# Statistical Thinking in Python Part 2

## 5). Putting it all together: a case study

**a). EDA of beak depths of Darwin's finches**

**# Create bee swarm plot**

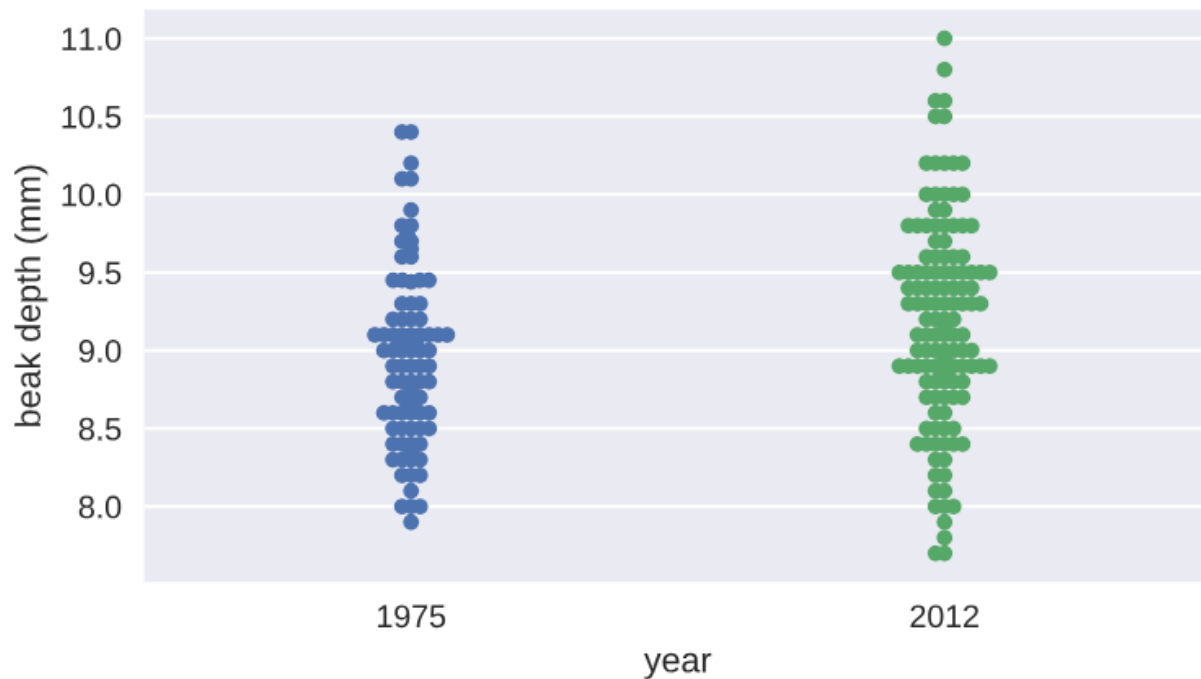**_ = sns.swarmplot(x='year', y='beak_depth',data=df)**


