

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

```
In [2]: 1 # Load the penguins dataset
        2 dataset = pd.read_csv("penguins.csv") # Make sure to provide the correct
        3
        4 # Display basic statistics about the dataset
        5 print(dataset.describe())
```

	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g \
count	342.000000	342.000000	342.000000	342.000000
mean	43.921930	17.151170	200.915205	4201.754386
std	5.459584	1.974793	14.061714	801.954536
min	32.100000	13.100000	172.000000	2700.000000
25%	39.225000	15.600000	190.000000	3550.000000
50%	44.450000	17.300000	197.000000	4050.000000
75%	48.500000	18.700000	213.000000	4750.000000
max	59.600000	21.500000	231.000000	6300.000000

	year
count	344.000000
mean	2008.029070
std	0.818356
min	2007.000000
25%	2007.000000
50%	2008.000000
75%	2009.000000
max	2009.000000

```
In [3]: 1 df = pd.DataFrame(dataset)
        2 df.head()
```

Out[3]:

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
0	Adelie	Torgersen	39.1	18.7	181.0	3750.0	male
1	Adelie	Torgersen	39.5	17.4	186.0	3800.0	female
2	Adelie	Torgersen	40.3	18.0	195.0	3250.0	female
3	Adelie	Torgersen	NaN	NaN	NaN	NaN	NaN
4	Adelie	Torgersen	36.7	19.3	193.0	3450.0	female

In [4]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 344 entries, 0 to 343
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  -
0   species               344 non-null   object
1   island                344 non-null   object
2   bill_length_mm        342 non-null   float64
3   bill_depth_mm         342 non-null   float64
4   flipper_length_mm     342 non-null   float64
5   body_mass_g           342 non-null   float64
6   sex                   333 non-null   object
7   year                  344 non-null   int64
dtypes: float64(4), int64(1), object(3)
memory usage: 21.6+ KB
```

In [5]: 1 df.dtypes

```
Out[5]: species                object
island                  object
bill_length_mm          float64
bill_depth_mm           float64
flipper_length_mm       float64
body_mass_g             float64
sex                     object
year                    int64
dtype: object
```

In [8]:

```
1  # Calculate the mean values of each column
2  column_means = df.mean()
3
4  # Fill missing values in the specified columns with their respective means
5  columns_to_impute = ['bill_length_mm', 'bill_depth_mm', 'flipper_length_mm']
6
7  df[columns_to_impute] = df[columns_to_impute].fillna(column_means[columns_to_impute])
```

In [9]:

