

Milestone1 Report

Machine Bias in Diabetes Dataset

By

647020016-1 Mr. Yordyod Limrostharn

655380002-2 Miss Natnarin Phonprapreut

655380003-0 Mr. Jakkrin Srihajak

665380009-9 Mr. Pipat Kumthinkaew

Present

Asst. Prof. Dr. Chitsutha Soomlek

SC348 810 Software Development and Project Management
for Data Science and Artificial Intelligence Semester 1/2566

College of Computing

Khon Kaen University

Milestone#1

Task 1: Study the characteristics of the data and identify the quality issues in the selected dataset.

Dataset details

1. This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective is to predict based on diagnostic measurements whether a patient has diabetes.
2. Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.

Features definition

- Pregnancies: Number of times pregnant
- Glucose: Plasma glucose concentration a 2 hours in an oral glucose tolerance test
- BloodPressure: Diastolic blood pressure (mm Hg)
- SkinThickness: Triceps skin fold thickness (mm)
- Insulin: 2-Hour serum insulin (mu U/ml)
- BMI: Body mass index (weight in kg/(height in m)²)
- DiabetesPedigreeFunction: Diabetes pedigree function
- Age: Age (years)
- Outcome: Class variable (0 or 1)

Sources:

(a) Original owners: National Institute of Diabetes and Digestive and Kidney Diseases

(b) Donor of database: Vincent Sigillito (vgs@aplcn.apl.jhu.edu)

Research Center, RMI Group Leader

Applied Physics Laboratory

The Johns Hopkins University

Johns Hopkins Road

Laurel, MD 20707

(301) 953-6231

(c) Date received: 9 May 1990

Identify the quality issues in the selected dataset

1. Dataset has some columns that include missing values.
2. Imbalance data Ex. Number of sampling in each classes class 0 (no diabetes) = 500, 1 (has diabetes) = 268
3. Sampling data were collected from specific 1 country, 1 gender and small size of sample in each class.
4. Some features in the dataset have many outlier values.

Task 2: Define the goals and a suitable measure for the quality issues.

Goal and Suitable measure for the quality issues

1. Arrange data to decrease bias in real dataset.
2. Develop model for predicting fairness predicted result.
3. Decrease data bias which is issues from collecting data processes.
 - a. Missing data
 - b. Data imbalance
 - c. Data outliers
4. Analyse model bias.

Suitable measure for the quality issues

1. Countplot : use for checking number of classes in this dataset.
2. Boxplot : use for checking outliers in data.
3. Dataframe : showing and deal with missing values in each features.
4. Confusion matrix : showing model predicted results comparing with actual results.
5. Classification report : Precision, Recall, F1-score, Support and Accuracy.

Task 3: Explore different kinds of machine learning models developed with different modeling techniques. Then, choose the machine learning techniques, implement the models using scikit-learn, and train the models.

1. Naïve Bayes

Naive Bayes is a probabilistic machine learning algorithm commonly used for classification and text analysis tasks, particularly in situations where you need to make predictions based on a set of features. It's based on Bayes' theorem and is considered "naive" because it makes a strong and often unrealistic assumption that the features used to make predictions are conditionally independent given the class label. Despite this simplification, Naive Bayes can work surprisingly well in many real-world applications and is computationally efficient

Bayes' Theorem: The foundation of the Naive Bayes model is Bayes' theorem, which describes the probability of an event based on prior knowledge of conditions related to that event. In a classification context, Bayes' theorem can be expressed as:

$$P(y | X) = (P(X | y) * P(y)) / P(X)$$

- $P(y | X)$: The probability of class y given features X (the posterior probability).
- $P(X | y)$: The probability of observing features X given class y (the likelihood).
- $P(y)$: The prior probability of class y .
- $P(X)$: The probability of observing the features X .

Naive Bayes classifiers are simple, fast to train, and can work surprisingly well in various tasks like spam email detection, sentiment analysis, and document categorization. However, their performance may suffer when the independence assumption is significantly violated, or when dealing with high-dimensional data. In such cases, more complex models like decision trees, random forests, or neural networks may be more suitable.

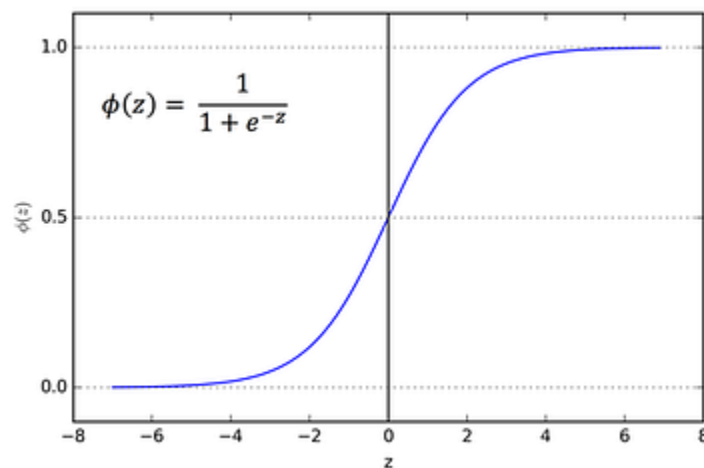
Table 1 : Naïve Bayes Classification Report

	precision	recall	f1-score	support
Class 0	0.77	0.81	0.79	93
Class 1	0.63	0.57	0.60	54
			accuracy	0.72

2. Logistic regression

Logistic Regression is a statistical and machine learning model used for binary classification tasks, where the goal is to predict one of two possible outcomes (classes), typically denoted as 0 and 1, based on one or more predictor variables or features. Despite its name, logistic regression is used for classification, not regression.

