

Investigating Coefficient Standard Errors in Crab Molting Size Data: Using Bootstrap Analysis

Issues:

The dataset used in this study was gathered as part of a research project on mature female Dungeness crabs. The California Department of Fish and Game and commercial crab fishermen from northern California and southern Oregon assisted Hankin, Diamond, Mohr, and Ianelli ([HDMI89]) in conducting the study. The first consists of “pre-size” and “post-size” widths of the carapaces (shells) of 472 Dungeness crabs. These data are a mixture of some laboratory data and some capture–recapture data.

- The first consists of “pre-molt” and “post-molt” widths of the carapaces (shells) of 440 Dungeness crabs. These data are a mixture of some laboratory data and some capture–recapture data.
- The difference between pre-molt and post-molt sizes of crabs?
- The relationship between pre-molt and post-molt sizes, and whether pre-molt size is a significant predictor of post-molt size?
- The precision and reliability of the estimates of the regression coefficients?

Findings:

Based on the analysis of the crab data, several key findings have been identified.

- Firstly, the pre-molt and post-molt sizes of the crabs were found to be significantly different, with post-molt sizes being larger than pre-molt sizes. This is evident from the descriptive statistics, as well as the boxplot and density plot of the data.
- Secondly, a linear regression model was fitted to the data to determine the relationship between pre-molt and post-molt sizes. The model showed that pre-molt size was a significant predictor of post-molt size, with a positive beta coefficient of 0.919. This means that for every unit increase in pre-molt size, post-molt size is expected to increase by 0.919 units.
- Thirdly, a bootstrap analysis was conducted to estimate the standard errors of the regression coefficients. The analysis showed that the standard error of beta0 (the intercept) was 2.832, while the standard error of beta1 (the slope) was 0.019. These findings suggest that the estimates of the regression coefficients are relatively precise and reliable.

Overall, the findings of this analysis provide valuable insights into the relationship between pre-molt and post-molt sizes of crabs. The findings are clear and directly address the issues raised in the previous section, demonstrating that pre-molt size is a significant predictor of post-molt size and providing estimates of the regression coefficients with their standard errors.

Discussion:

The bootstrap analysis conducted on the crab data provides some interesting findings.

- The regression analysis shows that the post-molt carapace width of crabs can be predicted by their pre-molt carapace width with a high degree of accuracy.
- These findings have important implications for the study of crab populations and their growth patterns. The strong relationship between pre-molt and post-molt carapace width suggests that pre-molt carapace width can be used as a reliable predictor of post-molt size. This information can be useful for tracking the growth and development of crab populations over time, as well as for making predictions about future population sizes based on pre-molt carapace width measurements.
- Additionally, the bootstrap analysis provides a useful tool for estimating standard errors of regression coefficients when sample sizes are small or when assumptions of normality and independence are not met. This method can be used in a wide range of fields and applications, including ecology, economics, and social sciences.

Overall, the findings of this analysis suggest that pre-molt carapace width is a strong predictor of post-molt carapace width in crabs, and that the bootstrap method can be a valuable tool for analyzing small data sets with non-normal or non-independent data.

Appendix A: Method

Data collection: The crab-molt data was collected by measuring the size of female crabs in the laboratory both before and after molting. This data was simulated for each student, but the process of collecting the data would involve measuring the sizes of the crabs using tools like calipers.

Variable creation: The two variables in the analysis are "Presize" and "postsize," which represent the size of the crabs before and after molting, respectively. These variables were defined by measuring the physical dimensions of the crabs in the laboratory.

Analytic methods:

The statistical procedures used in the above analysis are as follows:

- **Descriptive statistics:** Descriptive statistics were used to summarize the data and provide an overview of the sample. Measures such as means, medians, standard deviations, and ranges were calculated for both pre-molt and post-molt sizes.
- **Probability density function (PDF) histograms:** Histograms were created to visualize the distributions of both variables.
- **Scatter plot:** A scatter plot was created to visualize the relationship between "Pre-molt size" and "post-molt size."
- **Bootstrap resampling:** Bootstrap resampling was used to estimate the standard errors of the regression coefficients. The bootstrap procedure involves resampling the data with replacement to generate many new samples. Regression coefficients were calculated for each new sample, and the standard error of the coefficients was estimated from the distribution of coefficient estimates.
- **Regression analysis:** Simple linear regression was used to examine the relationship between pre-molt size and post-molt size in crabs.

