Elements Of Data Science - F2025

# Week 4: Hypothesis Testing

9/23/2025

# TODOs

- Readings
  - PDSH Chap 5: What is Machine Learning and Introduction to Scikit-Learn
    - PDSH Chap 5 In Depth: Linear Regression
    - PDSH Chap 5 In Depth: Decision Trees and Random Forests
  - Recommended PML Chap 3
  - Optional PML Chap 2
  - Optional PDSH Chap 5 In Depth: Support Vector Machines

- **Quiz 3**: due today, Sep 23nd, 11:59pm ET via Gradescope

- **Quiz 4**: due Tue Sep 30th, 11:59pm ET via Gradescope

- **HW1**: due Tue, October 14th at 11:59 pm EST via Gradescope

# Additional Resources

- Statistical Rules of Thumb, Gerald van Belle Chapter 2 online

- On the use of p-values

  - The ASA's Statement on p-Values: Context, Process, and Purpose

  - Moving to a World Beyond "p < 0.05"

  - "The 2019 ASA Guide to P-values and Statistical Significance: Don't Say What You Don't Mean" (Some Recommendations)(ii)

# Today

- Confidence Intervals
- Hypothesis Testing
- Multi-Armed Bandit (MAB)

# Questions?

# Environment Setup

# Environment Setup

```python
In [1]:   import numpy as np
          import pandas as pd
          import matplotlib.pyplot as plt
          import seaborn as sns

          sns.set_style('darkgrid')

          %matplotlib inline
```

# Confidence Intervals and Hypothesis Testing

- Random Sampling
- Confidence Intervals
- Hypothesis Testing
- Permutation Tests
- A/B Tests
- p-values
- Multi-Armed Bandit

# Questions and More Questions

- Have web conversions gone up?

- Which ad generates more sales?

- Which headline generates more clicks?

- Did the number of "likes" change?

# Example: What can we say about the trip distance of an average taxi trip in Jan 2017?

# Example: What can we say about the trip distance of an average taxi trip in Jan 2017?

```
In [2]:    1  df_taxi = (
           2      pd.read_csv('../data/yellowcab_demo_withdaycategories.csv',
           3                  header=1,
           4                  parse_dates=['pickup_datetime','dropoff_datetime'])
           5      .assign(
           6          weekpart = lambda df_: df_.is_weekend.apply(lambda x: 'Weekend' if x else 'Weekday'),
           7      )
           8      .loc[:,['trip_distance','is_weekend','weekpart']]
           9      .dropna()
          10  )
          11  print(df_taxi.shape)
          12  display(df_taxi.head(5))
```

(1000, 3)

|   | trip_distance | is_weekend | weekpart |
|---|---------------|------------|----------|
| 0 | 0.89          | False      | Weekday  |
| 1 | 2.70          | True       | Weekend  |
| 2 | 1.41          | True       | Weekend  |
| 3 | 0.40          | False      | Weekday  |
| 4 | 2.30          | False      | Weekday  |

# Mini Probability Review

- **Random Variable**

  - takes values from an associated probability distribution

  - Ex: trip_distance

- **Distribution**

  - describes probability of values of a Random Variable

- $P(x)$: **Probability**

  - probability of seeing $x$, takes value in [0,1]

  - Ex: P(trip_distance > 1)

- $P(x \mid y)$: **Conditional Probability**

  - probability of seeing $x$, given some $y$

  - Ex: P(trip_distance > 1 | is_weekend == True)

# Population Distributions and Sampling

- **"The World" or "Ground Truth"**

    - Ex: The length of taxi rides

- **"A Sample" or "Our Data"**

    - Ex: The length of taxi rides we saw in Jan 2017

# Population Distributions and Sampling

- **Population Distribution:** The actual distribution out in the world
  - Ex: Actual distribution of taxi trip lengths

- **Random Sample:** Our observations of the true population distrution
  - We hope this does not differ systematically from the true distribution
  - Ex: The taxi trip lengths recorded in Jan 2017

- **Sample Size (n):** The number of observations, the larger the better
  - Ex: We saw 1,000 trips

# Population Dists and Sampling

- **Population Mean vs. Sample Mean:**

  - Ex: The true mean trip length ($\mu$) vs the one we observed ($\bar{x}$)

- **Population Std. Dev. vs Sample Std. Dev.:**

  - Ex: The true spread of trip length ($\sigma$) vs the one we observed ($s$)

- **Sample Statistic:**

  - eg. mean, median, standard deviation

  - Ex: We're interested in mean trip length

- **Sampling Distribution:**

  - Distribution of the sample statistic

  - Ex: How is mean trip length distributed if we were to repeat our experiment many times?

# Things To Know First

- sample size

- shape (skewed?, multimodal?)

- location (central tendencies)

- spread (variance, standard deviation, IQR)

```
In [3]:  1  df_taxi.describe().round(2)
```

Out[3]:

|  | trip_distance |
|---|---|
| count | 1000.00 |
| mean | 2.88 |
| std | 3.68 |
| min | 0.00 |

# Sampling From the Population

# Sampling From the Population

In [4]:
```python
trip_distance_sample = df_taxi.trip_distance.sample(
    n=50,                 # our sample size
    random_state=123,    # needed for reproducability
    replace=False        # sample without replacement
)

print(trip_distance_sample.describe().round(2))
print()
print(f"sample skew = {trip_distance_sample.skew().round(2)}")
```
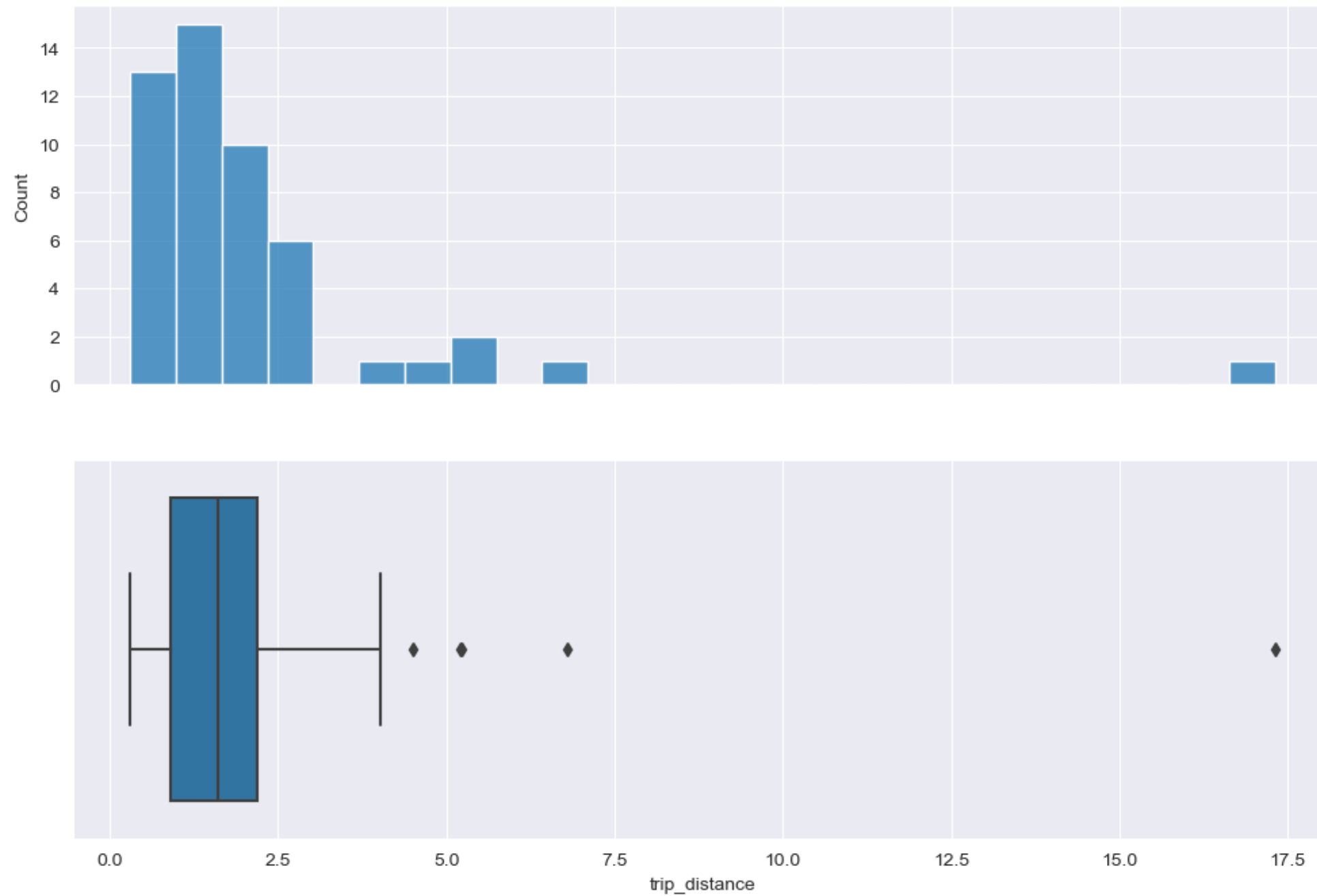
```
count     50.00
mean       2.14
std        2.56
min        0.30
25%        0.91
50%        1.60
75%        2.19
max       17.30
Name: trip_distance, dtype: float64

sample skew = 4.55
```

# Plot the distribution of our Sample

# Plot the distribution of our Sample

```python
fig,ax = plt.subplots(2,1,figsize=(12,8),sharex=True)
sns.histplot(x=trip_distance_sample, ax=ax[0]);
sns.boxplot(x=trip_distance_sample, ax=ax[1]);
```

# Define the Sample Statistic

# Define the Sample Statistic

```
1  trip_distance_sample_xbar = trip_distance_sample.mean()
2  print(f'sample mean: {trip_distance_sample_xbar:0.2f}')
```

```
sample mean: 2.14
```

# Define the Sample Statistic

```
In [6]:   1  trip_distance_sample_xbar = trip_distance_sample.mean()
          2  print(f'sample mean: {trip_distance_sample_xbar:0.2f}')

          sample mean: 2.14
```

- Is this sample statistic a good approximation?

- Let's take more samples!

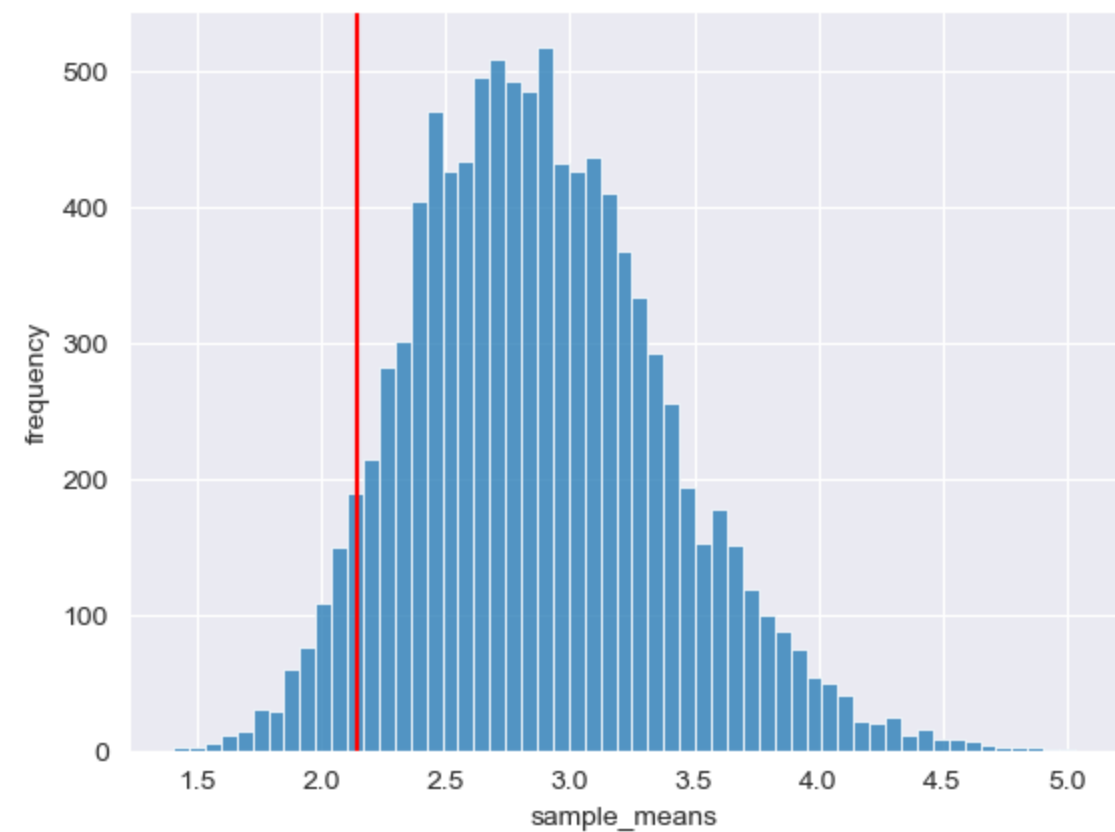# Generate Samples and Plot Distribution of Sample Stat.