**# Label the axes**

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

**_ = plt.ylabel('beak depth (mm)')**


**# Show the plot**

**plt.show()**

**b). <u>ECDFs of beak depths</u>**

**# Compute ECDFs**

**x_1975, y_1975 = ecdf(bd_1975)**

**x_2012, y_2012 = ecdf(bd_2012)**


**# Plot the ECDFs**

**_ = plt.plot(x_1975, y_1975, marker='.', linestyle='none')**

**_ = plt.plot(x_2012, y_2012, marker='.', linestyle='none')**


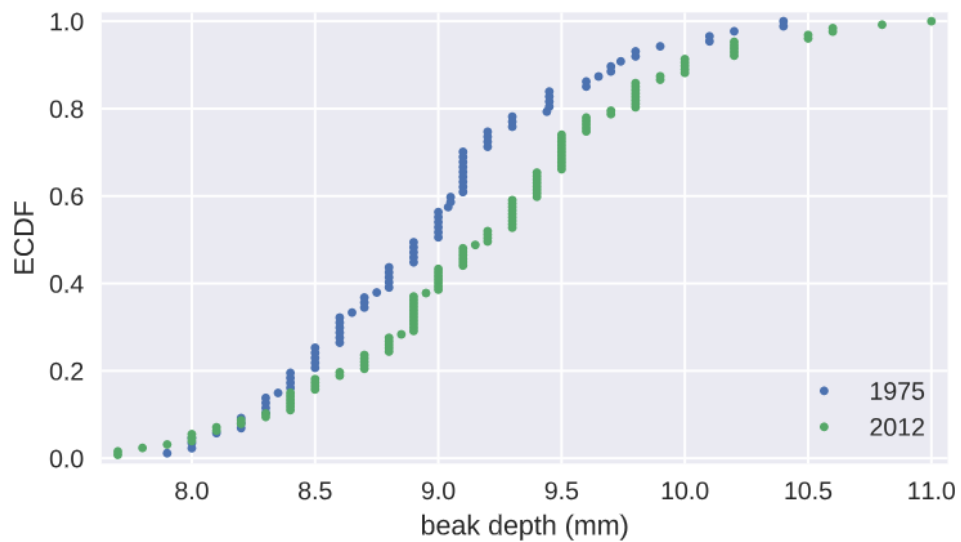**# Set margins**

**plt.margins(0.02)**


**# Add axis labels and legend**

**_ = plt.xlabel('beak depth (mm)')**

**_ = plt.ylabel('ECDF')**

**_ = plt.legend(('1975', '2012'), loc='lower right')**


**# Show the plot**

**plt.show()**

**c). <u>Parameter estimates of beak depths</u>**

**# Compute the difference of the sample means: mean_diff**

**mean_diff = np.mean(bd_2012) - np.mean(bd_1975)**

**# Get bootstrap replicates of means**

**bs_replicates_1975 = draw_bs_reps(bd_1975,np.mean, size=10000)**

**bs_replicates_2012 = draw_bs_reps(bd_2012,np.mean, size=10000)**

**# Compute samples of difference of means: bs_diff_replicates**

**bs_diff_replicates = bs_replicates_2012 - bs_replicates_1975**

**# Compute 95% confidence interval: conf_int**

**conf_int = np.percentile(bs_diff_replicates,[2.5,97.5])**

**# Print the results**

**print('difference of means =', mean_diff, 'mm')**

**print('95% confidence interval =', conf_int, 'mm')**

**<script.py> output:**

   **difference of means = 0.226220472441 mm**

  **95% confidence interval = [ 0.05633521  0.39190544] mm**

**d). <u>Hypothesis test: Are beaks deeper in 2012?</u>**

```python
# Compute mean of combined data set: combined_mean
combined_mean = np.mean(np.concatenate((bd_1975, bd_2012)))