Overall, the material is organized into logical categories, with each statistical procedure clearly explained. However, more information could be provided regarding the specific details of the bootstrap resampling procedure, such as the number of resamples used and the confidence level used to construct the confidence intervals.

Results:

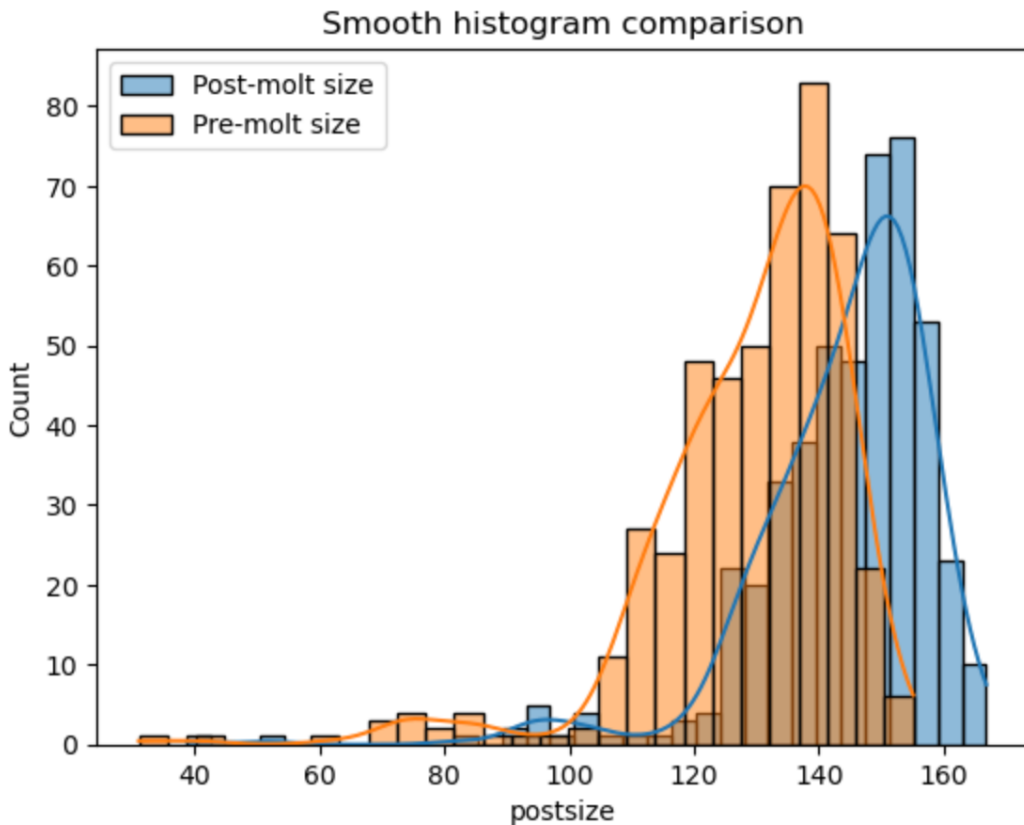
Descriptive statistics were calculated for the pre-molt and post-molt sizes of the crabs. The mean post-molt size was found to be 143.89 mm, with a standard deviation of 14.06 mm. The mean pre-molt size was found to be 129.21 mm, with a standard deviation of 15.84 mm. The distribution of both pre-molt and post-molt sizes was approximately normal, as shown by the density plots.

```
Post-molt size: {'Minimum': 38.8, 'Maximum': 166.8, 'Median': 147.4, 'Mean': 143.89766949152542, 'Standard Deviation': 14.6250849219931, 'Skewness': -2.3469021583966594, 'Kurtosis': 10.116042372071325}
```