# Generate Samples and Plot Distribution of Sample Stat.

```python
In [7]:
sample_means = []
for i in range(10000):
    sample_mean = df_taxi.trip_distance.sample(n=50,random_state=i).mean()
    sample_means.append(sample_mean)
```

# Generate Samples and Plot Distribution of Sample Stat.

In [7]:
```python
sample_means = []
for i in range(10000):
    sample_mean = df_taxi.trip_distance.sample(n=50,random_state=i).mean()
    sample_means.append(sample_mean)
```

In [8]:
```python
ax = sns.histplot(x=sample_means)
ax.set_xlabel('sample_means');
ax.set_ylabel('frequency');
ax.axvline(trip_distance_sample_xbar,color='red');
```

# Generate Samples and Plot Distribution of Sample Stat.

```
In [7]:  1  sample_means = []
         2  for i in range(10000):
         3      sample_mean = df_taxi.trip_distance.sample(n=50,random_state=i).mean()
         4      sample_means.append(sample_mean)
```
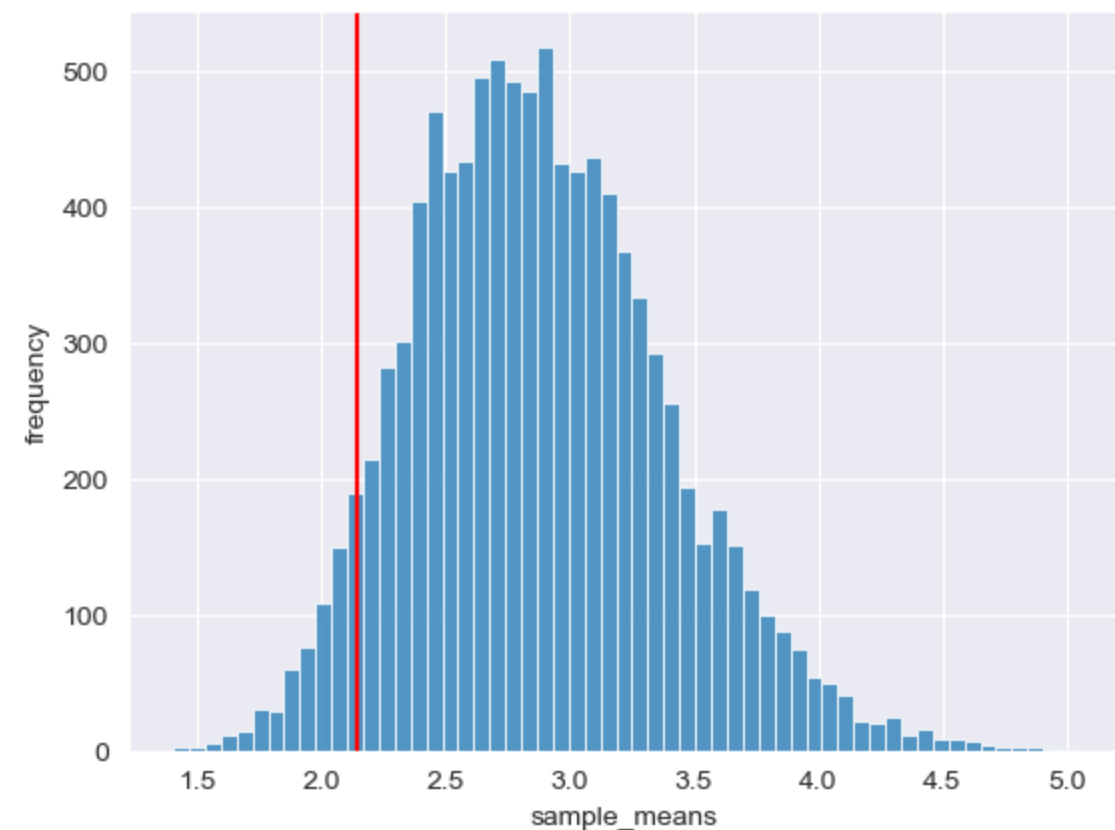
```
In [8]:  1  ax = sns.histplot(x=sample_means)
         2  ax.set_xlabel('sample_means');
         3  ax.set_ylabel('frequency');
         4  ax.axvline(trip_distance_sample_xbar,color='red');
```



But what if we can't generate additional samples? **Bootstrap Confidence Intervals**

# Confidence Intervals

Typically we only have one sample from the population (experimental results, survey results, etc.)

# Confidence Intervals

Typically we only have one sample from the population (experimental results, survey results, etc.)

```
In [9]:   1  n_trip = trip_distance_sample.shape[0]
          2  n_trip

Out[9]:  50
```

# Confidence Intervals

Typically we only have one sample from the population (experimental results, survey results, etc.)

In [9]:
```
1 n_trip = trip_distance_sample.shape[0]
2 n_trip
```

Out[9]: 50

In [10]:
```
1 trip_distance_sample_xbar = trip_distance_sample.mean()
2 print(f'sample mean: {trip_distance_sample_xbar:0.2f}')
```

sample mean: 2.14

# Confidence Intervals

Typically we only have one sample from the population (experimental results, survey results, etc.)

```
In [9]:  1  n_trip = trip_distance_sample.shape[0]
         2  n_trip
```

```
Out[9]:  50
```

```
In [10]:  1  trip_distance_sample_xbar = trip_distance_sample.mean()
          2  print(f'sample mean: {trip_distance_sample_xbar:0.2f}')
```

```
sample mean: 2.14
```

- What is the spread of our sample statistic?

- What other values would it be reasonable to observe?
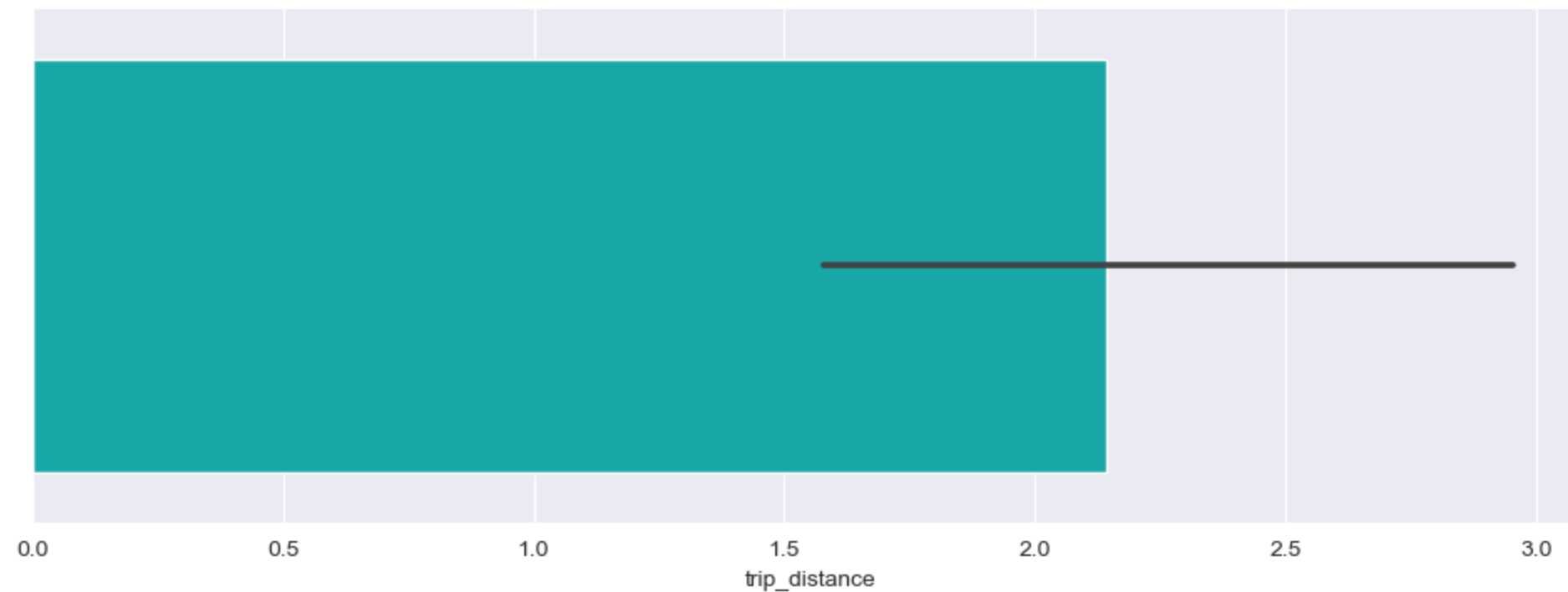
# Plotting Confidence Intervals with Seaborn

# Plotting Confidence Intervals with Seaborn

```python
fig,ax = plt.subplots(1,1,figsize=(12,4))

sns.barplot(x=trip_distance_sample,
            estimator=np.mean,    # default sample statistic
            ci=95,                # default 95% CI
            n_boot=1000,          # default number of bootstrap samples
            color='c',
            );
```

```
/var/folders/78/vhnqkq8n45dd4gj4f5qx8yb00000gn/T/ipykernel_53735/3208664956.py:3: FutureWarning:

The `ci` parameter is deprecated. Use `errorbar=('ci', 95)` for the same effect.

  sns.barplot(x=trip_distance_sample,
```
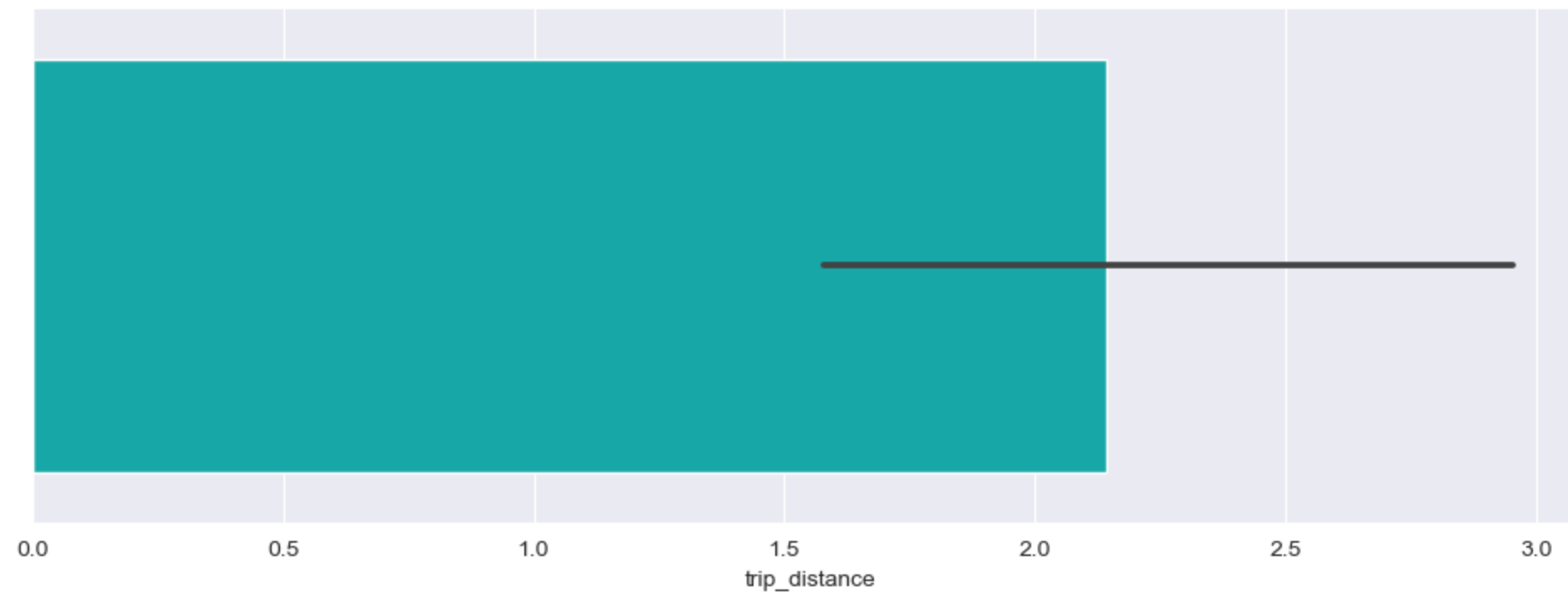
# Plotting Confidence Intervals with Seaborn

```
In [11]:  1  fig,ax = plt.subplots(1,1,figsize=(12,4))
          2
          3  sns.barplot(x=trip_distance_sample,
          4              estimator=np.mean,   # default sample statistic
          5              ci=95,               # default 95% CI
          6              n_boot=1000,         # default number of bootstrap samples
          7              color='c',
          8             );
```

```
/var/folders/78/vhnqkq8n45dd4gj4f5qx8yb00000gn/T/ipykernel_53735/3208664956.py:3: FutureWarning:

The `ci` parameter is deprecated. Use `errorbar=('ci', 95)` for the same effect.

  sns.barplot(x=trip_distance_sample,
```



- How are these confidence intervals generated from only one sample?

# Generate Confidence Intervals

**Bootstrapping:** sampling with replacement

**Bootstrap Confidence Interval:** create confidence interval using bootstrap samples

# Generate Confidence Intervals

**Bootstrapping:** sampling with replacement

**Bootstrap Confidence Interval:** create confidence interval using bootstrap samples

1. draw a random sample of size $n$ from the data

2. record the sample statistic from this random sample

3. repeat 1 and 2 many times

4. for an $x$% CI, find the trim points to remove $\frac{1}{2}\left(1 - \frac{x}{100}\right)$ of the data from both ends

5. those trim points are the endpoints of the the $x$% bootstrap CI

# 1. & 2. Draw a Random Sample and Record Statistic

# 1. & 2. Draw a Random Sample and Record Statistic

```
In [12]:   1  # 1. draw a random sample with replacement
           2  random_sample = trip_distance_sample.sample(
           3      n=trip_distance_sample.shape[0], # same size as number of observations (or frac=1)
           4      replace=True,                    # sample with replacement
           5      random_state=123                 # for reproducability
           6  )
           7  random_sample.head(3)

Out[12]:  691     0.7
          50      0.8
          882     6.8
          Name: trip_distance, dtype: float64
```

# 1. & 2. Draw a Random Sample and Record Statistic

In [12]:
```python
1  # 1. draw a random sample with replacement
2  random_sample = trip_distance_sample.sample(
3      n=trip_distance_sample.shape[0], # same size as number of observations (or frac=1)
4      replace=True,                    # sample with replacement
5      random_state=123                 # for reproducability
6  )
7  random_sample.head(3)
```

Out[12]:
```
691     0.7
50      0.8
882     6.8
Name: trip_distance, dtype: float64
```

In [13]:
```python
1  # 2. record sample statistic
2  sample_means = []
3  sample_means.append(random_sample.mean())
4  [x.round(2) for x in sample_means]
```

Out[13]: [1.82]

# 3. Repeat Many Times

# 3. Repeat Many Times

```
In [14]:    1  # tqdm gives us a progress bar when looping
            2  from tqdm.notebook import tqdm
```

```
In [15]:    1  df_taxi.trip_distance.shape[0]
```

Out[15]:  1000

# 3. Repeat Many Times

```
In [14]:   1  # tqdm gives us a progress bar when looping
           2  from tqdm.notebook import tqdm
```

```
In [15]:   1  df_taxi.trip_distance.shape[0]
```

