

Assignment 3

ML Class: CS 4375

April 10, 2025

1 Assignment Policies for CS 4375

The following are the policies regarding this assignment.

1. This assignment needs be done individually by everyone.
2. You are expected to work on the assignments on your own. If I find the assignments of a group of (two or more) students very similar, the group will get zero points towards this assignment. You may possibly also be reported to the judiciary committee.
3. Please use Python for writing code. You can submit the code as a Jupyter notebook
4. For the theory questions, please use either Latex or Word preferably.
5. This Assignment is for 50 points.
6. This will be due on Friday, May 9th. This IS A STRICT DEADLINE AND THERE WILL BE NO EXTENSIONS!!

2 Questions

1. **Naive Bayes (8 points):** Imagine you are given a dataset $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$, where the features $x^{(i)}$ are continuous-valued features. Answer the following questions:
 - (3 points) Explain how you will use the discrete Naive Bayes here. Please provide clear details with step-by-step procedure.
 - (5 points) Derive the expression for the continuous Naive Bayes. Describe what parameters you have, explain how you will derive the MLE estimates of those parameters, and describe how you will perform inference (i.e., given an instance, obtain the class label).
2. **Implement Logistic Regression (8 points):** Implement a regularized Logistic Regression from scratch. Pick any dataset of your choice for this (you can use the UCI Machine Learning repository: <https://archive.ics.uci.edu/ml/index.php> or you can use a dataset we used in the past). Please verify the following:
 - First, start with the regular feature set with regularization = 0. Determine if you are underfitting or overfitting.
 - If you are not overfitting, create polynomial features the way we did in the Linear Regression demo

- Again, train a logistic regression with zero regularization and vary the polynomial degree till you overfit on the train data.
 - Once you overfit on the train data, start adding regularization. What happens with regularization? Can you reduce the overfitting with regularization?
3. **Clustering (6 Points):** Implement k-means and k-medoids algorithms from scratch. To test your algorithm, just use the features from the digits dataset from scikit-learn:

```
sklearn.datasets.load_digits
```

The digits dataset consists of 10 classes (digits from 0 to 9). Run your k-means and k-medoids on this dataset with $k = 10$, and visualize the obtained clustering. Compare the performance of k-means and k-medoids clustering and check if the k-means and k-medoids both reduce the clustering loss.

4. **Variance/Bias Tradeoff (4 points):** Below, you are provided two classifiers, and you need to identify the tradeoff between variance and bias in each case (i.e. for e.g. compare the second classifier to the first and identify if the bias/variance is lower or higher). Provide a justification as to why that is the case.
- Logistic Regression vs Decision Tree (depth 5)
 - Decision Tree vs Random Forest vs Gradient Boosted Tree
 - Logistic Regression vs 1NN classifier
 - 5NN Classifier vs 1NN Classifier
 - Depth 5 Decision Tree vs Depth 10 Decision Tree
 - Random Forest with 100 Trees vs Random Forest with 10 Trees
 - AdaBoost with 10 Decision Stumps vs AdaBoost with 100 Decision Stumps.
 - Linear Regression vs Quadratic (Polynomial degree 2) Regression.

5. **Ensemble Models and Bias/Variance Tradeoff (6 points):** You are provided with a synthetic function defined as:

$$f(x) = \sin(2\pi x)$$

Your goal is to analyze the bias, variance, and mean squared error (MSE) of three models: (i) Decision Tree Regressor, (ii) Random Forest Regressor, (iii) Gradient Boosting Regressor.

First, generate 100 training datasets by sampling $x \sim \text{Uniform}(0, 1)$ and computing $y = f(x) + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.3^2)$. Create a fixed test set of 1000 evenly spaced points in the interval $[0, 1]$. For each model, train on each of the 100 datasets, make predictions on the test set, and compute bias^2 (squared difference between the average prediction and the true function), variance (variance of predictions across datasets), and $\text{mean squared error (MSE)}$.

Plot a few representative predictions along with the true function, the average prediction, and bar plots comparing bias^2 , variance, and MSE across models. Finally, reflect on and explain which model had the lowest bias, which had the lowest variance, and which achieved the lowest MSE. Comment on how *bagging* (Random Forests) and *boosting* (Gradient Boosting) affect the bias-variance tradeoff, and describe when you would expect one method to outperform the other in practice.

6. **Gradient Boosting (8 points):** In class, we studied gradient boosted decision trees with the squared L2 Loss. Specifically, $L(f(x), y) = 1/2(f(x) - y)^2$, and we compute the residual $y' = -\partial L / \partial f$. As a result, the residual in the case of the Square L2 Loss for Gradient Boosting was $y' = (y - f(x))$, where $f(x)$ is the sum of the trees learnt so far. Provide the algorithm for gradient boosting if instead of the squared L2 loss, we use the Logistic Loss and the Hinge Loss.
7. **Neural Networks (10 points):** In class, we looked at the expressions for back propagation in Neural Networks with the Sigmoid activation and the Mean Square Loss. Can you derive the expressions using a Rectified Linear Unit ($\max(x, 0)$) and using the Logistic Loss for Binary Classification.