```
1  import pandas as pd
2
3  # List of columns to convert to numeric
4  columns_to_convert = ['bill_length_mm', 'bill_depth_mm', 'flipper_length_mm']
5
6  # Loop through the list of columns and convert to numeric
7  for column in columns_to_convert:
8      df[column] = pd.to_numeric(df[column], errors='coerce')
```

In [10]: 1 df.head()

Out[10]:

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
0	Adelie	Torgersen	39.10000	18.70000	181.000000	3750.000000	male
1	Adelie	Torgersen	39.50000	17.40000	186.000000	3800.000000	female
2	Adelie	Torgersen	40.30000	18.00000	195.000000	3250.000000	female
3	Adelie	Torgersen	43.92193	17.15117	200.915205	4201.754386	NaN
4	Adelie	Torgersen	36.70000	19.30000	193.000000	3450.000000	female

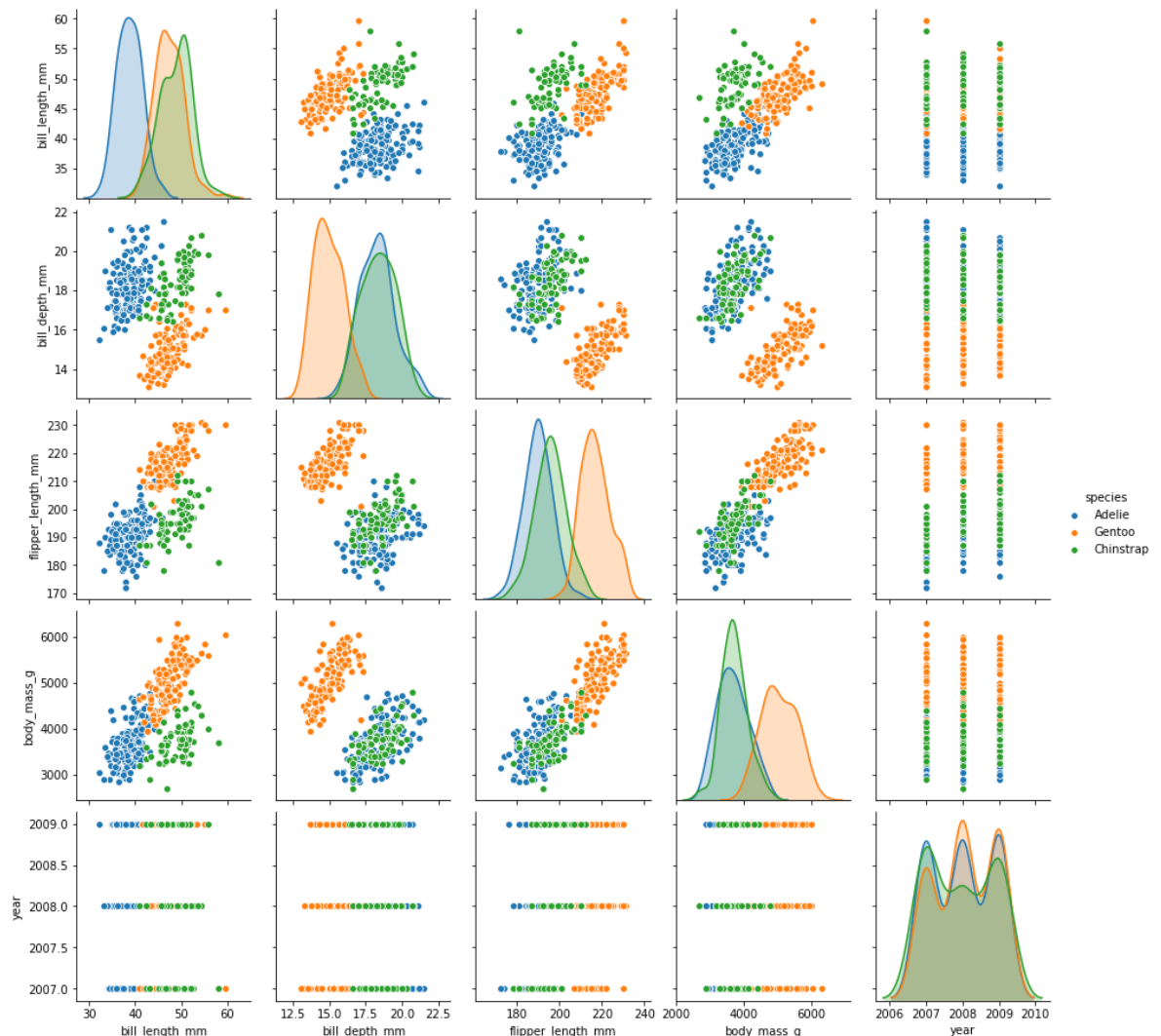
```

In [12]: 1 # Define a function to remove outliers using standard deviation
2 def remove_outliers_std(data_frame, column_name, std_threshold=3):
3     column_mean = data_frame[column_name].mean()
4     column_std = data_frame[column_name].std()
5     lower_bound = column_mean - std_threshold * column_std
6     upper_bound = column_mean + std_threshold * column_std
7
8     data_frame = data_frame[(data_frame[column_name] >= lower_bound) & (data_frame[column_name] <= upper_bound)]
9     return data_frame
10
11 # List of columns to remove outliers from
12 columns_to_remove_outliers = ['bill_length_mm', 'bill_depth_mm', 'flipper_length_mm']
13
14 # Iterate through the list of columns and remove outliers
15 for column in columns_to_remove_outliers:
16     df = remove_outliers_std(df, column)