Out[15]: 1000

```
In [16]:   1  # 3. repeat 1 and 2 many times
           2  num_iterations = 10000
           3  sample_means = []
           4
           5  for i in tqdm(range(num_iterations)):
           6      # 1. draw a random sample of size *n* from the data
           7      random_sample = df_taxi.trip_distance.sample(n=df_taxi.trip_distance.shape[0], # or frac=1
           8                                                   replace=True,    # sample with replacement
           9                                                   random_state=i   # for reproducability
          10                                                  )
          11      # 2. record the sample statistic from this random sample
          12      sample_means.append(random_sample.mean())
          13
          14  # convert into a numpy array
          15  sample_means = np.array(sample_means)
          16
          17  sample_means[:10].round(2)
```
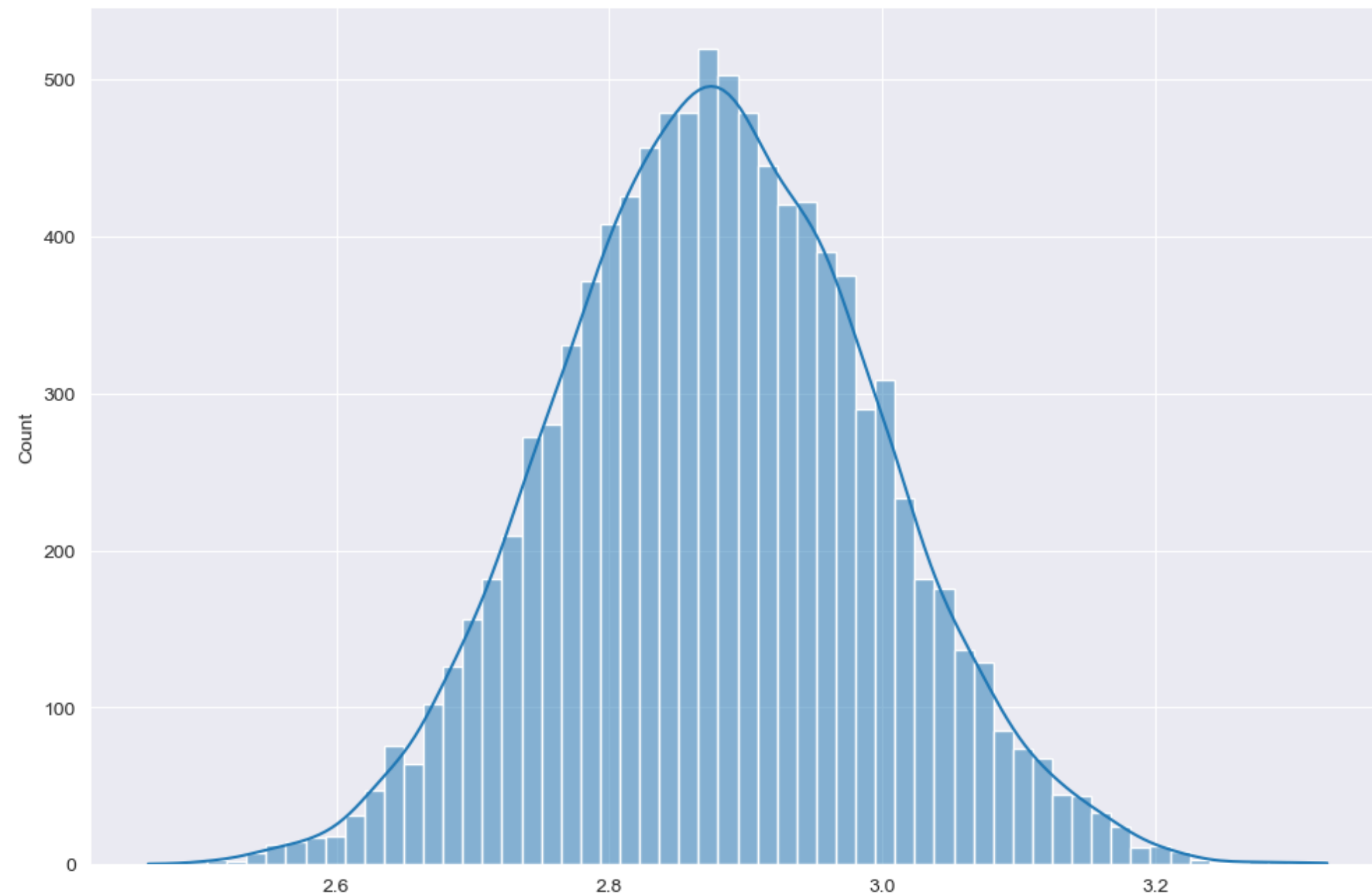
100%                          10000/10000 [00:03<00:00, 3032.86it/s]

Out[16]: array([2.98, 2.96, 3.02, 2.96, 3.01, 2.92, 2.74, 2.7 , 2.68, 2.82])
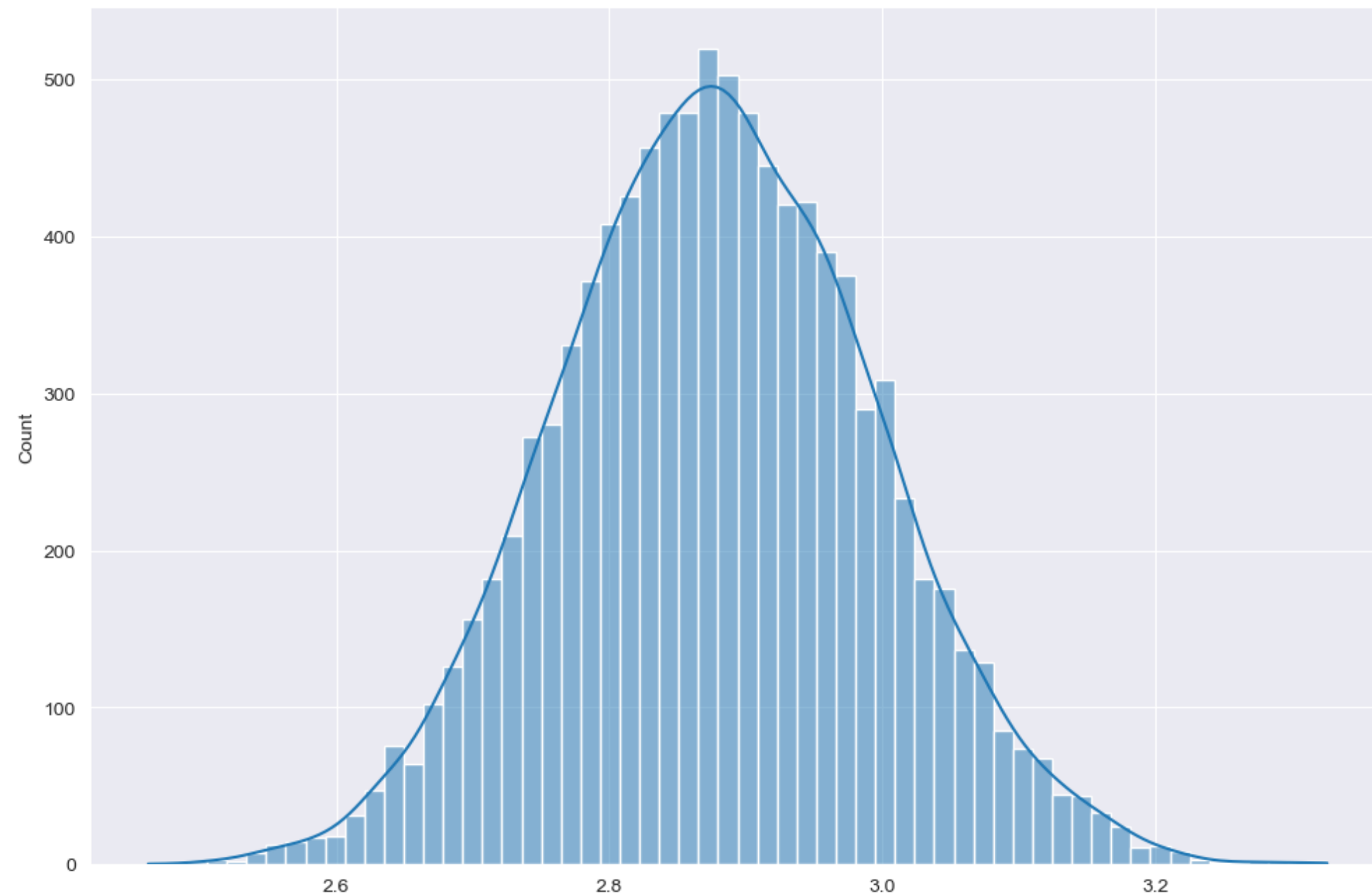
# Distribution of Sample Means?

# Distribution of Sample Means?

```
In [17]:  1  fig,ax = plt.subplots(1,1,figsize=(12,8))
          2  sns.histplot(x=sample_means,kde=True);
```

# Distribution of Sample Means?

```
In [17]:   1  fig,ax = plt.subplots(1,1,figsize=(12,8))
           2  sns.histplot(x=sample_means,kde=True);
```



- Between what two values do 95% of these samples fall?

# 4 & 5 Find CI Endpoints

# 4 & 5 Find CI Endpoints

```
In [18]:   1  # 4. For a 95% conf. int., trim off .5*(1-(95/100)) of the data from both ends
           2
           3  # calculate where to trim
           4  trim = .5*(1-.95) * num_iterations
           5
           6  # find the closest integer
           7  trim = int(np.round(trim))
           8  trim
```

Out[18]:   250

# 4 & 5 Find CI Endpoints

```
In [18]:  1  # 4. For a 95% conf. int., trim off .5*(1-(95/100)) of the data from both ends
          2
          3  # calculate where to trim
          4  trim = .5*(1-.95) * num_iterations
          5
          6  # find the closest integer
          7  trim = int(np.round(trim))
          8  trim
```
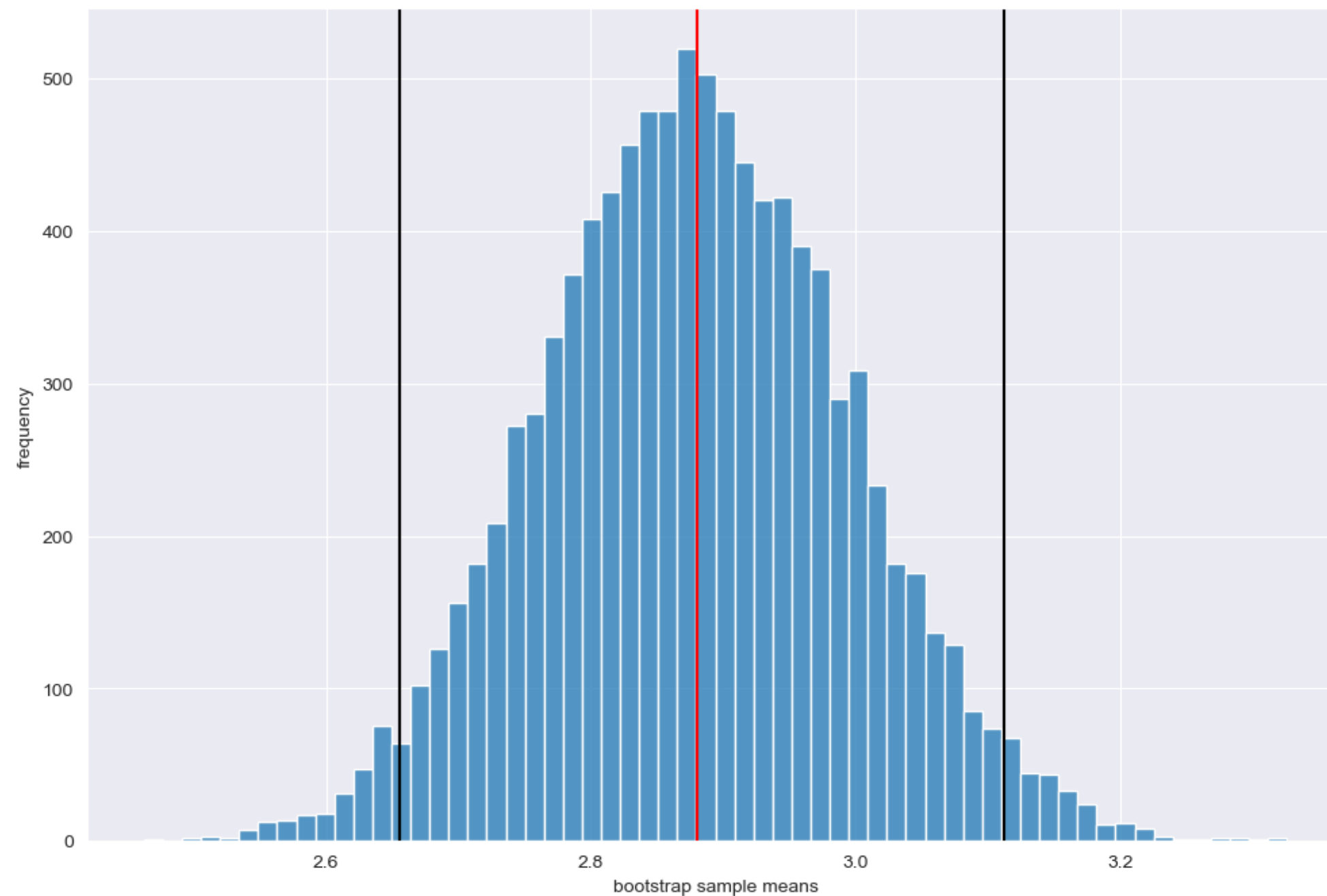
Out[18]:  250

```
In [19]:  1  # for 1000 iterations and a 95% CI, we want to find the 25th value and (1000-25)th value
          2
          3  # 5. those trim points are the endpoints of the the x% Bootstrap CI
          4
          5  ci = np.sort(sample_means)[[trim,-trim-1]] # sort the array first!
          6  ci.round(2)
```

Out[19]:  array([2.66, 3.11])

# Plotting Distribution of Sample Means With CIs

# Plotting Distribution of Sample Means With CIs

```
In [20]:  1  fig,ax = plt.subplots(1,1,figsize=(12,8))
          2  ax = sns.histplot(sample_means)
          3  ax.axvline(df_taxi.trip_distance.mean(), color='r');
          4  ax.axvline(ci[0],color='k');ax.axvline(ci[1],color='k')
          5  ax.set_xlabel('bootstrap sample means');
          6  ax.set_ylabel('frequency');
```
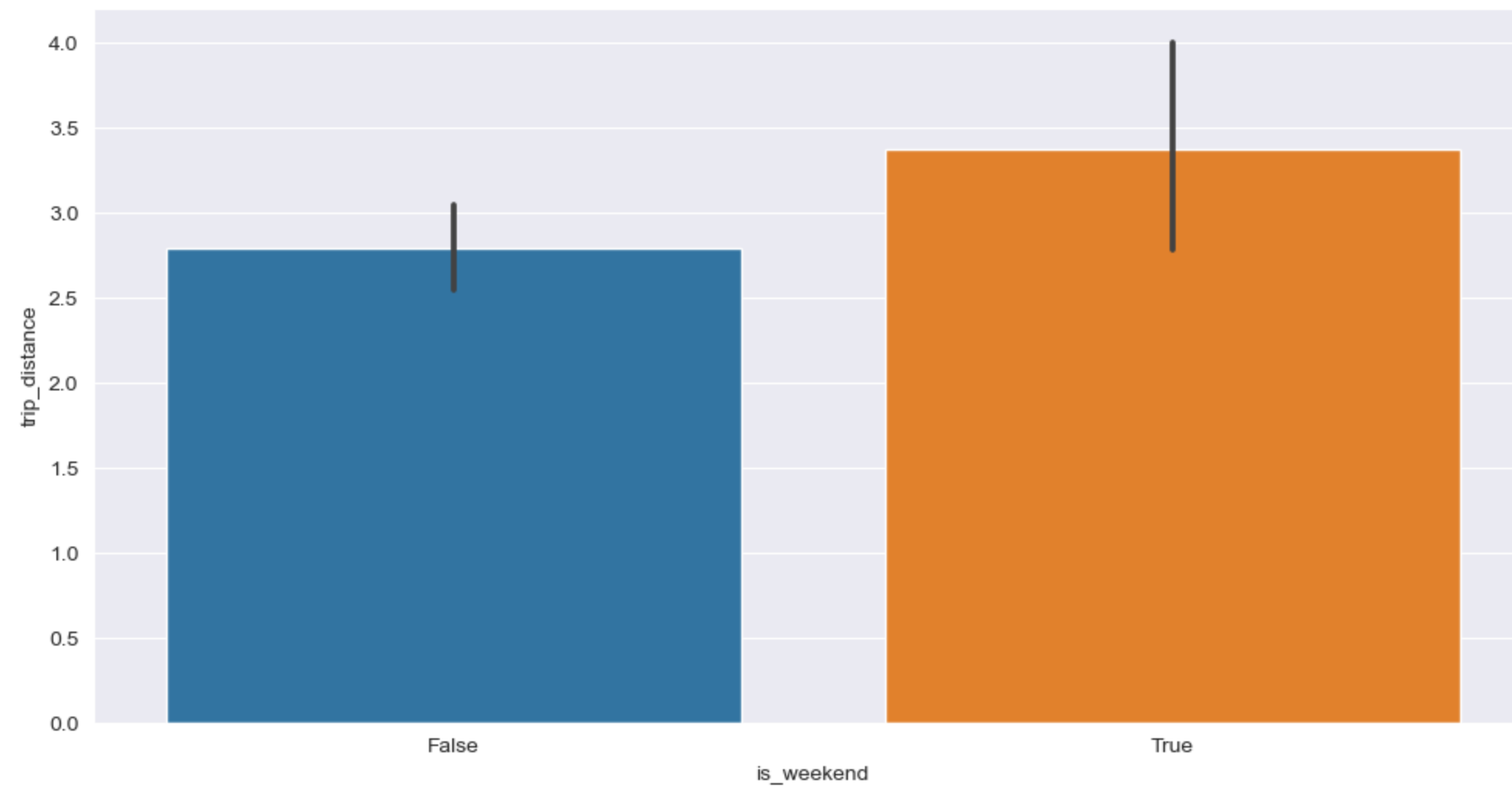
# Interpreting CIs