# Shift the samples
bd_1975_shifted = (bd_1975 - np.mean(bd_1975))+combined_mean
bd_2012_shifted = (bd_2012 - np.mean(bd_2012))+combined_mean


# Get bootstrap replicates of shifted data sets
bs_replicates_1975 = draw_bs_reps(bd_1975_shifted, np.mean, 10000)
bs_replicates_2012 = draw_bs_reps(bd_2012_shifted, np.mean, 10000)


# Compute replicates of difference of means: bs_diff_replicates
bs_diff_replicates = bs_replicates_2012 - bs_replicates_1975


# Compute the p-value
p = np.sum(bs_diff_replicates >= mean_diff) / len(bs_diff_replicates)


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

**<script.py> output:**

    **p = 0.0034**

**e). EDA of beak length and depth**

**# Make scatter plot of 1975 data**

**_ = plt.plot(bl_1975, bd_1975, marker='.',**

      **linestyle='none', color='blue', alpha=0.5)**
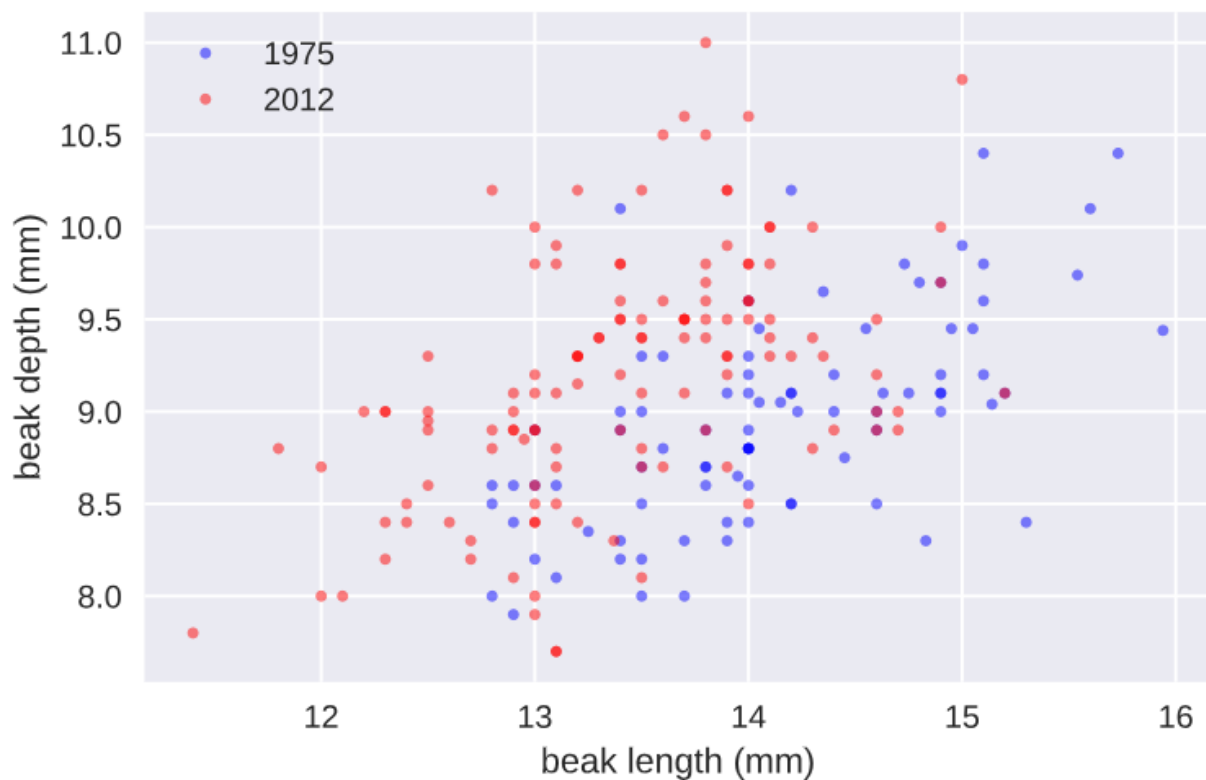
**# Make scatter plot of 2012 data**

**_ = plt.plot(bl_2012, bd_2012, marker='.',**

      **linestyle='none', color='red', alpha=0.5)**

**# Label axes and make legend**

**_ = plt.xlabel('beak length (mm)')**

**_ = plt.ylabel('beak depth (mm)')**

**_ = plt.legend(('1975', '2012'), loc='upper left')**

**# Show the plot**

**plt.show()**

**f). <u>Linear Regression</u>**

**# Compute the linear regressions**

**slope_1975, intercept_1975 = np.polyfit(bl_1975, bd_1975, 1)**

**slope_2012, intercept_2012 = np.polyfit(bl_2012, bd_2012, 1)**

**# Perform pairs bootstrap for the linear regressions**

**bs_slope_reps_1975, bs_intercept_reps_1975 = \**

    **draw_bs_pairs_linreg(bl_1975, bd_1975, 1000)**

**bs_slope_reps_2012, bs_intercept_reps_2012 = \**

    **draw_bs_pairs_linreg(bl_2012, bd_2012, 1000)**