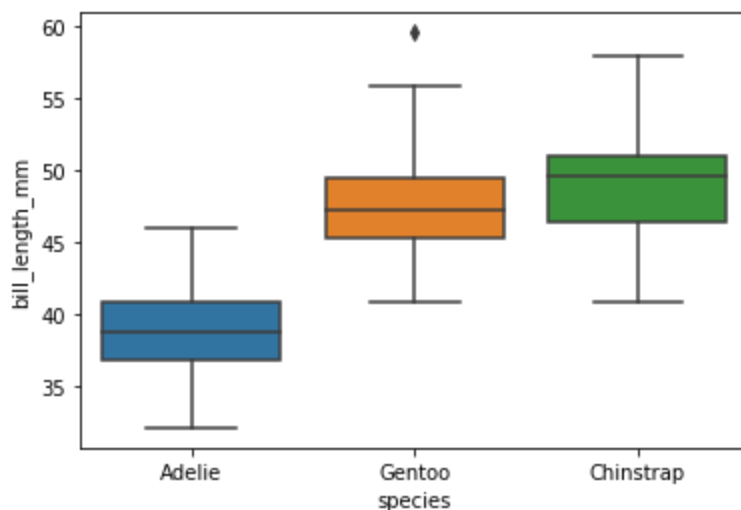
```

Visualization

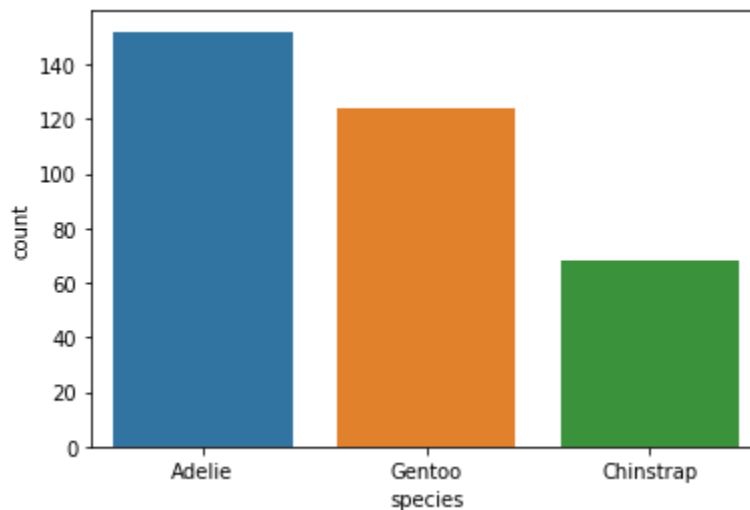
```
In [14]: 1 #A pairplot is a great way to visualize the relationships between numeric
2 sns.pairplot(df, hue="species")
3 plt.show()
```



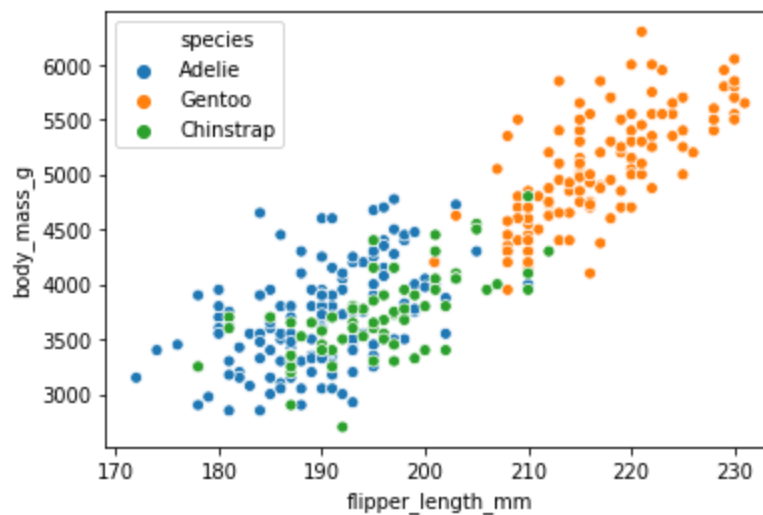
```
In [15]: 1 #Box plots can help you visualize the distribution and spread of numeric j
2 sns.boxplot(x="species", y="bill_length_mm", data=df)
3 plt.show()
```



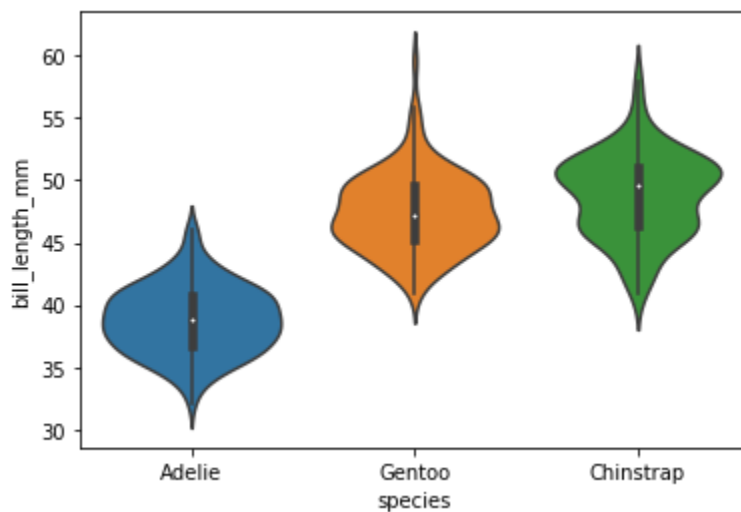
```
In [16]: 1 #Count plots can be used to visualize the distribution of categorical feat  
2 sns.countplot(x="species", data=df)  
3 plt.show()
```



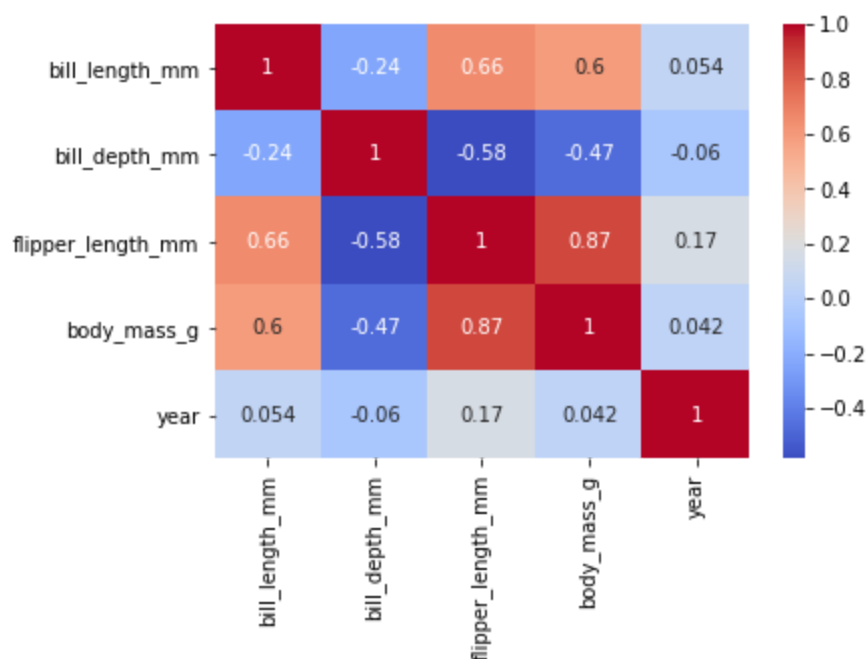
```
In [17]: 1 #Scatter plots can show relationships between two numeric variables.  
2 sns.scatterplot(x="flipper_length_mm", y="body_mass_g", data=df, hue="species")  
3 plt.show()
```