- **Does NOT tell us:** "the probability that the true value lies within that interval"

- **Tells us:** something about the variablity of this statistic

- **Tells us:** how confident we should be that our parameter lies in the interval

*If confidence intervals are constructed using a given confidence level from an infinite number of independent sample statistics, the proportion of those intervals that contain the true value of the parameter will be equal to the confidence level.*

# Interpreting CIs

# Interpreting CIs

```
In [21]:  1  fig,ax = plt.subplots(1,1,figsize=(12,6))
          2  sns.barplot(x='is_weekend',y='trip_distance',data=df_taxi);
```

# Questions re CIs?
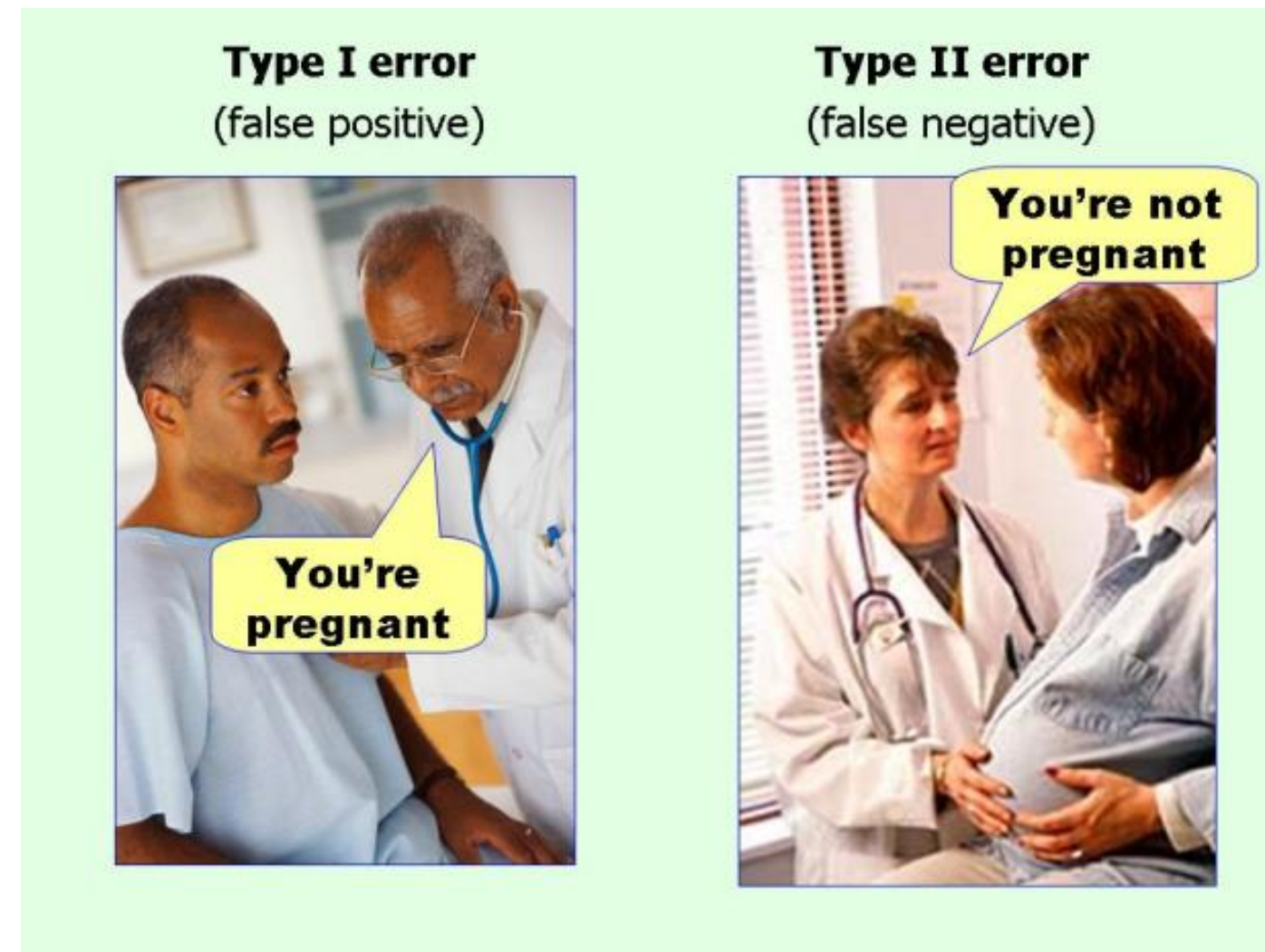
# Hypothesis Testing

- Ex: Is the average trip longer on weekends compared to weekdays?

- Ex: Does one advertisement lead to more sales than another?

- **Null Hypothesis:** $H_0$

    - the thing we're observing is happening due to random chance

    - there are no differences between two groups

- **Alternative Hypothesis:** $H_1$

    - the thing we're observing is happening not due to random chance

    - there is a difference between two groups

- Experiment: given data, do we **accept or reject** $H_0$?

    - Ex: can we say the difference between average trip on weekdays vs. weekends isn't random?

# Errors in Hypothesis Tests



https://www.gilliganondata.com/wp-content/uploads/2009/08/TypeI_TypeII1.JPG

# Errors in Hypothesis Tests

# Significance and Power

# Significance and Power

- $P$ (reject $H_0$ | $H_0$ true) = **Significance** of test or **p-value** (Type I Error)
  - Probablity of saying **things aren't by chance when they are**
  - Ex: Saying trips on weekends are longer, when the difference *is* random
  - Ex: Saying Ad A was correlated with more sales, when the difference *is* random

# Significance and Power

- $P$ (reject $H_0$ | $H_0$ true) = **Significance** of test or **p-value** (Type I Error)
    - Probablity of saying **things aren't by chance when they are**
    - Ex: Saying trips on weekends are longer, when the difference *is* random
    - Ex: Saying Ad A was correlated with more sales, when the difference *is* random


- $P$ (reject $H_0$ | $H_1$ true) = **Power** of test (1-Type II Error)
    - Probability of saying **things aren't by chance when they aren't**
    - Ex: Saying trips on weekends are longer, when the difference *is not* random
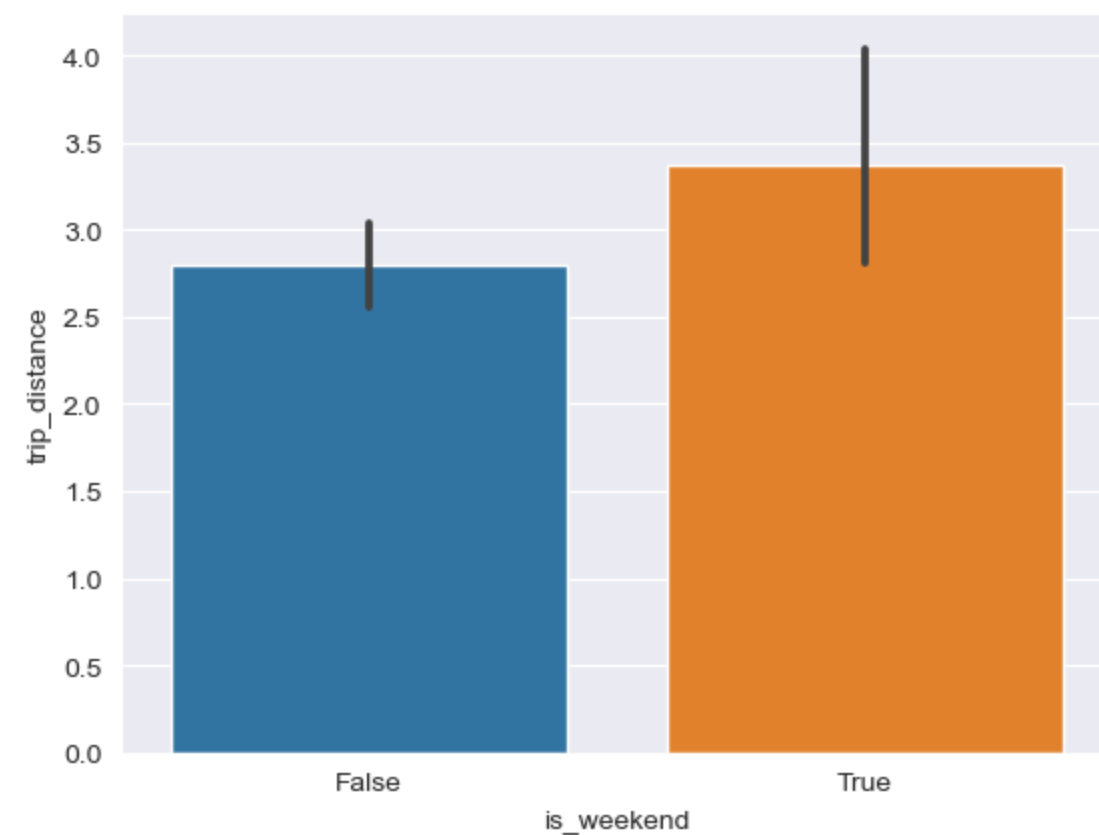    - Ex: Saying Ad A was correlated with more sales, when the difference *is not* random

# Ex: Trip-Distance by Weekday vs. Weekend

# Ex: Trip-Distance by Weekday vs. Weekend

- Question: Is the average trip_distance different on weekdays vs weekends?

```
In [22]:   1  sns.barplot(x='is_weekend',y='trip_distance',data=df_taxi);
```

# Ex: Trip-Distance by Weekday vs. Weekend, Define the Metric

- **Metric:** the measure we're interested in
    - Ex: We're interested in a difference of means: Weekday - Weekend

# Ex: Trip-Distance by Weekday vs. Weekend, Define the Metric

- **Metric:** the measure we're interested in
    - Ex: We're interested in a difference of means: Weekday - Weekend

```
In [23]:  1  mean_weekend = df_taxi.loc[df_taxi.is_weekend,'trip_distance'].mean()
          2  mean_weekday = df_taxi.loc[~df_taxi.is_weekend,'trip_distance'].mean()
          3  observed_trip_metric = mean_weekend-mean_weekday
          4  print(f'observed metric: {observed_trip_metric.round(2)}')

          observed metric: 0.58
```

# Ex: Trip-Distance by Weekday vs. Weekend, Define the Metric

- **Metric:** the measure we're interested in
  - Ex: We're interested in a difference of means: Weekday - Weekend

```
In [23]:   1  mean_weekend = df_taxi.loc[df_taxi.is_weekend,'trip_distance'].mean()
           2  mean_weekday = df_taxi.loc[~df_taxi.is_weekend,'trip_distance'].mean()
           3  observed_trip_metric = mean_weekend-mean_weekday
           4  print(f'observed metric: {observed_trip_metric.round(2)}')

           observed metric: 0.58
```

- Is this surprising? Should we reject the null?
  - Assuming that $H_0$ is true, is this observation surprising?

# Permutation Test

- How do we generate additional samples of the difference in means? Resampling!

- Need to repeatedly split the data into two groups and take the difference in means

- One way to do this: combine, permute (reorder) and split

# Permutation Test

1. combine groups together (assume $H_0$ is true)

2. permute (reorder) observations

3. create new groups (same sizes as original groups)

4. calculate metric

5. repeat many times

6. see where our original observation falls in the distribution of sample statistics

# Ex: Avg Trip-Distance on Weekday vs. Weekend, Permutation Test

## Ex: Avg Trip-Distance on Weekday vs. Weekend, Permutation Test

```python
In [24]:  1  # 0. get group sizes
          2  n_weekend = df_taxi.is_weekend.sum()
          3  n_weekday = (~df_taxi.is_weekend).sum()
          4  print(f'{n_weekend=} {n_weekday=}')
          5  assert n_weekday + n_weekend == df_taxi.shape[0]
```

```
n_weekend=150 n_weekday=850
```

# Ex: Avg Trip-Distance on Weekday vs. Weekend, Permutation Test

In [24]:
```python
1  # 0. get group sizes
2  n_weekend = df_taxi.is_weekend.sum()
3  n_weekday = (~df_taxi.is_weekend).sum()
4  print(f'{n_weekend=} {n_weekday=}')
5  assert n_weekday + n_weekend == df_taxi.shape[0]
```

```
n_weekend=150 n_weekday=850
```

In [25]:
```python
1  # 1. combine groups together (assume H0 is true)
2  trip_distances = df_taxi.trip_distance
3  trip_distances[:2]
```

Out[25]:
```
0    0.89
1    2.70
Name: trip_distance, dtype: float64
```

# Ex: Avg Trip-Distance on Weekday vs. Weekend, Permutation Test

In [24]:
```python
1  # 0. get group sizes
2  n_weekend = df_taxi.is_weekend.sum()
3  n_weekday = (~df_taxi.is_weekend).sum()
4  print(f'{n_weekend=} {n_weekday=}')
5  assert n_weekday + n_weekend == df_taxi.shape[0]
```

n_weekend=150 n_weekday=850

In [25]:
```python
1  # 1. combine groups together (assume H0 is true)
2  trip_distances = df_taxi.trip_distance
3  trip_distances[:2]
```

Out[25]:  0     0.89
          1     2.70
          Name: trip_distance, dtype: float64

In [26]:
```python
1  # 2. permute observations
2  permuted_trip_distances = trip_distances.sample(frac=1,replace=False,random_state=123)
3  permuted_trip_distances[:2]
```

/var/folders/78/vhnqkq8n45dd4gj4f5qx8yb00000gn/T/ipykernel_53735/3380519961.py:3: FutureWarning: The behavior of `series[i:j]`
with an integer-dtype index is deprecated. In a future version, this will be treated as *label-based* indexing, consistent with
e.g. `series[i]` lookups. To retain the old behavior, use `series.iloc[i:j]`. To get the future behavior, use `series.loc[i:j]`
`.
  permuted_trip_distances[:2]

Out[26]:  131    2.13
          203    2.15
          Name: trip_distance, dtype: float64

# Ex: Avg Trip-Distance on Weekday vs. Weekend, Permutation Test

In [24]:
```python
# 0. get group sizes
n_weekend = df_taxi.is_weekend.sum()
n_weekday = (~df_taxi.is_weekend).sum()
print(f'{n_weekend=} {n_weekday=}')
assert n_weekday + n_weekend == df_taxi.shape[0]
```

n_weekend=150 n_weekday=850

In [25]:
```python
# 1. combine groups together (assume H0 is true)
trip_distances = df_taxi.trip_distance
trip_distances[:2]
```