**# Compute confidence intervals of slopes**

**slope_conf_int_1975 = np.percentile(bs_slope_reps_1975,[2.5,97.5])**

**slope_conf_int_2012 = np.percentile(bs_slope_reps_2012,[2.5,97.5])**

**intercept_conf_int_1975 = np.percentile(bs_intercept_reps_1975,[2.5,97.5])**

**intercept_conf_int_2012 = np.percentile(bs_intercept_reps_2012,[2.5,97.5])**

**# Print the results**

**print('1975: slope =', slope_1975,**

    **'conf int =', slope_conf_int_1975)**

**print('1975: intercept =', intercept_1975,**

    **'conf int =', intercept_conf_int_1975)**

**print('2012: slope =', slope_2012,**

    **'conf int =', slope_conf_int_2012)**

**print('2012: intercept =', intercept_2012,**

    **'conf int =', intercept_conf_int_2012)**

<script.py> output:

    1975: slope = 0.465205169161 conf int = [ 0.33851226  0.59306491]

    1975: intercept = 2.39087523658 conf int = [ 0.64892945  4.18037063]

    2012: slope = 0.462630358835 conf int = [ 0.33137479  0.60695527]

    2012: intercept = 2.97724749824 conf int = [ 1.06792753  4.70599387]

**g). Displaying the linear regression result**

**# Make scatter plot of 1975 data**

**_ = plt.plot(bl_1975, bd_1975, marker='.',**

      **linestyle='none', color='blue', alpha=0.5)**

**# Make scatter plot of 2012 data**

**_ = plt.plot(bl_2012, bd_2012, marker='.',**

      **linestyle='none', color='red', alpha=0.5)**

**# Label axes and make legend**

**_ = plt.xlabel('beak length (mm)')**

**_ = plt.ylabel('beak depth (mm)')**

**_ = plt.legend(('1975', '2012'), loc='upper left')**


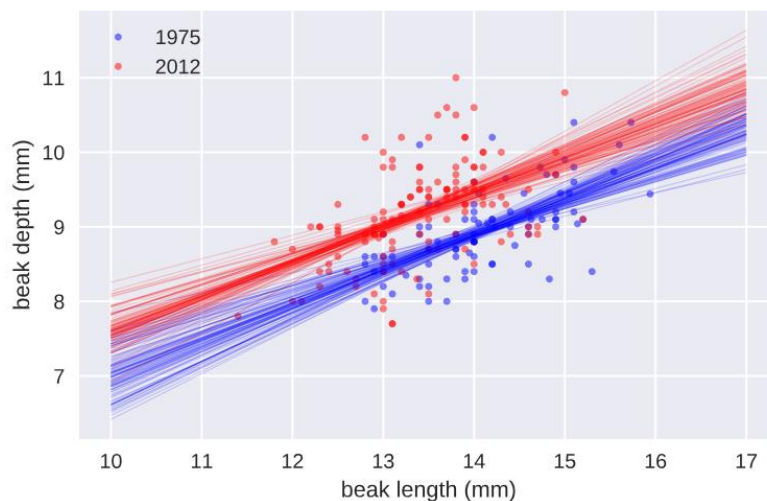**# Generate x-values for bootstrap lines: x**

**x = np.array([10, 17])**

**# Plot the bootstrap lines**

**for i in range(100):**

  **plt.plot(x, bs_slope_reps_1975[i]*x+bs_intercept_reps_1975[i],**

      **linewidth=0.5, alpha=0.2, color='blue')**

