# Statistical Thinking in Python Part 2

## 3). <u>Introduction to hypothesis testing</u>

**a).** <u>Generating a Permutation Sample:</u>

```
def permutation_sample(data1, data2):
    """Generate a permutation sample from two data sets."""

    # Concatenate the data sets: data
    data = np.concatenate((data1,data2))

    # Permute the concatenated array: permuted_data
    permuted_data = np.random.permutation(data)

    # Split the permuted array into two: perm_sample_1, perm_sample_2
    perm_sample_1 = permuted_data[:len(data1)]
    perm_sample_2 = permuted_data[len(data1):]

    return perm_sample_1, perm_sample_2
```

**b).** <u>Visualizing Permutation Sample:</u>

```
for x in range(50):

    # Generate permutation samples

    perm_sample_1, perm_sample_2 = permutation_sample(rain_june,rain_november)



    # Compute ECDFs

    x_1, y_1 = ecdf(perm_sample_1)

    x_2, y_2 = ecdf(perm_sample_2)
```

```
    # Plot ECDFs of permutation sample
  _ = plt.plot(x_1, y_1, marker='.', linestyle='none',
         color='red', alpha=0.02)
  _ = plt.plot(x_2, y_2, marker='.', linestyle='none',
         color='blue', alpha=0.02)


# Create and plot ECDFs from original data
x_1, y_1 = ecdf(rain_june)
x_2, y_2 = ecdf(rain_november)
_ = plt.plot(x_1, y_1, marker='.', linestyle='none', color='red')
_ = plt.plot(x_2, y_2, marker='.', linestyle='none', color='blue')


# Label axes, set margin, and show plot
plt.margins(0.02)
_ = plt.xlabel('monthly rainfall (mm)')
_ = plt.ylabel('ECDF')
plt.show()
```
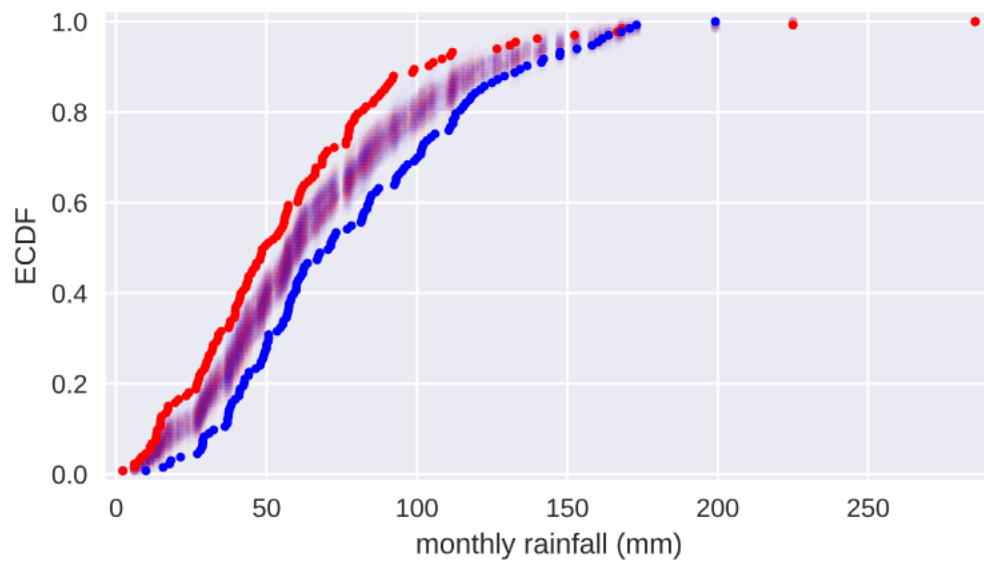
**c). <u>Generating Permutation Replicates</u>**

```python
def draw_perm_reps(data_1, data_2, func, size=1):
    """Generate multiple permutation replicates."""

    # Initialize array of replicates: perm_replicates
    perm_replicates = np.empty(size)

    for i in range(size):
        # Generate permutation sample
        perm_sample_1, perm_sample_2 = permutation_sample(data_1, data_2)

        # Compute the test statistic
        perm_replicates[i] = func(perm_sample_1, perm_sample_2)

    return perm_replicates
```