```
In [18]: 1 #Violin plots combine box plots with kernel density estimates.
2
3 sns.violinplot(x="species", y="bill_length_mm", data=df)
4 plt.show()
```



```
In [19]: 1 #A heatmap can show the correlation between numeric features.
2 correlation_matrix = df.corr()
3 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
4 plt.show()
```



```
In [20]: 1 # Remove the 'year' column from the DataFrame permanently as it is not used
2 df.drop(columns=['year'], inplace=True)
```

```

In [21]: 1 # Define a function to label encode and replace missing values with the mode
2 def label_encode_with_mode(data_frame, column_name):
3     # Calculate the mode of the column
4     mode = data_frame[column_name].mode().iloc[0] # Get the first mode (1
5
6     # Fill missing values with the mode
7     data_frame[column_name].fillna(mode, inplace=True)
8
9     # Use label encoding to convert the column to numeric
10    data_frame[column_name] = data_frame[column_name].astype('category').c
11
12    # List of categorical columns to process
13    categorical_columns = ['species', 'island', 'sex']
14
15    # Iterate through the list of columns and label encode with mode replacement
16    for column in categorical_columns:
17        label_encode_with_mode(df, column)

```

```

In [22]: 1 df.head()

```

Out[22]:

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
0	0	2	39.10000	18.70000	181.000000	3750.000000	1
1	0	2	39.50000	17.40000	186.000000	3800.000000	0
2	0	2	40.30000	18.00000	195.000000	3250.000000	0
3	0	2	43.92193	17.15117	200.915205	4201.754386	1
4	0	2	36.70000	19.30000	193.000000	3450.000000	0