Sigmoid Function: Logistic regression uses a sigmoid (logistic) function to model the relationship between the input features and the probability of belonging to the positive class (class 1). The sigmoid



function is defined as:

Where:

- **z**: The linear combination of input features and their corresponding coefficients (weights).
- **$\sigma(z)$** : The output of the sigmoid function, which represents the estimated probability that the input example belongs to the positive class.

Logistic Regression is widely used in various applications such as spam email classification, medical diagnosis, credit scoring, and many other binary classification tasks. It's simple, interpretable, and can serve as a baseline model for more complex classification problems.

Table 2 : Logistic regression Classification Report

	precision	recall	f1-score	support
Class 0	0.78	0.87	0.82	110
Class 1	0.62	0.46	0.53	50
			accuracy	0.74

3. Linear Discriminant Analysis (LDA)

A statistical method used for dimensionality reduction and classification tasks in machine learning and pattern recognition. LDA is particularly useful when we have a dataset with multiple classes, and we want to find a linear combination of features that best separates these classes. It's often used in supervised learning contexts, such as in pattern recognition and classification tasks.

Assume that the data follows a multivariate normal distribution, LDA works in these steps:

- **Compute Class Means:** For each class in our dataset, calculate the mean vector (centroid) of the feature values. These mean vectors represent the "average" data point for each class.
- **Compute Scatter Matrices:** There are two types of scatter matrices to calculate in LDA:
 - a. **Within-Class Scatter Matrix (Sw):** This matrix measures the spread of data within each class. It is computed by calculating the covariance matrix for each class and then summing them up. The formula is:

$$S_w = \sum (X - \mu_i) (X - \mu_i)^T \text{ for all classes } i$$

Where:

X represents the data points in a class.

μ_i is the mean vector for class i .

- b. **Between-Class Scatter Matrix (Sb):** This matrix measures the spread between class centroids. It is computed by finding the covariances between class means and weighting them by the number of samples in each class. The formula is:

$$S_b = \sum N_i (\mu_i - \mu) (\mu_i - \mu)^T \text{ for all classes } i$$

Where:

N_i is the number of samples in class i .

μ is the overall mean vector.

- **Compute Eigenvectors and Eigenvalues:** calculate the eigenvalues and eigenvectors of the matrix $S_w^{-1} * S_b$. These eigenvectors represent the directions (linear combinations of features) along which the data is most separable.
- **Select Discriminant Components:** Sort the eigenvectors by their corresponding eigenvalues in descending order. The eigenvectors with the highest eigenvalues contain the most discriminatory information. Choose a subset of these eigenvectors (discriminant components) to reduce the dimensionality of your data.
- **Project Data:** Project the original data onto the selected discriminant components to create a lower-dimensional representation of the data. This lower-dimensional space can be used for classification or visualization.

Table 3: LDA Classification Report

	precision	recall	f1-score	support
Class 0	0.85	0.88	0.87	95
Class 1	0.77	0.71	0.73	51
			accuracy	0.82

4. Support Vector Machine (SVM)

Support Vector Machine is a powerful and versatile machine learning algorithm used for both classification and regression tasks. SVM is particularly effective in scenarios where you need to find a hyperplane that best separates data points of different classes in a high-dimensional feature space. SVM aims to maximize the margin between the two classes while minimizing the classification error. Here's an explanation of how SVM works:

In this solution we use SVM for deal with classification data (Diabetes dataset). After process of explore data analysis we found that our selected dataset have some features that include missing value, outliers (Ex. Insulin and DiabetesPedigreeFunction). So we decided to remove some values in some features in dataset (column name “Glucose” and “BloodPressure”) before bring dataset to train in SVM model . Table 4 below showing result from SVM training model.

Table 4 : SVM Classification Report

	precision	recall	f1-score	support
Class 0	0.91	0.75	0.82	114
Class 1	0.46	0.76	0.57	33
			accuracy	0.75

Conclusion

In conclusion, our group explored and analyzed the diabetes dataset, and we can deal with some of the features that include missing values (zero value). We found that the dataset has many features having outlier values. However, we cannot find any solution or techniques to that problem because it is a real-world dataset and a source from the laboratory. We chose four models for training and prediction, including Naïve Bayes, Logistic regression, Linear Discriminant Analysis, and Support Vector Machine. The result of each model after training tells us that the imbalance of data made our model get the biased result by classification report of support every outcome of each model, class 0 (which have numbers of sampling data more than class 1) has learned by model more than class 1.

References

Diabetes dataset : <https://www.kaggle.com/datasets/akshaydattatraykhare/diabetes-dataset?datasetId=2527538&sortBy=voteCount>