  **plt.plot(x, bs_slope_reps_2012[i]*x+bs_intercept_reps_2012[i],**

      **linewidth=0.5, alpha=0.2, color='red')**

**# Draw the plot again**

**plt.show()**

**h). <u>Beak length to deapth ratio</u>**

**# Compute length-to-depth ratios**

**ratio_1975 = bl_1975/ bd_1975**

**ratio_2012 = bl_2012/ bd_2012**

**# Compute means**

**mean_ratio_1975 = np.mean(ratio_1975)**

**mean_ratio_2012 = np.mean(ratio_2012)**

**# Generate bootstrap replicates of the means**

**bs_replicates_1975 = draw_bs_reps(ratio_1975, np.mean, size=10000)**

**bs_replicates_2012 = draw_bs_reps(ratio_2012, np.mean, size=10000)**

**# Compute the 99% confidence intervals**

**conf_int_1975 = np.percentile(bs_replicates_1975, [0.5,99.5])**

**conf_int_2012 = np.percentile(bs_replicates_2012, [0.5,99.5])**

**# Print the results**

**print('1975: mean ratio =', mean_ratio_1975,**

   **'conf int =', conf_int_1975)**

**print('2012: mean ratio =', mean_ratio_2012,**

   **'conf int =', conf_int_2012)**

<script.py> output:

   1975: mean ratio = 1.57888237719 conf int = [ 1.55668803  1.60073509]

   2012: mean ratio = 1.46583422768 conf int = [ 1.44363932  1.48729149]

**i). <u>EDA od hieratibility</u>**

**# Make scatter plots**

**_ = plt.plot(bd_parent_fortis, bd_offspring_fortis,**

     **marker='.', linestyle='none', color='blue', alpha=0.5)**

**_ = plt.plot(bd_parent_scandens, bd_offspring_scandens,**

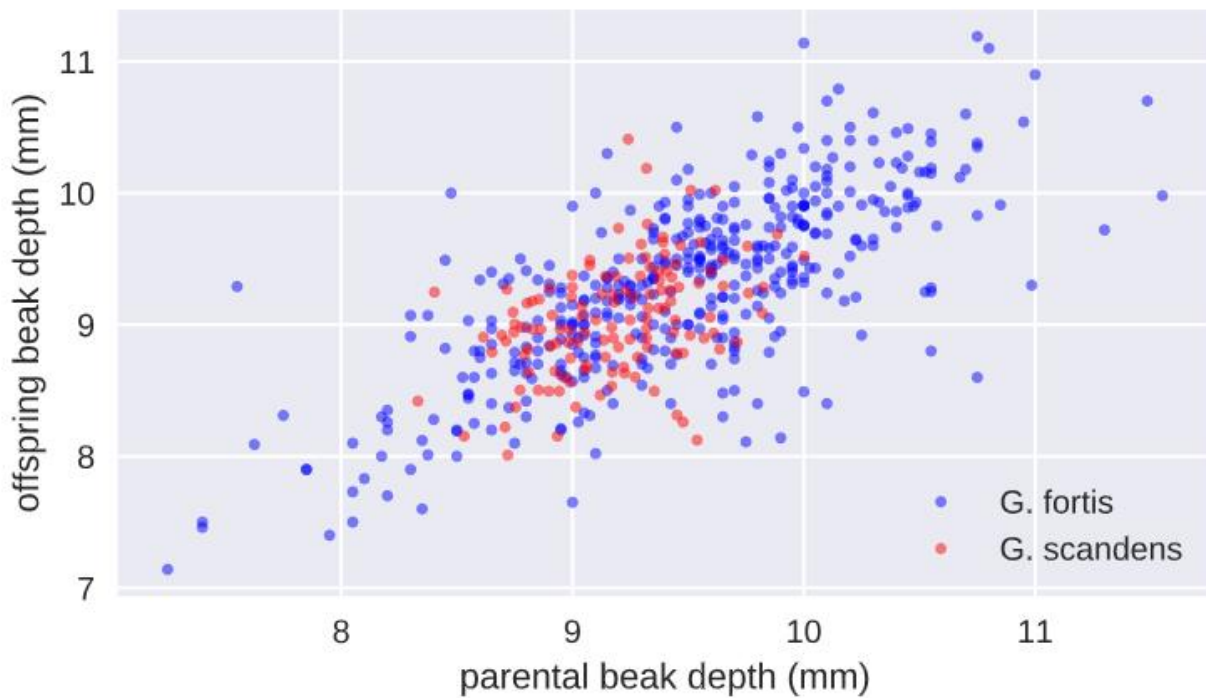     **marker='.', linestyle='none', color='red', alpha=0.5)**

**# Label axes**

**_ = plt.xlabel('parental beak depth (mm)')**

**_ = plt.ylabel('offspring beak depth (mm)')**

**# Add legend**

**_ = plt.legend(('G. fortis', 'G. scandens'), loc='lower right')**

**# Show plot**

**plt.show()**

**j). Corelation of offspring and Parental data**

