

Instructor: Alina Vereshchaka

# Assignment 1

## Classification and Regression Methods

Checkpoint: Sep 28, Thu, 11:59 pm

Due Date: Oct 12, Thu, 11:59 pm

Our Assignment 1 is focused on learning how to build classification and regression models from scratch. In the first part you will perform data preprocessing. In the second part, you will implement logistic regression and test it on a penguin dataset. In the third part, you will implement linear regression using ordinary least squares (OLS) and in the last part, you will extend it to ridge regression. There are also bonus tasks available that you can consider to complete.

The goal of this assignment is to give you hands-on experience with building classification and regression models and apply them to solve the tasks based on the real-world datasets.

**Note:** For this assignment, *scikit-learn* or any other libraries with in-built functions that help to implement ML methods cannot be used. Submissions with used ML libraries (e.g. *scikit-learn*) will not be evaluated.

### Part I: Data Analysis & Preprocessing [10 points]

In this section, we will focus on data analysis and preprocessing.

#### DATASETS:

Choose TWO datasets from the 'noisy\_datasets' folder provided at UBLearn:

1. Penguins Dataset from the 'noisy\_datasets' folder. This dataset is slightly different from the one provided in A0.
2. One more dataset from the 'noisy\_datasets' folder

#### TASK:

1. Import required libraries (not allowed: *scikit-learn* or any other libraries with in-built functions that help to implement ML methods)
2. Read, preprocess, and print the main statistics about the dataset (you can reuse your code from A0 with a proper citation)
3. Handle missing entries. Possible solutions:
  - Drop rows with missing entries. If you have a large dataset and only a few missing

features, it may be acceptable to drop the rows containing missing values.

- Impute missing data. Replace the missing entries with the mean/median/mode of the feature. You can use K-Nearest Neighbor algorithm to find the matching sample.

4. Handle mismatched string formats.

For example, in the penguins dataset, the "Species" feature might appear as "Adelie" or "adelie," both of which refer to the same penguin species. These variations should be standardized to a consistent format such as "Adelie" or "adelie" to ensure consistency.

5. Handle outliers. Detect and manage outliers within the dataset.

For example, in the penguins dataset, while flipper lengths typically fall within the range of [180 – 210], certain entries might exhibit values like [10, 30]. These can be considered outliers. Possible solutions:

- Remove outliers. If there are just a few outliers is limited, you may eliminate the rows containing these outliers.
- Impute outliers. Replace the outliers with the mean/median/mode of the feature.

6. Using any data visualization library (e.g. [matplotlib](#), [seaborn](#), [plotly](#)), provide at least 5 visualization graphs related to your dataset. You can utilize any columns or a combination of columns in your dataset to generate graphs. E.g. correlation matrix, features vs. the target, counts of categorical features vs. the target.

7. Identify uncorrelated or unrelated features.

Unrelated or uncorrelated features can introduce confusion to your model and negatively impact its performance. You can compute the correlation matrix between the features and the target variable. Features with a low correlation coefficient should be identified and subsequently dropped from the dataset to enhance model performance.

8. Convert features with string datatype to categorical (e.g., species, island, gender). Possible ways:

- One-hot encoding, creating binary columns for each category, denoting their presence or absence. E.g., in the "Species" feature, "Adelie," "Chinstrap," and "Gentoo" become binary columns with "1" for presence and "0" for absence.
- Label encoding assigns unique integers to distinct feature values, useful for ordinal relationships among categories. E.g., "Small" as 0, "Medium" as 1, and "Large" as 2 can represent a "Size" feature. However, it may introduce unintended patterns.

9. Normalize non-categorical features (e.g. bill\_length\_mm, bill\_depth\_mm, flipper\_length\_mm, body\_mass\_g).

- a. Find the min and max values for each column.
- b. Rescale dataset columns to the range from 0 to 1

Why do we do this? Normalization is to transform features to be on a similar scale. This improves the performance and training stability of the model.