```

In [32]: 1 # Define the columns to be normalized
2 columns_to_normalize = ['flipper_length_mm', 'body_mass_g']
3
4 # Calculate min and max values for each column
5 min_values = df[columns_to_normalize].min()
6 max_values = df[columns_to_normalize].max()
7
8 # Normalize the dataset columns to the range [0, 1]
9 for column in columns_to_normalize:
10     df[column] = (df[column] - min_values[column]) / (max_values[column] - min_values[column])

```

```

In [33]: 1 df.head()

```

Out[33]:

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
0	0	2	39.10000	18.70000	0.152542	0.291667	1
1	0	2	39.50000	17.40000	0.237288	0.305556	0
2	0	2	40.30000	18.00000	0.389831	0.152778	0
3	0	2	43.92193	17.15117	0.490088	0.417154	1
4	0	2	36.70000	19.30000	0.355932	0.208333	0

```
In [34]: 1 # Define your target 'y' and feature 'X' variables
2 X = df.drop('sex', axis=1) # Features: ALL columns except 'sex'
3 y = df['sex'] # Target variable: 'sex' column
```

```
In [35]: 1 # Define the split ratio
2 split_ratio = 0.8 # 80% training, 20% testing
3
4 # Calculate the number of samples for training and testing
5 total_samples = len(X)
6 num_train_samples = int(np.round(total_samples * split_ratio))
7 num_test_samples = total_samples - num_train_samples
8
9 # Slice the data for training and testing
10 X_train = X.iloc[:num_train_samples]
11 X_test = X.iloc[-num_test_samples:]
12 y_train = y.iloc[:num_train_samples]
13 y_test = y.iloc[-num_test_samples:]
14
15 # Print the shapes of the training and testing sets
16 print("X_train Shape: ", X_train.shape)
17 print("y_train Shape: ", y_train.shape)
18 print()
19 print("X_test Shape: ", X_test.shape)
20 print("y_test Shape: ", y_test.shape)
```

X_train Shape: (275, 6)

y_train Shape: (275,)

X_test Shape: (69, 6)

y_test Shape: (69,)

In [36]:

```

1  class LogitRegression:
2
3      def __init__(self, learning_rate, iterations):
4          self.learning_rate = learning_rate
5          self.iterations = iterations
6          self.weights = None
7          self.bias = None
8
9      def sigmoid(self, z):
10         return 1 / (1 + np.exp(-z))
11
12     def cost(self, y_train, X_train):
13         z = np.dot(X_train, self.weights) + self.bias
14         h = self.sigmoid(z)
15
16         parameter1 = -(y_train) * np.log(h)
17         parameter2 = (1 - y_train) * np.log(1 - h)
18
19         j = (1 / len(y_train)) * np.sum(parameter1 - parameter2)
20
21         return j
22
23     def gradient_descent(self, y_train, X_train):
24         z = np.dot(X_train, self.weights) + self.bias
25         pred = self.sigmoid(z)
26
27         difference_y = pred - y_train
28
29         update_weight = np.dot(X_train.T, difference_y) / len(y_train)
30         update_bias = np.sum(difference_y) / len(y_train)
31
32         return update_weight, update_bias
33
34     def scaling(self, X):
35         mean = np.mean(X, axis=0)
36         std = np.std(X, axis=0)
37         X_scaled = (X - mean) / std
38         return X_scaled
39
40     def fit(self, X_train, y_train):
41         X_scaled = self.scaling(X_train)
42         rows, features = X_scaled.shape
43         self.weights = np.random.uniform(0, 1, features)
44
45         self.bias = 0.5
46         loss = []
47
48         for i in range(self.iterations):
49             updated_weights, updated_bias = self.gradient_descent(y_train,
50                 loss.append(self.cost(y_train, X_scaled))
51
52             self.weights = self.weights - (self.learning_rate * updated_w
53             self.bias = self.bias - (self.learning_rate * updated_bias)
54
55         return loss
56
57     def predict(self, X_test):
58         X_scaled = self.scaling(X_test)
59
60         z = np.dot(X_scaled, self.weights) + self.bias
61         y_hat = self.sigmoid(z)

```

