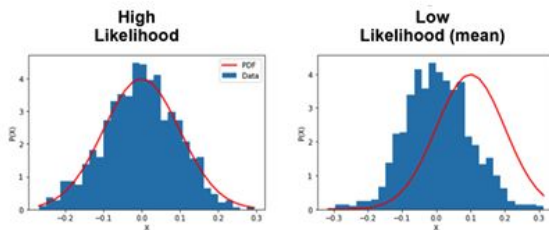
- Closed under affine transformation (a and b constant):

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

# Statistics

- Maximum likelihood (MLE): a general approach for point estimation
- Given  $X_1, \dots, X_n \sim \mathcal{D}_{\theta^*}$ , estimate  $\theta^*$  by finding the maximizer of the likelihood function

$$\mathcal{L}_n(\theta) = p(x_1, \dots, x_n; \theta) = p(x_1; \theta) \cdot \dots \cdot p(x_n; \theta)$$



- Intuition:  $\mathcal{L}_n(\theta)$  measures the “goodness of fit” of  $\mathcal{D}_{\theta}$  to data  $x_1, \dots, x_n$
- $\mathcal{D}_{\theta}$  can be general, e.g. Bernoulli, Gaussian, Poisson (in HW3)



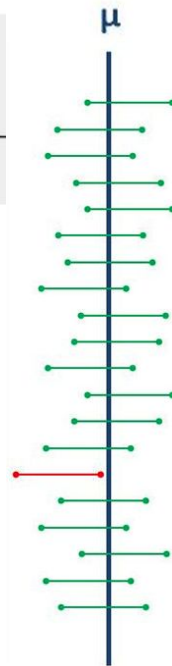
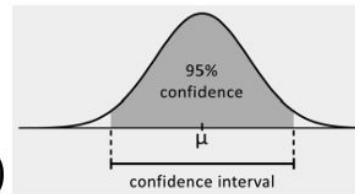
# Statistics

- Confidence interval (interval estimation)

- Definition of confidence intervals:

- Given data  $X_1, \dots, X_n \sim \mathcal{D}_\theta$  with unknown  $\theta$  (say,  $\mathcal{D}_\theta = \mathcal{N}(\theta, 1)$ )
- Construct  $a_n, b_n$  (that depends on  $X_1, \dots, X_n$ ), such that
$$P(\theta \in [a_n, b_n]) \geq 1 - \alpha$$

- Interpretation: unless we are extremely unlucky (in that we encounter an unrepresentative dataset, which happens with prob.  $\leq \alpha$ ), our confidence interval always contains the underlying parameter



# Statistics

- Confidence intervals for population mean:

- Gaussian(naive):

$$\left[ \hat{\mu} - \frac{z_{1-\alpha/2} \hat{\sigma}}{\sqrt{n}}, \hat{\mu} + \frac{z_{1-\alpha/2} \hat{\sigma}}{\sqrt{n}} \right], z_{1-\alpha/2} = 1 - \alpha/2\text{-quantile of } \mathcal{N}(0,1)$$

- Gaussian(corrected):

$$\left[ \hat{\mu} - \frac{t_{1-\alpha/2} \hat{\sigma}}{\sqrt{n}}, \hat{\mu} + \frac{t_{1-\alpha/2} \hat{\sigma}}{\sqrt{n}} \right], t_{1-\alpha/2} = 1 - \alpha/2\text{-quantile of } t \text{ distribution (degree of freedom=?)}$$

- We expect you to be able to compute them on a small dataset

- Confidence intervals for general population parameters: bootstrap