**Note:** `normalize()` is not allowed as it is a part of *scikit-learn* library.

10. Select another dataset and complete Steps 1-9. In total, you need to provide the data analysis and preprocessing for TWO datasets: 'penguins' and any of your choice from the provided folder 'noisy\_datasets'.

### In your report for Part I:

For each dataset:

1. Provide short description about your selected dataset (e.g. number of samples, the domain, main statistics)
2. Provide the short details of the methods you used for data preprocessing
3. Provide graphs and your short description

## Part II: Logistic Regression [40 points]

In this part, we will work on logistic regression and will use a logistic function to model a binomial (Binary / Bernoulli) output variable. The logistic regression model predicts that the observation belongs to a particular category. To generate these probabilities, logistic regression uses the sigmoid function that maps a real number to a value between 0 and 1.

### DATASET

For this Part, we will use 'penguins' dataset, that you have preprocessed in Part I.

### STEPS

1. Import required libraries (not allowed: *scikit-learn* or any other libraries with in-built functions that help to implement ML methods).
2. Choose your target Y. For this dataset, there are several options:
  - We can use a binary classifier to predict which gender a penguin belongs to (female or male). In this case, column gender can be used as Y (target)
  - We can use a binary classifier to predict if a penguin's location is Torgersen island or not. In this case, column island can be used as Y (target)
3. Create the data matrices for X (input) and Y (target) in a shape  $X = N \times d$  and  $Y = N \times 1$ , where  $N$  is a number of data samples and  $d$  has a number of features.
4. Divide the dataset into training and test, as 80% training and 20% testing dataset.
5. Print the shape of your X\_train, y\_train, X\_test, y\_test
6. Recommended structure of your code to define logistic regression:

```
class LogitRegression()  
  
    def __init__()
```

# Takes as an input hyperparameters: learning rate and the number of iterations.

```
def sigmoid():
```

# Define a sigmoid function as

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

```
def cost():
```

# Loss function for Logistic Regression can be defined as

$$h = \sigma(w^T x + b) = \left( \frac{1}{1 + e^{-(w^T x + b)}} \right)$$

$$J(w) = \frac{1}{N} (-y * \log(h) - (1 - y) * \log(1 - h))$$

```
def gradient_descent():
```

# Define current prediction  $\hat{y}$  for logistic regression as

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{X} + b)$$

# Gradient descent is just the derivative of the loss function with respect to its weights. Thus

$$\frac{\partial J(w)}{\partial w} = \frac{1}{N} X^T (\sigma(w^T x + b) - y)$$

To implement this formula, you can use intermediate variables, e.g.

```
pred =  $\sigma(w^T x + b)$ 
delta = pred - y_train
dW = (X^T * delta) / N
```

Update rule:  $w = w - \alpha \nabla J(w)$ , where  $\alpha$  is a learning rate

```
def fit():
```

# This method performs the training.

# Initialize weights

# For a number of iterations

# Call gradient\_descent function

# Call cost function and keep it in an array , e.g.

```
loss.append()
```

```
def predict(self, X):
```

# Return the predicted result in the binary form

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{X} + b) = \begin{cases} +1 & \text{if } \sigma(z) \geq 0.5 \\ 0 & \text{if } \sigma(z) < 0.5 \end{cases}$$

7. Train the model:

- Define a model by calling LogitRegression class and passing hyperparameters, e.g.  
`model = LogitRegression(learning_rate, iterations)`
- Train the model, by calling fit function and passing your training dataset, e.g.  
`model.fit(X_train, y_train)`
- Try at least THREE various hyperparameters. You can try different learning rates and number of iterations to improve your accuracy (accuracy of greater than 64% is expected). Suggested hyperparameters:

`learning_rate=1e-3`

`iterations=100000`

`weights = np.random.uniform(0, 1)`

8. Save the weights of the model, that returns the highest accuracy as pickle file. You will need to submit it along with your other files. [Check more details about Pickle](#)
9. Make a prediction on test dataset by counting how many correct/incorrect predictions your model makes and print your accuracy.