Out[25]:  0     0.89
          1     2.70
          Name: trip_distance, dtype: float64

In [26]:
```python
# 2. permute observations
permuted_trip_distances = trip_distances.sample(frac=1,replace=False,random_state=123)
permuted_trip_distances[:2]
```

/var/folders/78/vhnqkq8n45dd4gj4f5qx8yb00000gn/T/ipykernel_53735/3380519961.py:3: FutureWarning: The behavior of `series[i:j]` with an integer-dtype index is deprecated. In a future version, this will be treated as *label-based* indexing, consistent with e.g. `series[i]` lookups. To retain the old behavior, use `series.iloc[i:j]`. To get the future behavior, use `series.loc[i:j]`.
  permuted_trip_distances[:2]

Out[26]:  131     2.13
          203     2.15
          Name: trip_distance, dtype: float64

In [27]:
```python
# 3. create new groups
rand_mean_weekend = permuted_trip_distances[:n_weekend].mean()
rand_mean_weekday = permuted_trip_distances[n_weekend:].mean()

# 4. calculate metric
```

# Ex: Trip-Distance, Permutation Test Continued

# Ex: Trip-Distance, Permutation Test Continued

```python
In [28]:
 1  # 5. repeat many times
 2  rand_mean_trip_diffs = []
 3  iterations = 10_000
 4
 5  for i in tqdm(range(iterations)):
 6      permuted_trip_distances = trip_distances.sample(frac=1,replace=False,random_state=i)
 7
 8      rand_mean_weekend = permuted_trip_distances[:n_weekend].mean()
 9      rand_mean_weekday = permuted_trip_distances[n_weekend:].mean()
10
11      rand_mean_trip_diffs.append(rand_mean_weekend - rand_mean_weekday)
12
13  rand_mean_trip_diffs = np.array(rand_mean_trip_diffs) # convert list to numpy array
14
15  rand_mean_trip_diffs[:5].round(2)
```

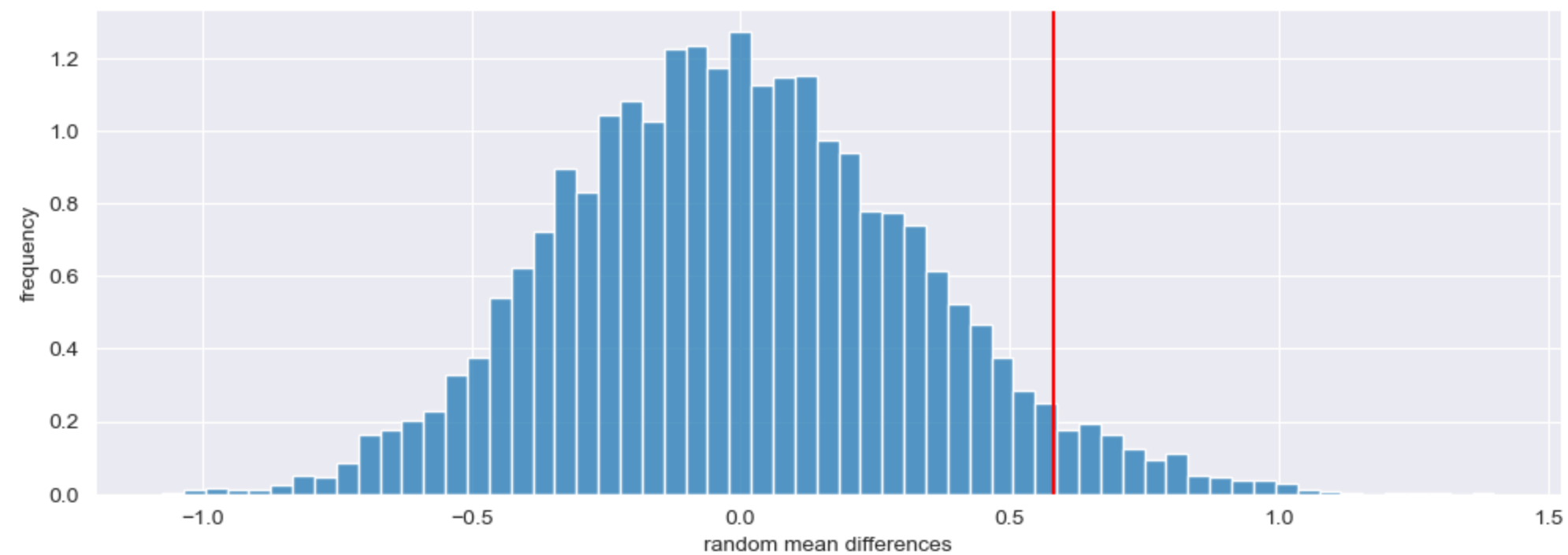100%                                        10000/10000 [00:06<00:00, 1761.25it/s]

```
/var/folders/78/vhnqkq8n45dd4gj4f5qx8yb00000gn/T/ipykernel_53735/1315191553.py:8: FutureWarning: The behavior of `series[i:j]`
with an integer-dtype index is deprecated. In a future version, this will be treated as *label-based* indexing, consistent with
e.g. `series[i]` lookups. To retain the old behavior, use `series.iloc[i:j]`. To get the future behavior, use `series.loc[i:j]
`.
  rand_mean_weekend = permuted_trip_distances[:n_weekend].mean()
/var/folders/78/vhnqkq8n45dd4gj4f5qx8yb00000gn/T/ipykernel_53735/1315191553.py:9: FutureWarning: The behavior of `series[i:j]`
with an integer-dtype index is deprecated. In a future version, this will be treated as *label-based* indexing, consistent with
e.g. `series[i]` lookups. To retain the old behavior, use `series.iloc[i:j]`. To get the future behavior, use `series.loc[i:j]
`.
  rand_mean_weekday = permuted_trip_distances[n_weekend:].mean()
```

Out[28]:  array([-0.49, -0.21,  0.58, -0.09, -0.37])

# Ex: Trip-Distance, Permutation Test Continued

# Ex: Trip-Distance, Permutation Test Continued

```
In [29]:    1  # 6. see where our original observation falls
            2  fig,ax = plt.subplots(1,1,figsize=(12,4))
            3  ax = sns.histplot(x=rand_mean_trip_diffs, stat='density')
            4  ax.set_xlabel('random mean differences');ax.set_ylabel('frequency');
            5  ax.axvline(observed_trip_metric, color='r');
```
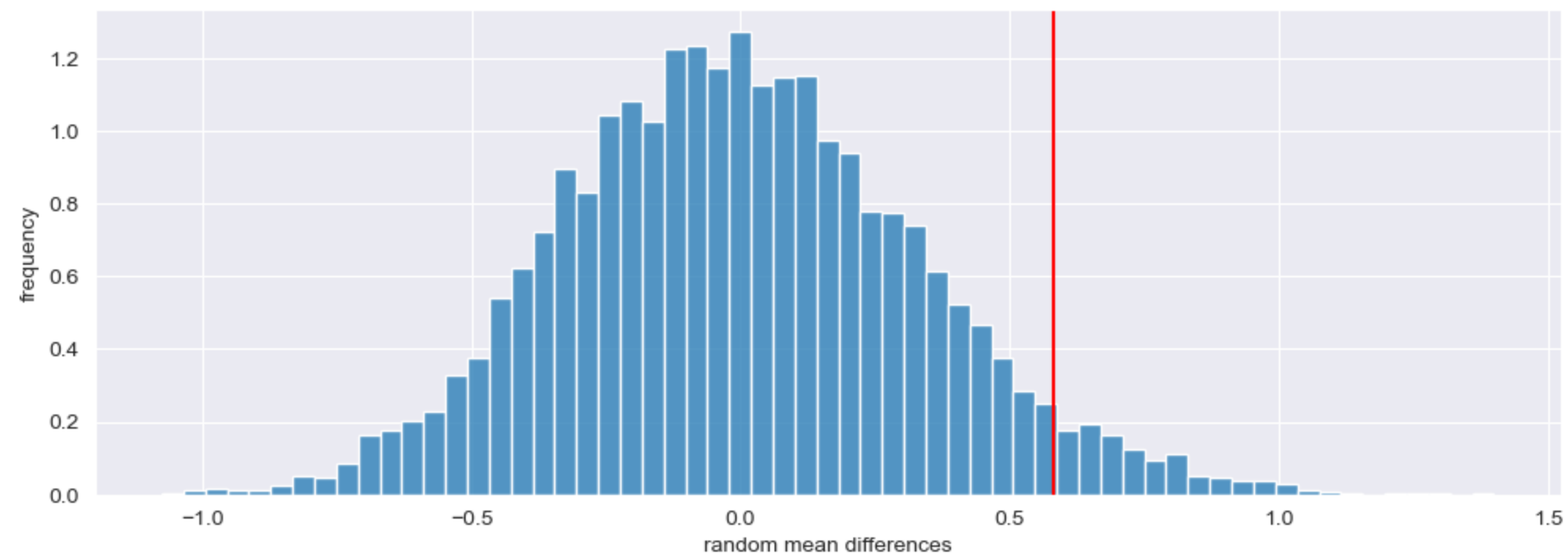
# Ex: Trip-Distance, Permutation Test Continued

```
In [29]:  1  # 6. see where our original observation falls
          2  fig,ax = plt.subplots(1,1,figsize=(12,4))
          3  ax = sns.histplot(x=rand_mean_trip_diffs, stat='density')
          4  ax.set_xlabel('random mean differences');ax.set_ylabel('frequency');
          5  ax.axvline(observed_trip_metric, color='r');
```



- This looks like a normal distribution?

- Why would that be?

- How can we turn this into a Standard Normal distribution...

# Aside: Central Limit Theorem

# Aside: Central Limit Theorem

If all samples are randomly drawn from the same sample population:

For reasonably large samples (usually $n \geq 30$), the distribution of sample mean $\bar{x}$ is normal regardless of the distribution of $X$.

The sampling distribution of $\bar{x}$ becomes approximately normal as the the sample size $n$ gets large.

Ex:

- $X$ = trip_distance
- $\bar{x}$ = mean trip_distance
- $n$ = 50

# Aside: What is Normal?

# Aside: What is Normal?

distribution defined by mean ($\mu$) and standard deviation ($\sigma$)

$$N(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma}}$$

**PDF (Probability Density Function):**

- function of a continuous random variable that provides a relative likelihood of seeing a particular sample of a random variable.

# Aside: Properties of a Normal Distribution



https://towardsdatascience.com/understanding-the-68-95-99-7-rule-for-a-normal-distribution-b7b7cbf760c2

# Aside: Scipy

- Routines for numerical integration, interpolation, optimization, linear algebra, and **statistics**.
- Useful for sampling from random distributions and equation based testing

# Aside: Scipy

- Routines for numerical integration, interpolation, optimization, linear algebra, and **statistics**.
- Useful for sampling from random distributions and equation based testing



```
In [30]:  1  import scipy as sp
```

# Aside: Plotting a Standard Normal Distribution

- **Standard Normal:** $\mu=0$, $\sigma=1$
- Often referred to as $Z$

# Aside: Plotting a Standard Normal Distribution

- **Standard Normal:** μ=0, σ=1
- Often referred to as $Z$

```
In [31]:  1  x = np.random.normal(0,1,size=100_000)                    # generate many random samples
          2  fig,ax = plt.subplots(1,1,figsize=(12,6))
          3  ax = sns.histplot(x=x,stat='density',kde=True);           # using density to normalize bin counts
          4  ax.set_xlabel('$x$');ax.set_ylabel('$N(x;\mu=0,\sigma=1)$'); # using latex in labels
          5  ax.vlines([-1,1],0,sp.stats.norm.pdf(1), colors='k');     # 1 standard deviation
          6  ax.vlines([-2,2],0,sp.stats.norm.pdf(2), colors='r');     # 2 standard deviations
```

# Normalization: z-score

Convert our distribution to an approximation of standard normal

1. shift mean to 0

2. transform to standard deviation of 1

$$z = \frac{x - \bar{x}}{s}$$

# Normalization: z-score

Convert our distribution to an approximation of standard normal

1. shift mean to 0

2. transform to standard deviation of 1

$$z = \frac{x - \bar{x}}{s}$$

```
In [32]:    1  rand_mean_trip_diffs_xbar = np.mean(rand_mean_trip_diffs)
            2  rand_mean_trip_diffs_s = np.std(rand_mean_trip_diffs)
            3
            4  rand_mean_trip_zscores = (rand_mean_trip_diffs - rand_mean_trip_diffs_xbar) / rand_mean_trip_diffs_s
            5  list(zip(rand_mean_trip_diffs[:3].round(2),rand_mean_trip_zscores[:3].round(2)))

Out[32]:  [(-0.49, -1.47), (-0.21, -0.64), (0.58, 1.77)]
```

# Normalization: z-score

Convert our distribution to an approximation of standard normal

1. shift mean to 0

2. transform to standard deviation of 1

$$z = \frac{x - \bar{x}}{s}$$

```
In [32]:  1  rand_mean_trip_diffs_xbar = np.mean(rand_mean_trip_diffs)
          2  rand_mean_trip_diffs_s = np.std(rand_mean_trip_diffs)
          3
          4  rand_mean_trip_zscores = (rand_mean_trip_diffs - rand_mean_trip_diffs_xbar) / rand_mean_trip_diffs_s
          5  list(zip(rand_mean_trip_diffs[:3].round(2),rand_mean_trip_zscores[:3].round(2)))
```

Out[32]: [(-0.49, -1.47), (-0.21, -0.64), (0.58, 1.77)]

```
In [33]:  1  observed_trip_metric_zscore =  (observed_trip_metric - rand_mean_trip_diffs_xbar) / rand_mean_trip_diffs_s
          2  observed_trip_metric.round(2),observed_trip_metric_zscore.round(2)
```

Out[33]: (0.58, 1.76)

# Ex: Trip-Distance, Permutation Test Continued

# Ex: Trip-Distance, Permutation Test Continued

```python
In [34]:  1  # 6. see where our original observation falls (normalized)
          2  fig,ax = plt.subplots(1,1,figsize=(12,2))
          3  ax = sns.histplot(rand_mean_trip_zscores, stat='density')
          4  ax.set_xlabel('random mean differences normed');ax.set_ylabel('frequency');
          5  ax.axvline(observed_trip_metric_zscore,color='r');
```