```
62  
63     y_pred = [1 if y_value >= 0.5 else 0 for y_value in y_hat]  
64  
65     return y_pred
```

In [37]:

```
1 def accuracy(y_test, y_pred):  
2     matched_values = np.sum(y_test == y_pred)  
3     return matched_values / len(y_test)
```

In [39]:

```

1  # Define learning rates and iterations to explore
2  learning_rates = [0.001, 0.01, 0.1]
3  iteration_values = [10000, 50000, 100000]
4
5  # Create dictionaries to store loss and accuracy for each combination
6  loss_dict = {}
7  accuracy_dict = {}
8
9  # Iterate through learning rates and iterations
10 for lr in learning_rates:
11     for iterations in iteration_values:
12         # Create and train the model
13         logistic_model = LogisticRegression(lr, iterations)
14         loss_values = logistic_model.fit(X_train, y_train)
15
16         # Make predictions
17         y_pred = logistic_model.predict(X_test)
18
19         # Calculate accuracy
20         accu = accuracy(y_test, y_pred)
21
22         # Store the results
23         key = (lr, iterations)
24         loss_dict[key] = loss_values
25         accuracy_dict[key] = accu
26
27         # Print the weight vector for each iteration
28         print(f"Learning Rate: {lr}, Iterations: {iterations}")
29         for i, weight in enumerate(logistic_model.weights):
30             print(f"Iteration {i}, Weight {i}: {weight:.6f}")
31
32 # Plot gradient progress for each combination of learning rate and iterations
33 plt.figure(figsize=(12, 6))
34 for lr, iterations in loss_dict.keys():
35     plt.plot(loss_dict[(lr, iterations)], label="LR: {}, It: {}".format(lr, iterations))
36 plt.xlabel("Iteration")
37 plt.ylabel("Loss")
38 plt.title("Gradient Descent Progress for Different Learning Rates and Iterations")
39 plt.legend()
40 plt.show()
41
42 # Print accuracy for each combination
43 for (lr, iterations), accu in accuracy_dict.items():
44     print(f"Learning Rate: {lr}, Iterations: {iterations}, Accuracy: {accu}")

```

Learning Rate: 0.001, Iterations: 10000

Iteration 0, Weight 0: 0.019368

Iteration 1, Weight 1: 0.215603

Iteration 2, Weight 2: 0.473955

Iteration 3, Weight 3: 1.386185

Iteration 4, Weight 4: 0.534324

Iteration 5, Weight 5: 0.707313

Learning Rate: 0.001, Iterations: 50000

Iteration 0, Weight 0: -0.398116

Iteration 1, Weight 1: 0.199784

Iteration 2, Weight 2: 1.337713

Iteration 3, Weight 3: 2.476437

Iteration 4, Weight 4: 0.057959

Iteration 5, Weight 5: 1.939160