**Accepted accuracy rate for penguin dataset is above 64%**

10. Plot the loss graph and print out the loss values over each iteration.

**In your report for Part II:**

1. Provide your best accuracy.
2. Include a loss graph and provide a short analysis of the results.
3. Provide at least 3 different setups with learning rate and #iterations and discuss the results along with plotting of graphs. Plot the graph to discuss impact and loss over the iterations. Explain how hyperparameters influence the accuracy of the model.
4. Discuss the benefits/drawbacks of using a Logistic Regression model.

**Checkpoint Submission (Part I and Part II):**

- Report (as a pdf file) named as TEAMMATE1\_TEAMMATE2 \_assignment1\_report.pdf  
e.g., avereshc\_ pinazmoh\_ assignment1\_report.pdf
- Code (as separate ipynb file with saved outputs) named as  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part\_1.ipynb  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part\_2.ipynb  
e.g., avereshc\_ pinazmoh\_ assignment1\_ part\_1.ipynb
- Pickle files with saved weights that generate the best results for your model named as  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part2.pickle
- Preprocesses datasets, named as DATASET\_ preprocessed.csv  
e.g. penguins\_preprocessed.csv

## Part III: Linear Regression [25 points]

In this part, we implement a linear regression model and apply this model to solve a problem based on the real-world dataset.

Implement linear regression using the ordinary least squares (OLS) method to perform direct minimization of the squared loss function.

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

In matrix-vector notation, the loss function can be written as:

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

where  $\mathbf{X}$  is the input data matrix,  $\mathbf{y}$  is the target vector, and  $\mathbf{w}$  is the weight vector for regression.

### DATASET

For this Part, you can work with a dataset that you preprocessed in Part I or any other dataset from 'noisy\_dataset' folder except from 'penguin' dataset.

### STEPS

1. Import required libraries (not allowed: scikit-learn or any other libraries with in-built functions that help to implement ML).
2. Perform a data analysis and data preprocessing for your selected dataset. You can reuse your code from Part I. The dataset should come from 'noisy\_dataset' and it should not be 'penguins'
3. Choose your target  $Y$ .
4. Create the data matrices for  $X$  (input) and  $Y$  (target) in a shape  $X = N \times d$  and  $Y = N \times 1$ , where  $N$  is a number of data samples and  $d$  is a number of features.
5. Divide the dataset into training and test, as 80% training, 20% testing dataset.
6. Print the shape of your  $X_{\text{train}}$ ,  $y_{\text{train}}$ ,  $X_{\text{test}}$ ,  $y_{\text{test}}$ .
7. Calculate the weights with the OLS equation:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

8. Get the predictions and calculating the sum of squared errors:

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

9. Plot the predictions vs the actual targets.

### In your report for Part III:

1. Data analysis (you can reuse Part I):
  - a. Provide brief details about the nature of your dataset. What is it about? What type of data are we encountering? How many entries and variables does the dataset comprise?
  - b. Provide the main statistics about the entries of the dataset (mean, std, number of missing values, etc.)
  - c. Provide at least 5 visualization graphs with a brief description for each graph, e.g. discuss if there are any interesting patterns or correlations.
2. Provide your loss value.
3. Show the plot comparing the predictions vs the actual test data.
4. Discuss the benefits/drawbacks of using OLS estimate for computing the weights.
5. Discuss the benefits/drawbacks of using a Linear Regression model.

## Part IV: Ridge Regression [25 points]

**Use your implementation from Part III and extend it to Ridge Regression.**

Implement parameter estimation for ridge regression by minimizing the regularized squared loss as follows:

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \frac{1}{2} \lambda \mathbf{w}^T \mathbf{w}$$

In matrix-vector notation, the squared loss can be written as:

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \frac{1}{2} \lambda \mathbf{w}^T \mathbf{w}$$

OLS equation for Ridge regression can be estimated as

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

### DATASET

For this Part, you can work with the same dataset that you used for Part III.