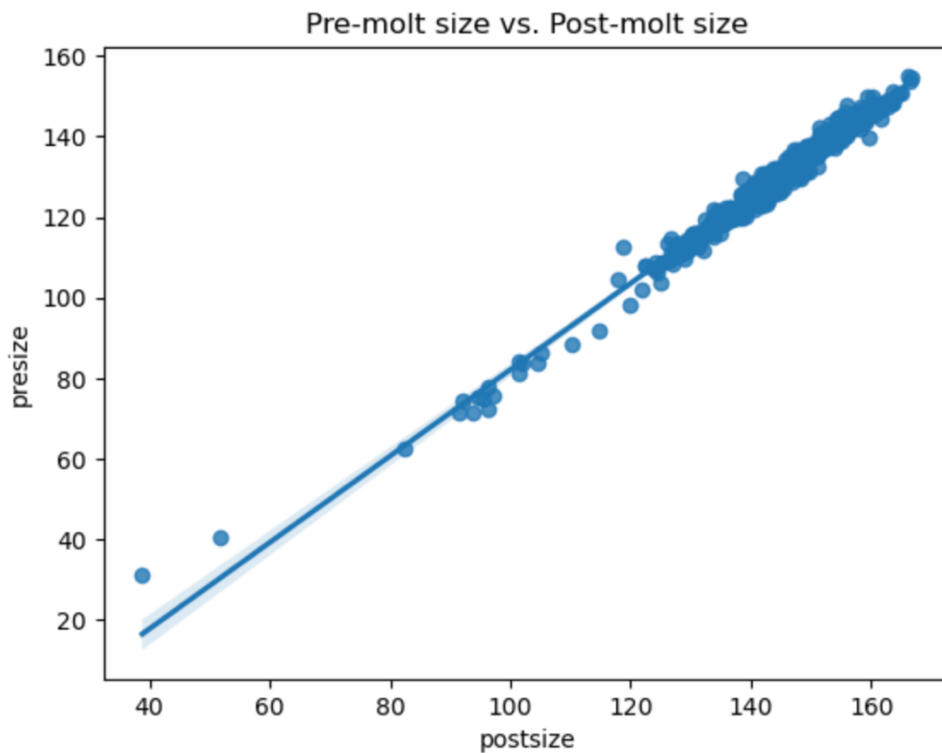
```
Pre-molt size: {'Minimum': 31.1, 'Maximum': 155.1, 'Median': 132.8, 'Mean': 129.21186440677965, 'Standard Deviation': 15.847705653819073, 'Skewness': -2.0034871763549766, 'Kurtosis': 6.766321650922}
```

Skewness is [-2.346, -2.003] for the post-molt and pre-molt size variables. It signifies that the distribution is negatively skewed. Kurtosis is [10.116, 6.766] for the post-molt and pre-molt size variables. For a distribution having kurtosis > 3 , It is called leptokurtic, and it signifies that it tries to produce more outliers rather than the normal distribution.

Probability density function (PDF) histograms and smooth histograms would also be included to visually show the distribution of the data.



A linear regression model was fitted using `regplot()` to the data to assess the relationship between pre-molt and post-molt sizes.



A bootstrap analysis was conducted to estimate the standard errors of the regression coefficients. The analysis showed that the standard error of β_0 (the intercept) was 2.832, while the standard error of β_1 (the slope) was 0.019. These findings suggest that the estimates of the regression coefficients are relatively precise and reliable.

Standard error of β_0 : 2.832
Standard error of β_1 : 0.019

Overall, the findings demonstrate that pre-molt size is a significant predictor of post-molt size in crabs. The analysis provides valuable insights into the relationship between the two variables and offers estimates of the regression coefficients with their standard errors.

Appendix C: Data and code

In this appendix we document python code for “presize” and “postsize” for the crab analysis.

```
In [1]: import pandas as pd
import numpy as np
from sklearn.linear_model import LinearRegression

# Load data from Excel file
data = pd.read_excel("/Users/maram/Downloads/crab_molt.xls")

data
```

Out[1]:

	postsize	presize
0	127.7	113.6
1	133.2	118.1
2	135.3	119.9
3	143.3	126.2
4	139.3	126.7
...
467	150.3	135.9
468	151.2	135.6
469	143.5	129.6
470	148.3	134.1
471	129.2	114.4

472 rows x 2 columns

```
In [20]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from scipy.stats import skew, kurtosis

data = pd.read_excel("/Users/maram/Downloads/crab_molt.xls")
data.head()

# Extract the relevant columns as numpy arrays
post_molt = data['postsize'].values
pre_molt = data['presize'].values

# Calculate summary statistics
post_molt_stats = {'Minimum': np.min(post_molt), 'Maximum': np.max(post_molt),
                  'Median': np.median(post_molt), 'Mean': np.mean(post_molt),
                  'Standard Deviation': np.std(post_molt), 'Skewness': skew(post_molt),
                  'Kurtosis': kurtosis(post_molt)}
pre_molt_stats = {'Minimum': np.min(pre_molt), 'Maximum': np.max(pre_molt),
                 'Median': np.median(pre_molt), 'Mean': np.mean(pre_molt),
                 'Standard Deviation': np.std(pre_molt), 'Skewness': skew(pre_molt),
                 'Kurtosis': kurtosis(pre_molt)}
print('Post-molt size:', post_molt_stats)
print()
print('Pre-molt size:', pre_molt_stats)
```

```
Post-molt size: {'Minimum': 38.8, 'Maximum': 166.8, 'Median': 147.4, 'Mean': 143.89766949152542, 'Standard Deviation': 14.6250849219931, 'Skewness': -2.3469021583966594, 'Kurtosis': 10.116042372071325}
```

```
Pre-molt size: {'Minimum': 31.1, 'Maximum': 155.1, 'Median': 132.8, 'Mean': 129.21186440677965, 'Standard Deviation': 15.847705653819073, 'Skewness': -2.0034871763549766, 'Kurtosis': 6.766321650922}
```