# Ex: Trip-Distance, Permutation Test Continued

In [34]:
```python
# 6. see where our original observation falls (normalized)
fig,ax = plt.subplots(1,1,figsize=(12,2))
ax = sns.histplot(rand_mean_trip_zscores, stat='density')
ax.set_xlabel('random mean differences normed');ax.set_ylabel('frequency');
ax.axvline(observed_trip_metric_zscore,color='r');
```
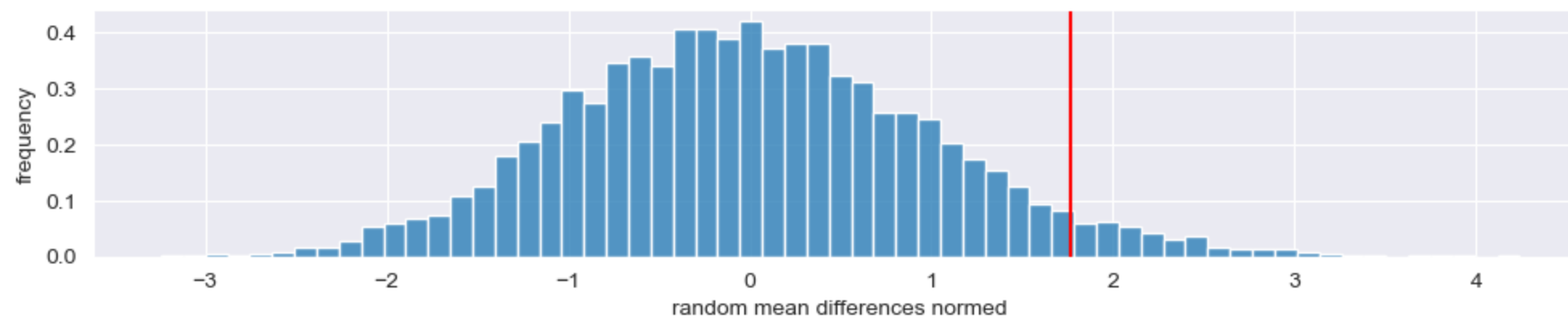


In [35]:
```python
# Compared to our original distribution
fig,ax = plt.subplots(1,1,figsize=(12,2))
ax = sns.histplot(x=rand_mean_trip_diffs, stat='density')
ax.set_xlabel('random mean differences');ax.set_ylabel('frequency');
ax.axvline(observed_trip_metric, color='r');
```

# Why Use Permutation Tests?
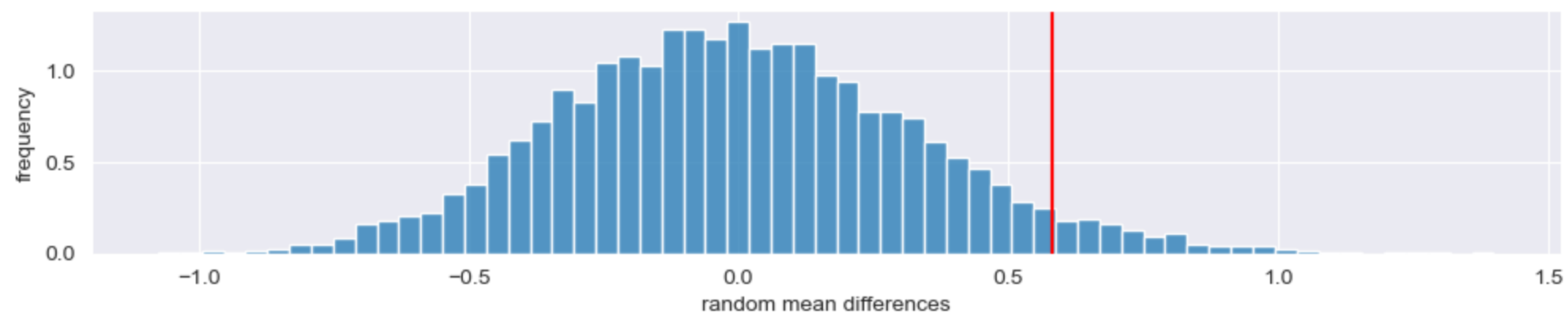
- data can be numeric or boolean (ex. temperature, conversion, etc)

- group sizes can be different

- assumptions about normally distributed data are not needed (with many permutations)

# A/B Tests

- Do one of two treatments produce superior results?

  - testing two prices to determine which generates more profit
  - testing two web headlines to determine which produces more clicks
  - testing two advertisements to see which produces more conversions

- Often Used Test Statistics

  - difference in means
  - difference in counts

# Ex: Webpages and Sales

- Question: Which webpage leads to more sales?

- Potential Issue: what if sales are large but infrequent?

- **Proxy Variable:** stand in for true value of interest

    - Ex: Assume 'time on page' is correlated with sales

# Ex: Webpages and Sales

- Question: Which webpage leads to more sales?

- Potential Issue: what if sales are large but infrequent?

- **Proxy Variable:** stand in for true value of interest

    - Ex: Assume 'time on page' is correlated with sales

```python
1  session_times = pd.read_csv('../data/web_page_data.csv')
2  print(session_times.shape)
3
4  session_times.head(3)
```

In [36]:

(36, 2)

Out[36]:

|   | Page | Time |
|---|------|------|
| 0 | Page A | 12.6 |
| 1 | Page B | 151.7 |
| 2 | Page A | 21.0 |

# Ex: Webpages and Sales

# Ex: Webpages and Sales

```
In [37]:  1  fig,ax = plt.subplots(1,3,figsize=(12,4),sharey=True)
          2  sns.boxplot(x='Page',y='Time',data=session_times,ax=ax[0]);
          3  sns.swarmplot(x='Page',y='Time',data=session_times,ax=ax[1]);
          4  sns.barplot(x='Page',y='Time',data=session_times,ax=ax[2]);
```

# Ex: Webpages and Sales, Define the Metric

- **Metric:** the measure we're interested in
  - Ex: We're interested in a difference of means (Page A - Page B)

# Ex: Webpages and Sales, Define the Metric

- **Metric:** the measure we're interested in
    - Ex: We're interested in a difference of means (Page A - Page B)

```python
In [38]:   1  mean_a = session_times.loc[session_times.Page == 'Page A','Time'].mean()
           2  mean_b = session_times[session_times.Page == 'Page B'].Time.mean()
           3  observed_ad_metric = mean_a-mean_b
           4  print('observed metric: {:0.2f}'.format(observed_ad_metric))

observed metric: -49.77
```

# Ex: Websites and Sales, Permutation Test

# Ex: Websites and Sales, Permutation Test

```
In [39]:   1  # 0. get group sizes
           2  n_a = (session_times.Page == 'Page A').sum()
           3  n_b = session_times.shape[0] - n_a
           4  print(f'{n_a=} {n_b=}')
```

```
n_a=21 n_b=15
```

# Ex: Websites and Sales, Permutation Test

```
In [39]:  1  # 0. get group sizes
          2  n_a = (session_times.Page == 'Page A').sum()
          3  n_b = session_times.shape[0] - n_a
          4  print(f'{n_a=} {n_b=}')
```

n_a=21 n_b=15

```
In [40]:  1  # 1. combine groups together (assume H0 is true)
          2  session_times.Time[:2]
```

```
Out[40]:  0      12.6
          1     151.7
          Name: Time, dtype: float64
```

# Ex: Websites and Sales, Permutation Test

```
In [39]:  1  # 0. get group sizes
          2  n_a = (session_times.Page == 'Page A').sum()
          3  n_b = session_times.shape[0] - n_a
          4  print(f'{n_a=} {n_b=}')
```

n_a=21 n_b=15

```
In [40]:  1  # 1. combine groups together (assume H0 is true)
          2  session_times.Time[:2]
```

```
Out[40]:  0     12.6
          1    151.7
          Name: Time, dtype: float64
```

```
In [41]:  1  # 2. permute observations
          2  session_times_permuted = session_times.Time.sample(frac=1,replace=False,random_state=123)
          3  session_times_permuted[:2]
```

```
/var/folders/78/vhnqkq8n45dd4gj4f5qx8yb00000gn/T/ipykernel_53735/1235933868.py:3: FutureWarning: The behavior of `series[i:j]`
with an integer-dtype index is deprecated. In a future version, this will be treated as *label-based* indexing, consistent with
e.g. `series[i]` lookups. To retain the old behavior, use `series.iloc[i:j]`. To get the future behavior, use `series.loc[i:j]
`.
  session_times_permuted[:2]
```

```
Out[41]:  6    50.5
          8    79.2
          Name: Time, dtype: float64
```

# Ex: Websites and Sales, Permutation Test

```
In [39]:   1  # 0. get group sizes
           2  n_a = (session_times.Page == 'Page A').sum()
           3  n_b = session_times.shape[0] - n_a
           4  print(f'{n_a=} {n_b=}')
```

```
n_a=21 n_b=15
```

```
In [40]:   1  # 1. combine groups together (assume H0 is true)
           2  session_times.Time[:2]
```

```
Out[40]:   0      12.6
           1     151.7
           Name: Time, dtype: float64
```

```
In [41]:   1  # 2. permute observations
           2  session_times_permuted = session_times.Time.sample(frac=1,replace=False,random_state=123)
           3  session_times_permuted[:2]
```

```
/var/folders/78/vhnqkq8n45dd4gj4f5qx8yb00000gn/T/ipykernel_53735/1235933868.py:3: FutureWarning: The behavior of `series[i:j]`
with an integer-dtype index is deprecated. In a future version, this will be treated as *label-based* indexing, consistent with
e.g. `series[i]` lookups. To retain the old behavior, use `series.iloc[i:j]`. To get the future behavior, use `series.loc[i:j]
`.
   session_times_permuted[:2]
```

```
Out[41]:   6      50.5
           8      79.2
           Name: Time, dtype: float64
```

```
In [42]:   1  # 3. create new groups
           2  rand_mean_a = session_times_permuted[:n_a].mean()
           3  rand_mean_b = session_times_permuted[n_a:].mean()
           4
           5  # 4. calculate metric
           6  rand_mean_ad_diff = (rand_mean_a - rand_mean_b)
```

# Ex: Websites and Sales, Permutation Test Continued

# Ex: Websites and Sales, Permutation Test Continued

```
In [43]:    1  # 5. repeat many times
            2  rand_mean_ad_diffs = []
            3  iterations = 10_000
            4
            5  for i in tqdm(range(iterations)):
            6      session_times_permuted = session_times.Time.sample(frac=1,replace=False,random_state=i)
            7
            8      rand_mean_a = session_times_permuted.iloc[:n_a].mean()
            9      rand_mean_b = session_times_permuted.iloc[n_a:].mean()
           10
           11     rand_mean_ad_diffs.append(rand_mean_a - rand_mean_b)
           12
           13  rand_mean_ad_diffs = np.array(rand_mean_ad_diffs)
           14  rand_mean_ad_diffs[:5].round(2)
```

100%                                    10000/10000 [00:05<00:00, 2128.79it/s]

Out[43]: array([  9.79, -13.71, -15.83,  34.99,  -4.67])

# Ex: Websites and Sales, Permutation Test Continued

# Ex: Websites and Sales, Permutation Test Continued

```python
1  # 6. see where our original observation falls
2  fig,ax = plt.subplots(1,1,figsize=(12,8))
3  ax = sns.histplot(x=rand_mean_ad_diffs, stat='density')
4  ax.set_xlabel('random mean differences');ax.set_ylabel('frequency');
5  ax.axvline(observed_ad_metric, color='r');
```

# Ex: Websites and Sales, Permutation Test Continued

# Ex: Websites and Sales, Permutation Test Continued

```
In [45]:   1  # Normalize our values
           2  rand_mean_ad_diffs_xbar = np.mean(rand_mean_ad_diffs)
           3  rand_mean_ad_diffs_s    = np.std(rand_mean_ad_diffs)
           4
           5  rand_mean_ad_zscores = (rand_mean_ad_diffs - rand_mean_ad_diffs_xbar) / rand_mean_ad_diffs_s
           6  list(zip(rand_mean_ad_diffs[:3].round(2),rand_mean_ad_zscores[:3].round(2)))
```

Out[45]:  [(9.79, 0.5), (-13.71, -0.69), (-15.83, -0.8)]

# Ex: Websites and Sales, Permutation Test Continued

In [45]:
```python
# Normalize our values
rand_mean_ad_diffs_xbar = np.mean(rand_mean_ad_diffs)
rand_mean_ad_diffs_s    = np.std(rand_mean_ad_diffs)

rand_mean_ad_zscores = (rand_mean_ad_diffs - rand_mean_ad_diffs_xbar) / rand_mean_ad_diffs_s
list(zip(rand_mean_ad_diffs[:3].round(2),rand_mean_ad_zscores[:3].round(2)))
```

Out[45]: [(9.79, 0.5), (-13.71, -0.69), (-15.83, -0.8)]

In [46]:
```python
observed_ad_metric_zscore =  (observed_ad_metric - rand_mean_ad_diffs_xbar) / rand_mean_ad_diffs_s
observed_ad_metric.round(2),observed_ad_metric_zscore.round(2)
```

Out[46]: (-49.77, -2.52)

# Ex: Websites and Sales, Permutation Test Continued

# Ex: Websites and Sales, Permutation Test Continued

```python
# 6. see where our original observation falls (normalized)
fig,ax = plt.subplots(1,1,figsize=(12,8))
ax = sns.histplot(rand_mean_ad_zscores, stat='density')
ax.set_xlabel('random mean differences normed');ax.set_ylabel('frequency');
ax.axvline(observed_ad_metric_zscore,color='r');
```

# How sure are we?

- **p-value**

*The probability of finding the observed result, or one more extreme, when the null hypothesis ($H_0$) is true.*

# How sure are we?

- **p-value**

*The probability of finding the observed result, or one more extreme, when the null hypothesis ($H_0$) is true.*

- **does mean** $: P \, (\text{data} \mid H_0 \text{ is true})$

# How sure are we?

- **p-value**

*The probability of finding the observed result, or one more extreme, when the null hypothesis ($H_0$) is true.*

- **does mean** : $P(\text{data} \mid H_0 \text{ is true})$

- **does NOT mean** : $P(H_0 \text{ is not true} \mid \text{data})$

# How sure are we?

- **p-value**

> *The probability of finding the observed result, or one more extreme, when the null hypothesis ($H_0$) is true.*

- **does mean** : $P\,(\text{data} \mid H_0 \text{ is true})$

- **does NOT mean** : $P\,(H_0 \text{ is not true} \mid \text{data})$

- Our question about significance becomes:

> *"How often did we see a value as or more extreme than our observed metric?"*

# Choosing One-Tailed vs Two-Tailed

- Do we have a strong reason for a one-directional question? One-Tailed
  - Ex: $H_0$ is "difference is less than or equal to 0"
  - Need a strong reason

- Otherwise? Two-tailed
  - Ex: $H_0$ is "there is no real difference between groups"
  - More conservative
  - Usually a better choice

# Calculating p for Two-Tailed Test

# Calculating p for Two-Tailed Test

```
In [48]:  1  # find absolute values greater than our observed_metric
          2  ad_gt = np.abs(rand_mean_ad_diffs) >= np.abs(observed_ad_metric)
```

# Calculating p for Two-Tailed Test

```
In [48]:  1  # find absolute values greater than our observed_metric
          2  ad_gt = np.abs(rand_mean_ad_diffs) >= np.abs(observed_ad_metric)
```

```
In [49]:  1  # how many are greater than or equal to?
          2  num_ad_gt = ad_gt.sum()
          3
          4  # proportion of total that are as or more extreme
          5  p = num_ad_gt / len(rand_mean_ad_diffs)
          6  print(f'{p = :}')
```

```
p = 0.0078
```

# One-Tailed Test Example

# One-Tailed Test Example

```
In [50]:  1  # one-tailed test
          2  sum(np.array(rand_mean_ad_diffs) <= observed_ad_metric) / len(rand_mean_ad_diffs)

Out[50]:  0.0037
```

# One-Tailed Test Example

```
In [50]:  1  # one-tailed test
          2  sum(np.array(rand_mean_ad_diffs) <= observed_ad_metric) / len(rand_mean_ad_diffs)

Out[50]: 0.0037
```

Note that this is less than our Two-Tailed value!