### STEPS

Reuse steps from Part III while using OLS equation for Ridge regression to learn the parameters.

**In your report for Part III:**

1. Provide your loss value.
2. Show the plot comparing the predictions vs the actual test data.
3. Discuss the difference between Linear and Ridge regressions. What is the main motivation for using  $l_2$  regularization?
4. Discuss the benefits/drawbacks of using a Ridge Regression model.

**Final Submission (reupload your Part I & Part II, completed Part III & Part IV and report):**

- Report (as a pdf file): combine the reports for all parts into one pdf file named as TEAMMATE1\_TEAMMATE2 \_assignment1\_report.pdf (e.g., avereshc\_ pinazmoh\_ assignment1\_report.pdf)
- Code (as ipynb file with saved outputs): separate file for Part I, Part II, and a single file for both Part III & IV named as  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part\_1.ipynb  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part\_2.ipynb  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part\_3\_4.ipynb  
(e.g., avereshc\_ pinazmoh\_ assignment1\_part\_1.ipynb)
- Pickle files with saved weights that generate the best results for your model for Part II, Part III & Part IV named as  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part2.pickle  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part3.pickle  
TEAMMATE1\_TEAMMATE2 \_assignment1\_part4.pickle
- Preprocesses datasets, named as DATASET\_ preprocessed.csv  
e.g. penguins\_preprocessed.csv
- Add all the assignment-related files in a zip folder
- Notes:
  - Ensure that your code follows a clear structure and contains comments for the main functions and specific attributes related to your solution. You can submit multiple files, but they all need to be labeled with a clear name.
  - After executing command `python main.py` in the first level directory or Jupyter Notebook, it should generate all the results and plots you used in your report and print them out in a clear manner.



## Bonus points [max 10 points]

### Gradient Descent from Scratch [5 points]

- Based on your implementation of Ridge Regression in Part IV, implement a gradient descent method from scratch.
- Discuss the results.
- Provide training and testing errors and time to train.
- Try at least THREE different weight initializations. How do the different weight initialization techniques affect performance?

Try the following methods to stop gradient descent and compare their performance:

1. Stop after a fixed number of iterations: Run the gradient descent algorithm for a fixed number of iterations e.g., 10,000 or 100,000.
2. Stop based on gradient: Stop when the gradient value is equal to 0 or close to it. E.g.,  $-0.01 < \text{gradient} < 0.01$

### Elastic Net Regularization from Scratch [5 points]

Based on your implementation of Ridge regression in Part IV, implement an elastic net regularization and compare the results with Part IV.

### Bonus part submission:

- Create a separate Jupyter Notebook (.ipynb) named as TEAMMATE1\_TEAMMATE2\_assignment1\_bonus.ipynb  
e.g., avereshc\_pinazmoh\_assignment1\_bonus.ipynb
- You can duplicate code from your Part IV if needed
- Submit Jupyter Notebook (.ipynb) with saved outputs
- A file with saved weights that generate the best results for your model (.pickle).
- Report is not required, you can include all the analysis as part of your Jupyter Notebook.

## ASSIGNMENT STEPS

### 1. Register your team (September 18)

You may work individually or in a team of up to 2 people. The evaluation will be the same for a team of any size.

Register your team at UBLearns > Groups. In case you joined the wrong group, make a private post on piazza.

## 2. Submit checkpoint (September 28)

- Complete Part I and Part II of the assignment
- For the checkpoint report for Part I and Part II is not mandatory, you can complete the report and submit it along with the final submission
- Add all your assignment files in a zip folder including .ipynb files, the report (if available) and .pickle file at UBLearns > Assignments
- Name zip folder with all the files as  
TEAMMATE1\_TEAMMATE2 \_assignment1\_checkpoint.zip  
e.g., avereshc\_ pinazmoh\_ assignment1\_checkpoint.zip
- Submit to UBLearns > Assignments
- If you are working in a team, we expect equal contribution for the assignment. Each team member is expected to make a code-related contribution. Provide a contribution summary by each team member in the form of a table below. If the contribution is highly skewed, then the scores of the team members may be scaled w.r.t the contribution.