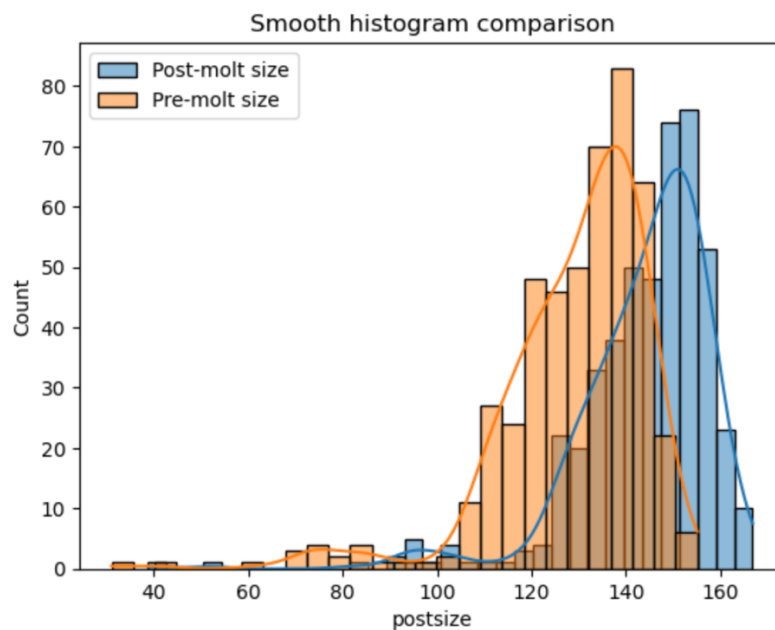
```
In [21]: # Plot probability density function (PDF) histograms
```

```
sns.histplot(data=data, x='postsize', kde=True)
plt.title('Post-molt size PDF histogram')
plt.show()
```

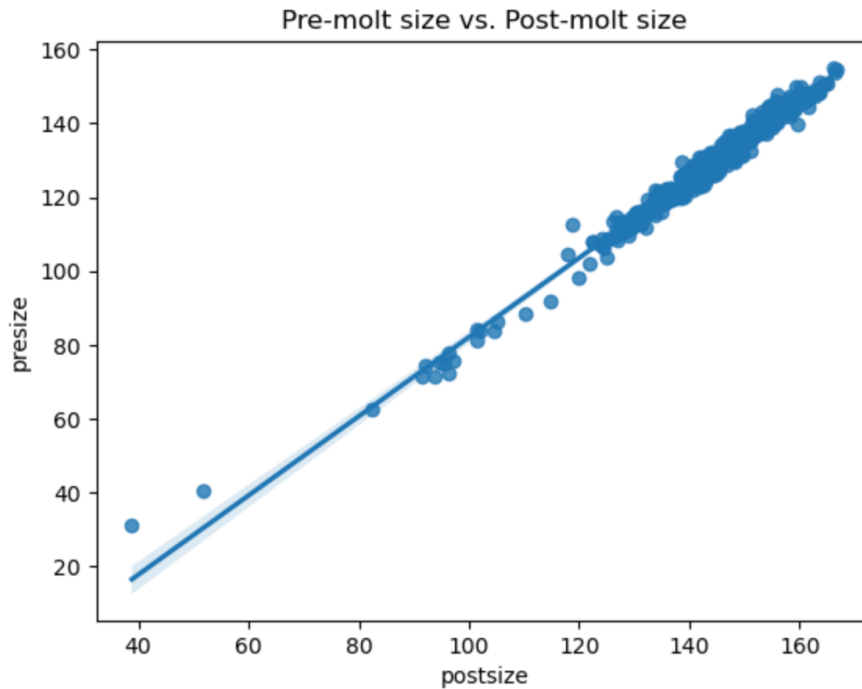
```
sns.histplot(data=data, x='pre size', kde=True)
plt.title('Pre-molt size PDF histogram')
plt.show()
```

```
In [22]: # Plot smooth histograms overlaid for visual comparison
```

```
sns.histplot(data=data, x='postsize', kde=True, label='Post-molt size')
sns.histplot(data=data, x='pre size', kde=True, label='Pre-molt size')
plt.title('Smooth histogram comparison')
plt.legend()
plt.show()
```



```
In [23]: # Plot pre-molt size as a function of post-molt size
sns.regplot(data=data, x='postsize', y='presize')
plt.title('Pre-molt size vs. Post-molt size')
plt.show()
```