# One-Tailed Test Example

```
In [50]:   1  # one-tailed test
           2  sum(np.array(rand_mean_ad_diffs) <= observed_ad_metric) / len(rand_mean_ad_diffs)
```

Out[50]:  0.0037

Note that this is less than our Two-Tailed value!

```
In [51]:   1  # two-tailed test
           2  sum(np.abs(rand_mean_ad_diffs) >= np.abs(observed_ad_metric)) / len(rand_mean_ad_diffs)
```

Out[51]:  0.0078

# Equation Based Hypothesis Test: t-Test

# Equation Based Hypothesis Test: t-Test

- based on the Student-t distribution

- more involved to describe

- works for numeric data (can't use it for the next example)

# Equation Based Hypothesis Test: t-Test

- based on the Student-t distribution

- more involved to describe

- works for numeric data (can't use it for the next example)

```python
In [52]:  1  # using our session_times example
          2  t = sp.stats.ttest_ind(session_times[session_times.Page == 'Page A'].Time.values,
          3                         session_times[session_times.Page == 'Page B'].Time.values,
          4                         equal_var=False)
          5  print(f'{t.pvalue = :0.3f}')

t.pvalue = 0.010
```

# Equation Based Hypothesis Test: t-Test

- based on the Student-t distribution

- more involved to describe

- works for numeric data (can't use it for the next example)

```
In [52]:  1  # using our session_times example
          2  t = sp.stats.ttest_ind(session_times[session_times.Page == 'Page A'].Time.values,
          3                         session_times[session_times.Page == 'Page B'].Time.values,
          4                         equal_var=False)
          5  print(f'{t.pvalue = :0.3f}')

          t.pvalue = 0.010
```

- close to the 0.008 value we found via permutation test

# Choosing $\alpha$

- **alpha** ($\alpha$): significance level

    - What we compare our p-value to

    - Best to choose this before calculating metrics

    - Probability of rejecting the null when it is true (Type I Error)

- Common values:

    - .1 (Error 1 out of 10 times)

    - .05 (Error 1 out of 20 times)

    - .01 (Error 1 out of 100 times)

- Should depend on how bad a Type I (False Positive) Error is

# Another Example: Price vs Conversion

- Does Price A lead to higher conversions than Price B?

- **Conversion:** Turning a visit into a sale

- $H_0$: conversions for Price A ≤ conversions for Price B

  - Price A does not lead to more conversions

- $H_1$: conversions for Price A > conversions for Price B

  - Price A leads to more conversions

# Another Example: Price vs Conversion

- Does Price A lead to higher conversions than Price B?

- **Conversion:** Turning a visit into a sale

- $H_0$: conversions for Price A ≤ conversions for Price B

  - Price A does not lead to more conversions

- $H_1$: conversions for Price A > conversions for Price B

  - Price A leads to more conversions

```
In [53]:   1  # Counts of observations
           2  df = pd.DataFrame({'Price A':[200,23539],
           3                     'Price B':[182,22406]},
           4                    index=['Conversion','No Conversion'])
           5  df
```

Out[53]:

|               | Price A | Price B |
|---------------|---------|---------|
| **Conversion**    | 200     | 182     |
| **No Conversion** | 23539   | 22406   |

# Another Example: Price vs Conversion Continued

- Metric of Interest?
  - difference in percent conversion

# Another Example: Price vs Conversion Continued

- Metric of Interest?
  - difference in percent conversion

```
In [54]:   1  pct_conv = df.loc['Conversion'] / df.sum(axis=0) * 100
           2  pct_conv.round(2)

Out[54]:   Price A    0.84
           Price B    0.81
           dtype: float64
```

# Another Example: Price vs Conversion Continued

- Metric of Interest?
  - difference in percent conversion

In [54]:
```
1  pct_conv = df.loc['Conversion'] / df.sum(axis=0) * 100
2  pct_conv.round(2)
```

Out[54]:
```
Price A    0.84
Price B    0.81
dtype: float64
```

In [55]:
```
1  diff_pct_conv = pct_conv['Price A'] - pct_conv['Price B']
2  print(f'{diff_pct_conv.round(3)}%')
```

```
0.037%
```

# Another Example: Price vs Conversion Continued

- First: Choose our $\alpha$: 0.05

- Reminder of Permutation Test:

  0. get group sizes

  1. combine groups together

  2. permute observations

  3. create two new groups (same sizes as originals)

  4. calculate metric

  5. repeat many times

  6. see where our original observation falls

# Another Example: Price vs Conversion Continued

- What are our samples?

    - 1 = Conversion

    - 0 = No conversion

- How many samples are there?

# Another Example: Price vs Conversion Continued

- What are our samples?

  - 1 = Conversion
  - 0 = No conversion

- How many samples are there?

```
In [56]:   1  n = df.sum().sum()
           2  n

Out[56]:   46327
```

# Another Example: Price vs Conversion Continued

- Turning counts into samples

# Another Example: Price vs Conversion Continued

- Turning counts into samples

```
In [57]:  1  n_conversion = df.loc['Conversion'].sum()
          2  n_conversion
          3

Out[57]:  382
```

# Another Example: Price vs Conversion Continued

- Turning counts into samples

In [57]:
```python
1  n_conversion = df.loc['Conversion'].sum()
2  n_conversion
3
```

Out[57]: 382

In [58]:
```python
1  conv_samples = np.zeros(n)
2  conv_samples[:n_conversion] = 1
3
4  assert sum(conv_samples) == n_conversion
```

# Another Example: Price vs Conversion Continued

- Turning counts into samples

```
In [57]:  1  n_conversion = df.loc['Conversion'].sum()
          2  n_conversion
          3
```

```
Out[57]:  382
```

```
In [58]:  1  conv_samples = np.zeros(n)
          2  conv_samples[:n_conversion] = 1
          3
          4  assert sum(conv_samples) == n_conversion
```

```
In [59]:  1  conv_samples
```

```
Out[59]:  array([1., 1., 1., ..., 0., 0., 0.])
```

# Another Example: Price vs Conversion Continued

```
In [60]:   1  df
```

Out[60]:

|  | Price A | Price B |
|---|---|---|
| **Conversion** | 200 | 182 |
| **No Conversion** | 23539 | 22406 |

# Another Example: Price vs Conversion Continued

In [60]:
```
1 df
```

Out[60]:

|  | Price A | Price B |
|---|---|---|
| **Conversion** | 200 | 182 |
| **No Conversion** | 23539 | 22406 |

In [61]:
```
1 n_pricea, n_priceb = df.sum(axis=0)
2 print(f'{n_pricea=} {n_priceb=} {n=}')
3
4 assert n_pricea + n_priceb == n
```

```
n_pricea=23739 n_priceb=22588 n=46327
```

# Another Example: Price vs Conversion Continued

In [60]:
```
1 df
```

Out[60]:

|  | Price A | Price B |
|---|---|---|
| **Conversion** | 200 | 182 |
| **No Conversion** | 23539 | 22406 |

In [61]:
```
1 n_pricea, n_priceb = df.sum(axis=0)
2 print(f'{n_pricea=} {n_priceb=} {n=}')
3
4 assert n_pricea + n_priceb == n
```

n_pricea=23739 n_priceb=22588 n=46327

In [62]:
```
1 np.random.seed(123)
2 rand_conv_diffs = []
3 for i in tqdm(range(1000)):
4     conv_permuted = np.random.permutation(conv_samples)
5     rand_conv_a = sum(conv_permuted[:n_pricea]) / n_pricea
6     rand_conv_b = sum(conv_permuted[n_pricea:]) / n_priceb
7     rand_conv_diffs.append(100 * (rand_conv_a - rand_conv_b))
```
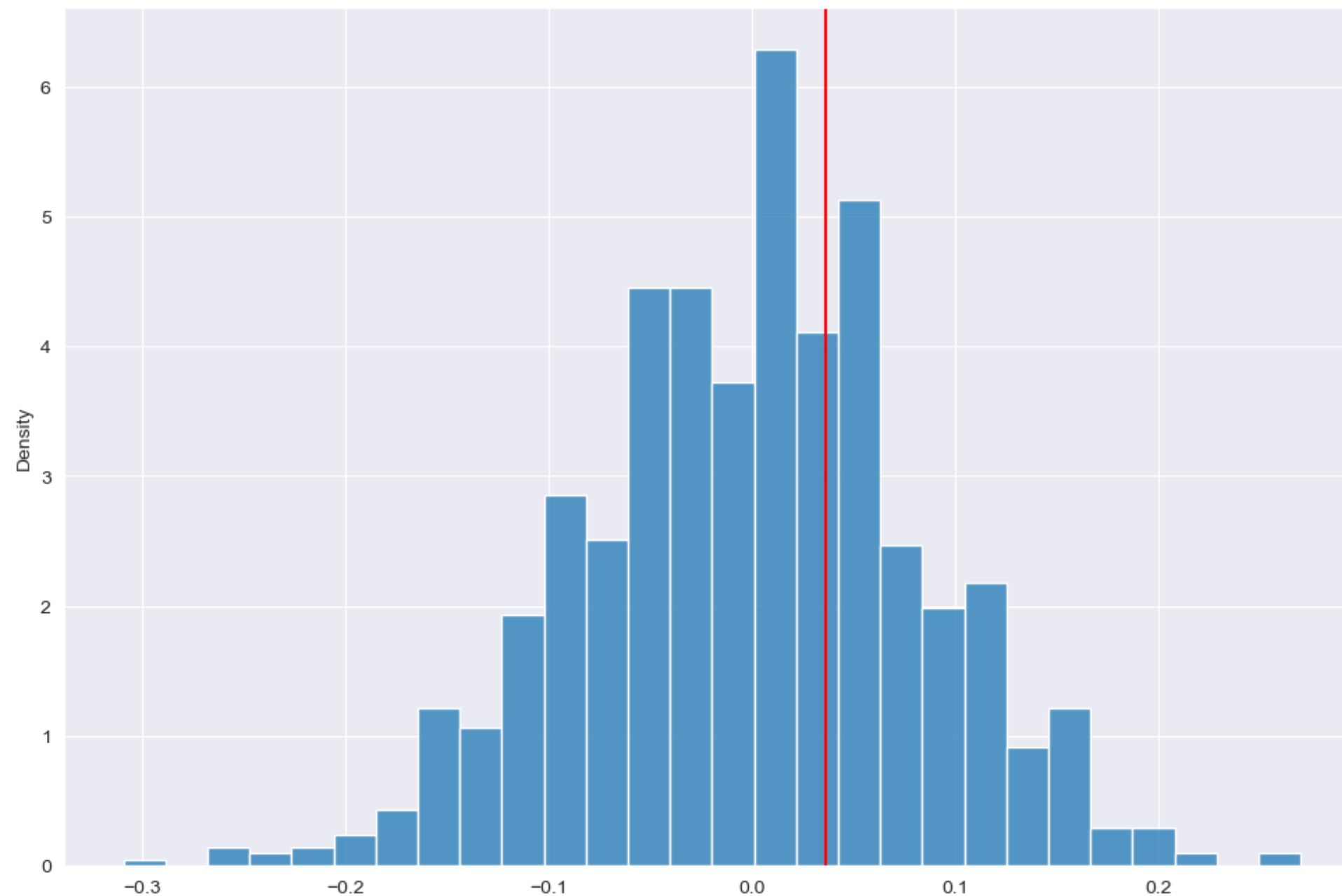
100%              1000/1000 [00:03<00:00, 282.07it/s]

# Another Example: Price vs Conversion Continued

# Another Example: Price vs Conversion Continued

```
In [63]:  1  fig,ax = plt.subplots(1,1,figsize=(12,8))
          2  ax = sns.histplot(x=rand_conv_diffs, stat='density')
          3  ax.axvline(diff_pct_conv,color='r');
```
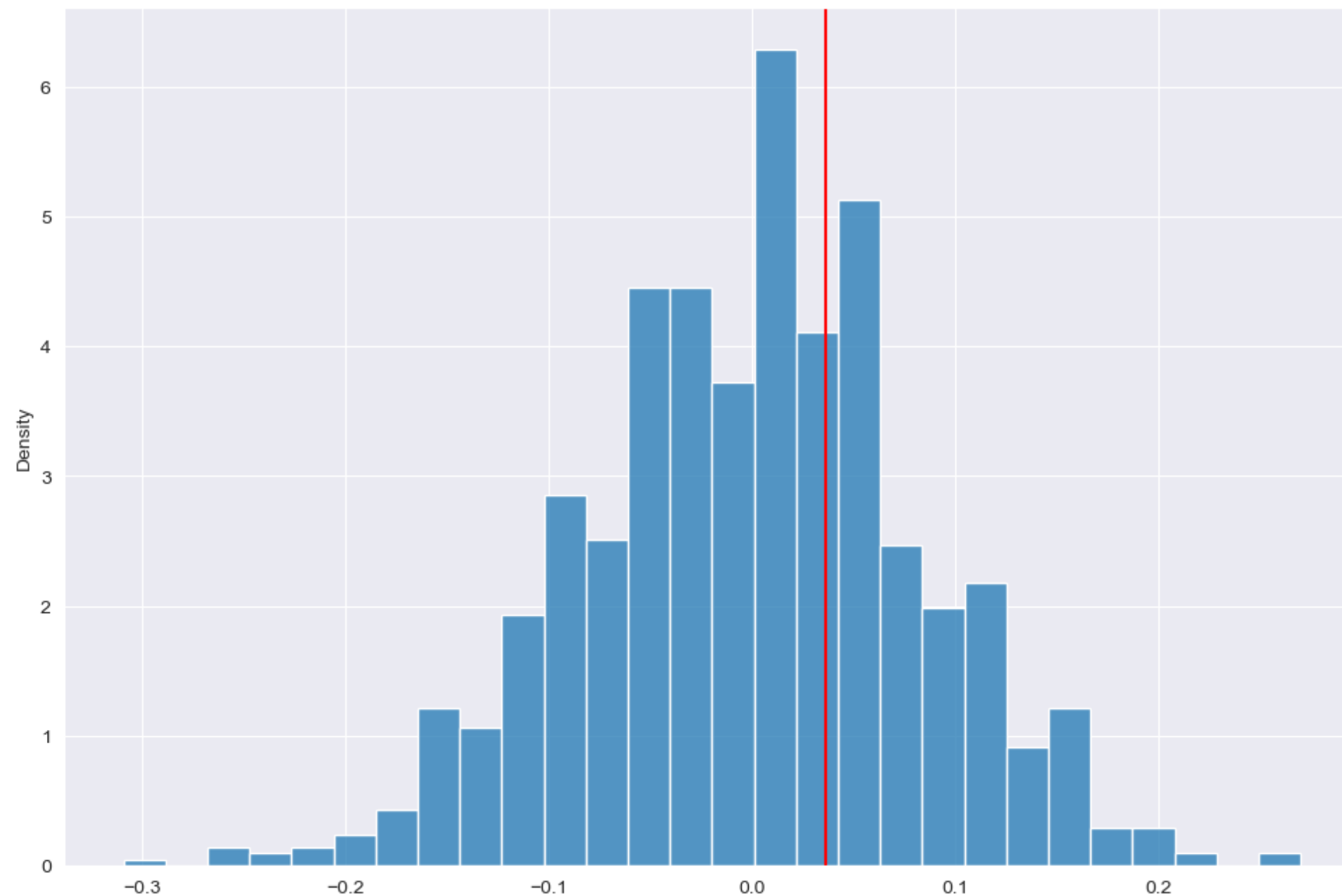
# Another Example: Price vs Conversion Continued

```
In [63]:   1  fig,ax = plt.subplots(1,1,figsize=(12,8))
           2  ax = sns.histplot(x=rand_conv_diffs, stat='density')
           3  ax.axvline(diff_pct_conv,color='r');
```



```
In [64]:   1  # calculate a two-tailed p-value
```

# Equation Based Proportion Test

# Equation Based Proportion Test

```
In [65]:  1  from statsmodels.stats.proportion import proportions_ztest
          2
          3  z,p = proportions_ztest(df.loc['Conversion'].values,
          4                          df.sum(),
          5                          alternative='two-sided')
          6  print(f'{p = :0.3f}')

          p = 0.662
```

# Statistically Significant?

- Don't base your conclusions solely on whether an association or effect was found to be "statistically significant" (i.e., the p-value passed some arbitrary threshold such as p < 0.05).