Team Member	Assignment Part	Contribution (%)

## 3. Submit final results (October 12)

- Fully complete all parts of the assignment
- Submit to UBLearns > Assignments
- Add all your assignment files in a zip folder including ipynb files for Part I, Part II, Part III, Part IV, & Bonus part (optional), the report and .pickle file at UBLearns > Assignments
- Name zip folder with all the files as  
TEAMMATE1\_TEAMMATE2 \_assignment1\_final.zip  
e.g. avereshc\_ pinazmoh\_ assignment1\_final.zip
- Your Jupyter notebook should be saved with the results. If you are submitting python scripts, after extracting the ZIP file and executing command `python main.py` in the first level directory, all the generated results and plots you used in your report should appear printed out in a clear manner.
- Include all the references that have been used to complete the assignment.

- If you are working in a team, we expect equal contribution for the assignment. Each team member is expected to make a code-related contribution. Provide a contribution summary by each team member in the form of a table below. If the contribution is highly skewed, then the scores of the team members may be scaled w.r.t the contribution.

Team Member	Assignment Part	Contribution (%)

## Academic Integrity

The standing policy of the Department is that all students involved in any academic integrity violation (e.g., plagiarism in any way, shape, or form) will receive an F grade for the course. The catalog describes plagiarism as “Copying or receiving material from any source and submitting that material as one’s own, without acknowledging and citing the particular debts to the source, or in any other manner representing the work of another as one’s own.”. Refer to the [Office of Academic Integrity](#) for more details.

## Important Information

This project can be done in a team of up to two people.

- All team members are responsible for the project files submission
- No collaboration, cheating, and plagiarism is allowed in assignments, quizzes, the midterms or final project.
- All the submissions will be checked using SafeAssign as well as other tools. SafeAssign is based on the submitted works for the past semesters as well the current submissions. We can see all the sources, so you don't need to worry if there is a high similarity with your Checkpoint submission.
- The submissions should include all the references. Kindly note that referencing the source does not mean you can copy/paste it fully and submit as your original work. Updating the hyperparameters or modifying the existing code is a subject to plagiarism. Your work has to be original. If you have any doubts, send a private post on piazza to confirm.
- All group members and parties involved in any suspicious cases will be officially reported using the Academic Dishonesty Report form. What does that mean?
  - In most cases, the grade for the assignment/quiz/final project/midterm will be 0 and all bonus points will be subject to removal from the final evaluation for all students involved.
  - Those found violating academic integrity more than once throughout their program will receive an immediate F in the course.

Please refer to the [Academic Integrity Policy](#) for more details.

- The report should be delivered as a separate pdf file. You can combine report for Part I, Part II, Part III, & Part IV into the same pdf file. You can follow the [NIPS template](#) as a report structure (NOT mandatory). You may include comments in the Jupyter Notebook; however, you will need to duplicate the results in a separate pdf file.
- All the references can be listed at the end of the report. There is no minimum requirement for the report size, just make sure it includes all the information required.
- For the Bonus part, no report is needed.

## Late Days Policy

You can use up to 7 late days throughout the course that can be applied to any assignment-related due dates. You do not have to inform the instructor, as the late submission will be tracked in UBLearn.

If you work in teams, the late days used will be subtracted from both partners. In other words, you have 4 late days, and your partner has 3 late days left. If you submit one day after the due date, you will have 3 days and your partner will have 2 late days left.

## Important Dates

**Sep 28, Thursday** - Checkpoint is Due

**Oct 12, Thursday** - Final Submission is Due