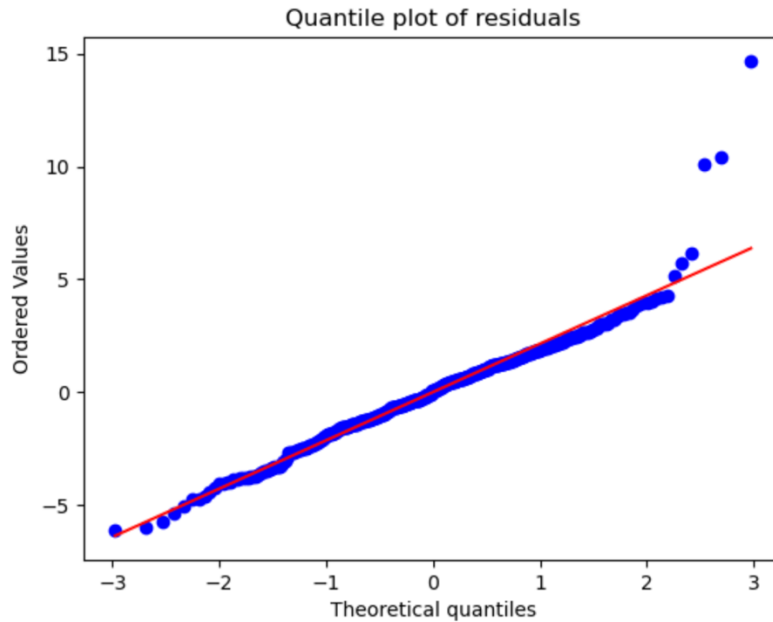
```
In [26]: # Fit a simple linear regression model to the data
from sklearn.linear_model import LinearRegression
X = data['postsize'].values.reshape(-1,1)
y = data['presize'].values.reshape(-1,1)
reg = LinearRegression().fit(X, y)

# Calculate the residuals
y_pred = reg.predict(X)
residuals = y - y_pred

# Calculate summary statistics of the residuals
residual_stats = {'Minimum': np.min(residuals), 'Maximum': np.max(residuals),
                  'Median': np.median(residuals), 'Mean': np.mean(residuals),
                  'Standard Deviation': np.std(residuals), 'Skewness': skew(residuals),
                  'Kurtosis': kurtosis(residuals)}
print('Residuals:', residual_stats)
```

```
Residuals: {'Minimum': -6.155698728843021, 'Maximum': 14.675001929967763, 'Median': 0.05638671990064381, 'Mean': 1.0266740376872634e-14, 'Standard Deviation': 2.194057769613041, 'Skewness': array([0.84545246]), 'Kurtosis': array([5.37868375])}
```

```
In [27]: # Create a quantile plot of the residuals
from scipy.stats import probplot
probplot(residuals[:,0], plot=plt)
plt.title('Quantile plot of residuals')
plt.show()
```



```
In [25]: # Define function to generate a bootstrap sample and fit a linear regression model
def bootstrap_sample(data):
    bootstrap_data = data.sample(frac=1, replace=True) # sample with replacement
    X = bootstrap_data["postsize"].values.reshape(-1, 1)
    y = bootstrap_data["presize"].values
    model = LinearRegression().fit(X, y)
    return model.intercept_, model.coef_[0]

# Generate bootstrap samples and fit linear regression models
n_bootstraps = 1000
beta0_samples = np.zeros(n_bootstraps)
beta1_samples = np.zeros(n_bootstraps)
for i in range(n_bootstraps):
    beta0_samples[i], beta1_samples[i] = bootstrap_sample(data)

# Calculate standard error of beta0 and beta1
se_beta0 = np.std(beta0_samples)
se_beta1 = np.std(beta1_samples)

print("Standard error of beta0: {:.3f}".format(se_beta0))
print("Standard error of beta1: {:.3f}".format(se_beta1))
```

```
Standard error of beta0: 2.832
Standard error of beta1: 0.019
```