**d). <u>Look Before You Leap: EDA Before Hypothesis Testing</u>**

**# Make bee swarm plot**

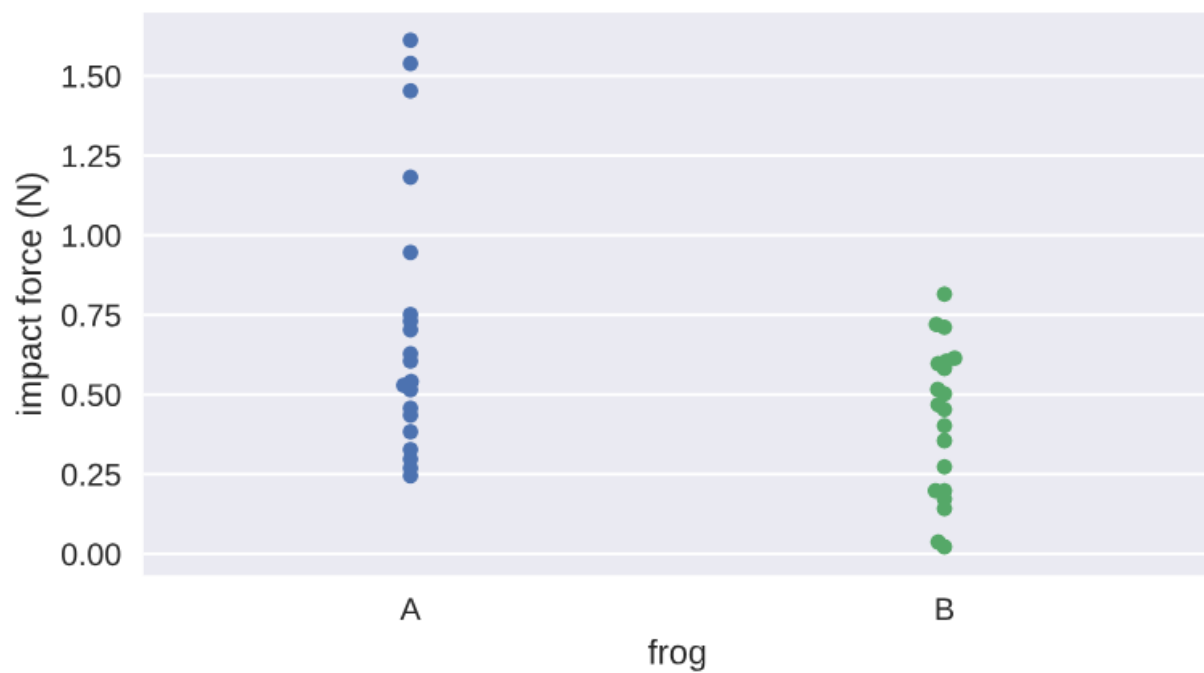**_ = sns.swarmplot(x='ID',y='impact_force',data='df')**

**# Label axes**

**_ = plt.xlabel('frog')**

**_ = plt.ylabel('impact force (N)')**

**# Show the plot**

**plt.show()**

**e).  Permutation Test on Frog data**

```python
def diff_of_means(data_1, data_2):
    """Difference in means of two arrays."""

    # The difference of means of data_1, data_2: diff
    diff = np.mean(data_1)-np.mean(data_2)

    return diff


# Compute difference of mean impact force from experiment: empirical_diff_means
empirical_diff_means = diff_of_means(force_a, force_b)


# Draw 10,000 permutation replicates: perm_replicates
perm_replicates = draw_perm_reps(force_a, force_b,
                 func(diff_of_means), size=10000)


# Compute p-value: p
p = np.sum(perm_replicates >= empirical_diff_means) / len(perm_replicates)


# Print the result
print('p-value =', p)
```

<script.py> output:
    p-value = 0.0063

**f). <u>A one sample bootstrap hypothesis test</u>**

**# Make an array of translated impact forces: translated_force_b**

**translated_force_b = force_b-np.mean(force_b)+0.55**

**# Take bootstrap replicates of Frog B's translated impact forces: bs_replicates**

**bs_replicates = draw_bs_reps(translated_force_b, np.mean, 10000)**

**# Compute fraction of replicates that are less than the observed Frog B force: p**

**p = np.sum(bs_replicates <= np.mean(force_b)) / 10000**

**# Print the p-value**

**print('p = ', p)**

<script.py> output:

  p =  0.0046

**f). <u>A one sample bootstrap hypothesis test</u>**

**g). <u>A bootstrap test for identical distribution</u>**

**# Compute difference of mean impact force from experiment: empirical_diff_means**

**empirical_diff_means = diff_of_means(force_a, force_b)**

**# Concatenate forces: forces_concat**

**forces_concat = np.concatenate((force_a, force_b))**

**# Initialize bootstrap replicates: bs_replicates**

**bs_replicates = np.empty(10000)**

**for i in range(10000):**

   **# Generate bootstrap sample**

   **bs_sample = np.random.choice(forces_concat, size=len(forces_concat))**

   **# Compute replicate**

   **bs_replicates[i] = diff_of_means(bs_sample[:len(force_a)],**

                        **bs_sample[len(force_a):])**

**# Compute and print p-value: p**

**p = np.sum(bs_replicates >= empirical_diff_means) / float(len(bs_replicates))**

**print('p-value =', p)**

<script.py> output:

   p-value = 0.0055

**h). <u>A two sample bootstrap hypothesis test for difference of means</u>**

**# Compute mean of all forces: mean_force**

**mean_force = np.mean(forces_concat)**

**# Generate shifted arrays**

**force_a_shifted = force_a - np.mean(force_a) + mean_force**

**force_b_shifted = force_b - np.mean(force_b) + mean_force**

**# Compute 10,000 bootstrap replicates from shifted arrays**

**bs_replicates_a = draw_bs_reps(force_a_shifted, np.mean, 10000)**

**bs_replicates_b = draw_bs_reps(force_b_shifted, np.mean, 10000)**

**# Get replicates of difference of means: bs_replicates**

**bs_replicates = bs_replicates_a - bs_replicates_b**

**# Compute and print p-value: p**

**p = np.sum(bs_replicates>= empirical_diff_means) / len(bs_replicates)**

**print('p-value =', p)**

&lt;script.py&gt; output:

   p-value = 0.0043