Learning Rate: 0.001, Iterations: 100000

Iteration 0, Weight 0: -0.947539

Iteration 1, Weight 1: 0.136512

Iteration 2, Weight 2: 1.640533

Iteration 3, Weight 3: 2.703672

Iteration 4, Weight 4: 0.154033

Iteration 5, Weight 5: 2.362995

Learning Rate: 0.01, Iterations: 10000

Iteration 0, Weight 0: -0.835204

Iteration 1, Weight 1: 0.136277

Iteration 2, Weight 2: 1.616419

Iteration 3, Weight 3: 2.784613

Iteration 4, Weight 4: 0.225171

Iteration 5, Weight 5: 2.257634

Learning Rate: 0.01, Iterations: 50000

Iteration 0, Weight 0: -1.579083

Iteration 1, Weight 1: 0.154800

Iteration 2, Weight 2: 1.967877

Iteration 3, Weight 3: 2.745980

Iteration 4, Weight 4: -0.228035

Iteration 5, Weight 5: 3.282394

Learning Rate: 0.01, Iterations: 100000

Iteration 0, Weight 0: -1.798479

Iteration 1, Weight 1: 0.145665

Iteration 2, Weight 2: 1.993314

Iteration 3, Weight 3: 2.659150

Iteration 4, Weight 4: -0.207990

Iteration 5, Weight 5: 3.399491

Learning Rate: 0.1, Iterations: 10000

Iteration 0, Weight 0: -1.810592

Iteration 1, Weight 1: 0.144939

Iteration 2, Weight 2: 1.994632

Iteration 3, Weight 3: 2.654287

Iteration 4, Weight 4: -0.204945

Iteration 5, Weight 5: 3.403569

Learning Rate: 0.1, Iterations: 50000

Iteration 0, Weight 0: -1.854058

Iteration 1, Weight 1: 0.142556

Iteration 2, Weight 2: 1.999390

Iteration 3, Weight 3: 2.637078

Iteration 4, Weight 4: -0.195432

Iteration 5, Weight 5: 3.420235

Learning Rate: 0.1, Iterations: 100000

Iteration 0, Weight 0: -1.854058

Iteration 1, Weight 1: 0.142556

Iteration 2, Weight 2: 1.999390

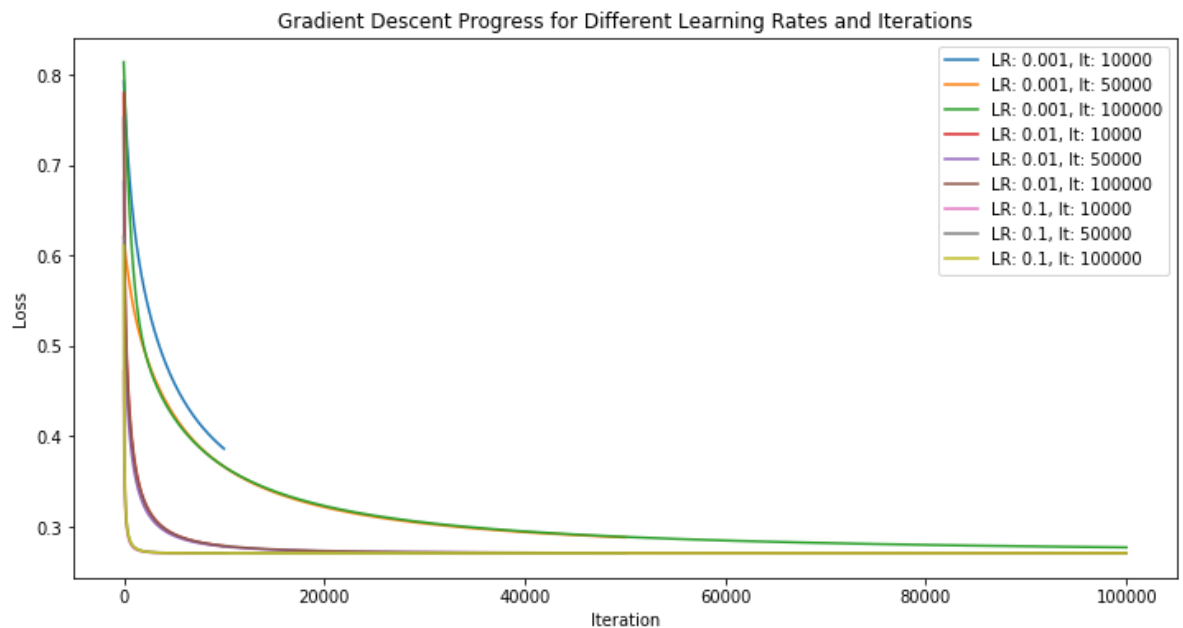
Iteration 3, Weight 3: 2.637078

Iteration 4, Weight 4: -0.195432

Iteration 5, Weight 5: 3.420235

D:\anaconda\lib\site-packages\IPython\core\pylabtools.py:132: UserWarning: Creating legend with loc="best" can be slow with large amounts of data.

```
fig.canvas.print_figure(bytes_io, **kw)
```



Learning Rate: 0.001, Iterations: 10000, Accuracy: 86.96%

Learning Rate: 0.001, Iterations: 50000, Accuracy: 88.41%

Learning Rate: 0.001, Iterations: 100000, Accuracy: 89.86%

Learning Rate: 0.01, Iterations: 10000, Accuracy: 88.41%

Learning Rate: 0.01, Iterations: 50000, Accuracy: 88.41%

Learning Rate: 0.01, Iterations: 100000, Accuracy: 88.41%

Learning Rate: 0.1, Iterations: 10000, Accuracy: 88.41%

Learning Rate: 0.1, Iterations: 50000, Accuracy: 88.41%

Learning Rate: 0.1, Iterations: 100000, Accuracy: 88.41%

In []:

1