```python
def draw_bs_pairs(x, y, func, size=1):
    """Perform pairs bootstrap for single statistic."""

    # Set up array of indices to sample from: inds
    inds = np.arange(len(x))

    # Initialize replicates: bs_replicates
    bs_replicates = np.empty(size)

    # Generate replicates
    for i in range(size):
        bs_inds = np.random.choice(inds, len(inds))
        bs_x, bs_y = x[bs_inds], y[bs_inds]
        bs_replicates[i] = func(bs_x, bs_y)

    return bs_replicates
```

**k). <u>Pearson correlation of offspring and parental data</u>**

**# Compute the Pearson correlation coefficients**

**r_scandens = pearson_r(bd_parent_scandens, bd_offspring_scandens)**

**r_fortis = pearson_r(bd_parent_fortis, bd_offspring_fortis)**

**# Acquire 1000 bootstrap replicates of Pearson r**

**bs_replicates_scandens = draw_bs_pairs(bd_parent_scandens, bd_offspring_scandens, pearson_r, 1000)**

**bs_replicates_fortis = draw_bs_pairs(bd_parent_fortis,bd_offspring_fortis, pearson_r, 1000)**

**# Compute 95% confidence intervals**

**conf_int_scandens = np.percentile(bs_replicates_scandens,[2.5,97.5])**

**conf_int_fortis = np.percentile(bs_replicates_fortis,[2.5,97.5])**

**# Print results**

**print('G. scandens:', r_scandens, conf_int_scandens)**

**print('G. fortis:', r_fortis, conf_int_fortis)**

<script.py> output:

   G. scandens: 0.41170636294 [ 0.26564228  0.54388972]

   G. fortis: 0.728341239552 [ 0.6694112   0.77840616]

**l). <u>Measuring Heritability</u>**

```python
def heritability(parents, offspring):
    """Compute the heritability from parent and offspring samples."""
    covariance_matrix = np.cov(parents, offspring)
    return covariance_matrix[0,1] / covariance_matrix[0,0]


# Compute the heritability
heritability_scandens = heritability(bd_parent_scandens, bd_offspring_scandens)
heritability_fortis = heritability(bd_parent_fortis, bd_offspring_fortis)


# Acquire 1000 bootstrap replicates of heritability
replicates_scandens = draw_bs_pairs(bd_parent_scandens, bd_offspring_scandens, heritability, 1000)


replicates_fortis = draw_bs_pairs(bd_parent_fortis, bd_offspring_fortis, heritability, 1000)



# Compute 95% confidence intervals
conf_int_scandens = np.percentile(replicates_scandens, [2.5,97.5])
conf_int_fortis = np.percentile(replicates_fortis, [2.5,97.5])


# Print results
print('G. scandens:', heritability_scandens, conf_int_scandens)
print('G. fortis:', heritability_fortis, conf_int_fortis)
```

&lt;script.py&gt; output:

   G. scandens: 0.548534086869 [ 0.34395487  0.75638267]

   G. fortis: 0.722905191144 [ 0.64655013  0.79688342]

**j). <u>Is beak depth heritable at all in G. scandens?</u>**

**# Initialize array of replicates: perm_replicates**

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

**# Draw replicates**

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

   **# Permute parent beak depths**

   **bd_parent_permuted = np.random.permutation(bd_parent_scandens)**

   **perm_replicates[i] = heritability(bd_parent_permuted, bd_offspring_scandens)**

**# Compute p-value: p**

**p = np.sum(perm_replicates >= heritability_scandens) / len(perm_replicates)**

**# Print the p-value**

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

&lt;script.py&gt; output:

   p-val = 0.0

**j). <u>Is beak depth heritable at all in G. scandens?</u>**