- Don't believe an association/effect **exists** just because it **was statistically significant**.

- Don't believe an association/effect **is absent** just because it **was not stat. significant**.

- Don't believe that your p-value:
    1. gives the **probability that chance alone** produced the observed association/effect or
    2. the probability that your **test hypothesis is true**.]

- Don't conclude anything about **scientific or practical importance** based on statistical significance (or lack thereof).

# Statistically Significant?

- **<u>Moving to a World Beyond "p < 0.05"</u>** Wasserstein, Schirm & Lazar, 20 Mar 2019

    - Try to avoid "Statistically Significant"

    - "**A**ccept uncertainty. Be **t**houghtful, **o**pen, and **m**odest." Remember "**ATOM**."

# Statistically Significant?

- [Moving to a World Beyond "p < 0.05"](#) Wasserstein, Schirm & Lazar, 20 Mar 2019

  - Try to avoid "Statistically Significant"

  - "**A**ccept uncertainty. Be **t**houghtful, **o**pen, and **m**odest." Remember "**ATOM**."

- ATOM

  - **A:** Seek better measures, more sensitive designs, larger samples

  - **T:** Begin with clearly expressed objectives

  - **T:** Ask "What are the practical implications?"

  - **O::** Be open/transparent in analysis and communication

  - **M:** Accept limititaions, assumptions, reproduction, recognizing differences in stakes

# Issues with Multiple Testing

- p-hacking: keep trying comparisons till you find something that works

- multiple tests: the more tests you run, the more likely a Type 1 Error

- One simple solution:
  - **Bonferonni correction** $\frac{\alpha}{m}$ where $m$ is the number of tests

# Comparing More Than 2 Groups

- ANOVA

  - need more stats than we have time for

- Multi-Armed Bandit (MAB)

  - can compare many distributions
  - don't need to make assumptions about underlying distributions
  - can also be used for early stopping of experiment

# Multi-Armed Bandit



**Question:** Which arm should we pull?
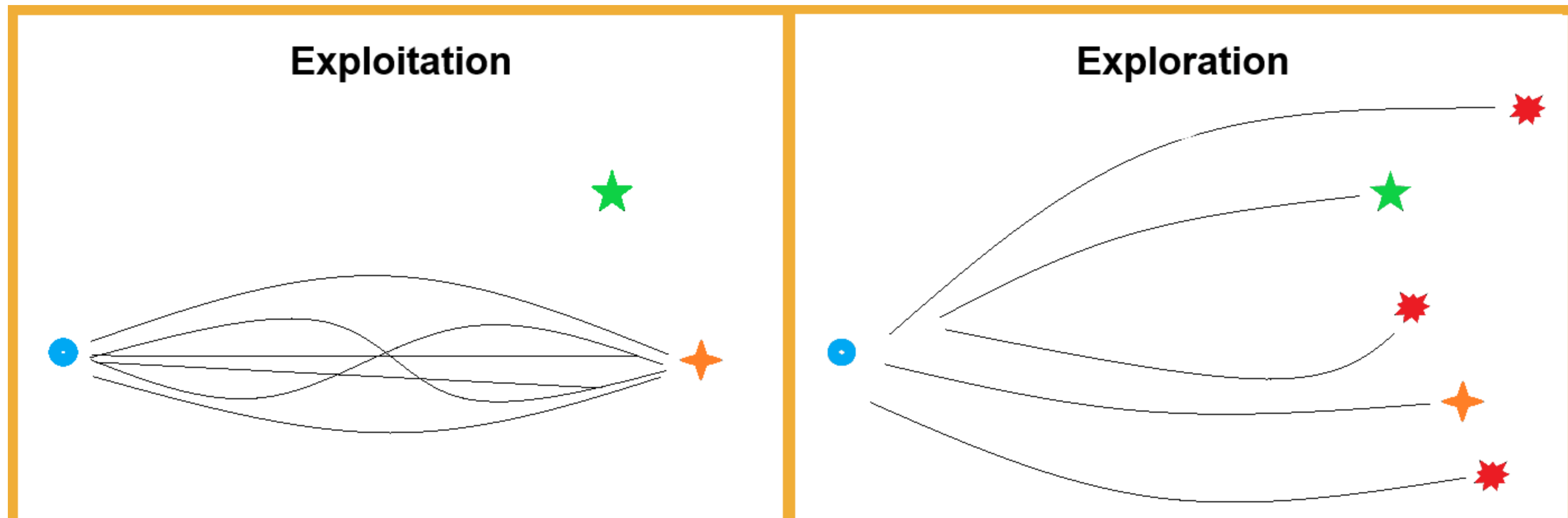
# Greedy MAB

**greedy:** do something simple that heads towards the goal

    1. pull arm with highest payout

But what if there's a better choice, we just haven't seen it yet?

# Exploration Vs Exploitation

- **Exploration:** There might be a better arm

  - keep choosing different arms randomly

- **Exploitation:** We want to make use of the best

  - keep pulling the best arm

# $\epsilon$-Greedy MAB

- choose a small epsilon ($\epsilon$) between 0 and 1

# $\epsilon$-Greedy MAB

- choose a small epsilon ($\epsilon$) between 0 and 1

      1. generate random number between 0 and 1

      2. if < $\epsilon$, choose arm randomly

      3. if ≥ $\epsilon$, choose best arm

      4. GOTO 1

- larger $\epsilon$ -> more exploration

# MAB Example

- We have three ads
- We don't know how often each will lead to a response
- We need to decide which ad to add to each page request

# MAB Example

- We have three ads

- We don't know how often each will lead to a response

- We need to decide which ad to add to each page request

```
In [66]:   1  # creating three ads (distributions) with unknown response rate
           2  np.random.seed(13)
           3  ad_A = sp.stats.bernoulli(p=np.random.rand())
           4  ad_B = sp.stats.bernoulli(p=np.random.rand())
           5  ad_C = sp.stats.bernoulli(p=np.random.rand())
```

# MAB Example

- We have three ads

- We don't know how often each will lead to a response

- We need to decide which ad to add to each page request

```
In [66]:   1  # creating three ads (distributions) with unknown response rate
           2  np.random.seed(13)
           3  ad_A = sp.stats.bernoulli(p=np.random.rand())
           4  ad_B = sp.stats.bernoulli(p=np.random.rand())
           5  ad_C = sp.stats.bernoulli(p=np.random.rand())
```

- We'll use an $\epsilon$-greedy MAB to decide which ad to show

# MAB Example

- We have three ads

- We don't know how often each will lead to a response

- We need to decide which ad to add to each page request

```
In [66]:  1  # creating three ads (distributions) with unknown response rate
          2  np.random.seed(13)
          3  ad_A = sp.stats.bernoulli(p=np.random.rand())
          4  ad_B = sp.stats.bernoulli(p=np.random.rand())
          5  ad_C = sp.stats.bernoulli(p=np.random.rand())
```

- We'll use an $\epsilon$-greedy MAB to decide which ad to show

```
In [67]:  1  # epsilon probability
          2  epsilon = 0.40
```

# MAB Example Continued

# MAB Example Continued

- Rounds 1,2,3
  - Pull each arm once

# MAB Example Continued

- Rounds 1,2,3
  - Pull each arm once

```
In [68]:  1  rewards_A = [ad_A.rvs()]
          2
          3  rewards_B = [ad_B.rvs()]
          4
          5  rewards_C = [ad_C.rvs()]
          6
          7  rewards_A,rewards_B,rewards_C

Out[68]:  ([0], [1], [1])
```

# MAB Example Continued

# MAB Example Continued

- Round 3
  - With probability $1 - \epsilon$, choose the best arm (randomly if tied)

# MAB Example Continued

- Round 3
  - With probability $1 - \epsilon$, choose the best arm (randomly if tied)

```
In [69]:   1  be_greedy = np.random.rand() > epsilon
           2  be_greedy

Out[69]:   True
```

# MAB Example Continued

- Round 3
  - With probability $1 - \epsilon$, choose the best arm (randomly if tied)

```
In [69]:  1 be_greedy = np.random.rand() > epsilon
          2 be_greedy
```

Out[69]:  True

```
In [70]:  1 best_arms = ['B','C']
          2 best_arms[np.random.randint(2)]
          3 # np.random.randint(2) randomly choose 0 or 1
```

Out[70]:  'B'

# MAB Example Continued

- Round 3
    - With probability $1 - \epsilon$, choose the best arm (randomly if tied)

```
In [69]:   1  be_greedy = np.random.rand() > epsilon
           2  be_greedy
```

Out[69]:   True

```
In [70]:   1  best_arms = ['B','C']
           2  best_arms[np.random.randint(2)]
           3  # np.random.randint(2) randomly choose 0 or 1
```

Out[70]:   'B'

```
In [71]:   1  rewards_B.append(ad_B.rvs())
           2
           3  rewards_A,rewards_B,rewards_C
```

Out[71]:   ([0], [1, 1], [1])

# MAB Example Continued

# MAB Example Continued

```python
1  def mab(arms = [],rewards = [],arm_names = [],epsilon=0.4):
2      n_arms = len(arms)
3  #      if not rewards:
4  #          for i in range(n_arms):
5  #              pulls.append(list)
6      be_greedy = np.random.rand() > epsilon
7      if not be_greedy: # randomly choose
8          arm_idx = np.random.randint(0,n_arms)
9          rewards[arm_idx].append(arms[arm_idx].rvs())
10     else: # be greedy
11         reward_means = np.array([sum(x)/len(x) for x in rewards])
12         best_arms = np.where(reward_means == np.amax(reward_means))[0]
13         arm_idx = best_arms[np.random.randint(0,best_arms.shape[0])]
14         rewards[arm_idx].append(arms[arm_idx].rvs())
15     return rewards, be_greedy, arm_names[arm_idx]
16
17 def print_mab_results(be_greedy,choice,rewards):
18     print(f'greedy:{str(be_greedy):5s} choice:{choice} => '+
19         f"[{':'.join([str(round(sum(x)/len(x),1)) for x in rewards])}] "+
20         f"| {','.join([str(x).ljust(20,' ') for x in rewards])}")
```

# MAB Example Continued

# MAB Example Continued

- Round 4

```
In [73]:  1  arms = [ad_A,ad_B,ad_C]
          2  rewards = [rewards_A,rewards_B,rewards_C]
          3  labels = ['A','B','C']
          4
          5  rewards, be_greedy, choice = mab(arms,rewards,labels,epsilon)
          6
          7  print_mab_results(be_greedy,choice,rewards)
```

```
greedy:False choice:C => [0.0:1.0:1.0] | [0]                ,[1, 1]                ,[1, 1]
```

# MAB Example Continued

# MAB Example Continued

```
In [74]:  1  for i in range(10):
          2      rewards, be_greedy, choice = mab(arms,rewards,labels,epsilon)
          3      print_mab_results(be_greedy,choice,rewards)
```

```
greedy:False choice:A => [0.5:1.0:1.0] | [0, 1]           ,[1, 1]              ,[1, 1]
greedy:True  choice:B => [0.5:0.7:1.0] | [0, 1]           ,[1, 1, 0]           ,[1, 1]
greedy:False choice:A => [0.3:0.7:1.0] | [0, 1, 0]        ,[1, 1, 0]           ,[1, 1]
greedy:True  choice:C => [0.3:0.7:1.0] | [0, 1, 0]        ,[1, 1, 0]           ,[1, 1, 1]
greedy:True  choice:C => [0.3:0.7:0.8] | [0, 1, 0]        ,[1, 1, 0]           ,[1, 1, 1, 0]
greedy:True  choice:C => [0.3:0.7:0.8] | [0, 1, 0]        ,[1, 1, 0]           ,[1, 1, 1, 0, 1]
greedy:True  choice:C => [0.3:0.7:0.8] | [0, 1, 0]        ,[1, 1, 0]           ,[1, 1, 1, 0, 1, 1]
greedy:False choice:B => [0.3:0.5:0.8] | [0, 1, 0]        ,[1, 1, 0, 0]        ,[1, 1, 1, 0, 1, 1]
greedy:True  choice:C => [0.3:0.5:0.9] | [0, 1, 0]        ,[1, 1, 0, 0]        ,[1, 1, 1, 0, 1, 1, 1]
greedy:False choice:C => [0.3:0.5:0.9] | [0, 1, 0]        ,[1, 1, 0, 0]        ,[1, 1, 1, 0, 1, 1, 1, 1]
```

# MAB Example Continued

# MAB Example Continued

- Which arm seems best?

# MAB Example Continued

- Which arm seems best?

```python
In [75]:  1  rates = ' '.join([f"{label}:{np.mean(reward).round(1)}" for label,reward in zip(labels,rewards)])
          2  print(f"conversion rates: {rates}")

conversion rates: A:0.3 B:0.5 C:0.9
```

# MAB Example Continued

- Which arm seems best?

```
In [75]:  1  rates = ' '.join([f"{label}:{np.mean(reward).round(1)}" for label,reward in zip(labels,rewards)])
          2  print(f"conversion rates: {rates}")

conversion rates: A:0.3 B:0.5 C:0.9
```

- Did we pick the best one?

# MAB Example Continued

- Which arm seems best?

In [75]:
```python
1  rates = ' '.join([f"{label}:{np.mean(reward).round(1)}" for label,reward in zip(labels,rewards)])
2  print(f"conversion rates: {rates}")
```

conversion rates: A:0.3 B:0.5 C:0.9

- Did we pick the best one?

In [76]:
```python
1  ground_truth_rates = ' '.join([f"{label}:{arm.pmf(1).round(1)}" for label,arm in zip(labels,arms)])
2  print(f'ground truth: {ground_truth_rates}')
```

ground truth: A:0.8 B:0.2 C:0.8

# MAB Variations

- Thompson's Sampling: uses Baysian approach

- UCB1: maximize expected reward using Upper Confidence Bounds

- ...

# Questions?