

Comparative Approaches for Classification of Diabetes Mellitus Data

Study of machine learning algorithms and application on the Pima Indians Diabetes
Dataset. **Alexander Roque Rodrigues**

A dissertation presented for the degree of
Bachelor in Computer Science Department of Computer Science

Smt. Parvatibai Chowgule College of Arts and Science
India

18 August 2019

Contents

I	Acknowledgements	3
II	Background	4
III	The Pima Indians	5
1	Introduction	5
2	Diabetes Mellitus	6
IV	Review of Literature	7
V	Data Preprocessing	8
1	Feature Extraction	8
1.1	Filters	8
1.2	Information Gain	8
1.3	Chi-square test	8
1.4	Fisher score	8
1.5	Correlation coefficient	8
1.6	Variance threshold	8
2	Wrappers	8
VI	Machine Learning Algorithms	9
1	Linear Regression	10
1.1	Introduction to Linear Regression	10
2	Logistic Regression	11
3	K-Nearest Neighbours	12
3.1	Finding Optimum Number of Clusters	12
4	Decision Tree	13
4.1	Introduction	13
4.2	Overfitting in Trees	13
5	Random Forest	15
6	Gradient Boosting	16

7	Support Vector Machine	17
7.1	Linear Support Vector Machines	17
8	Perceptron	18
9	Multilayered Perceptron	19
VII	Real World Application	20
VIII	Building Information Systems for Prediction	21
1	Introduction	21
2	Feature Selection	21
3	Models	22
4	Conclusion	23
IX	Conclusion	24
X	Bibliography	25

Part I

Acknowledgements

I, Alexander Roque Rodrigues, would like to acknowledge the role of the Department of Computer Science in guiding me and helping me to achieve the completion of this project. I would also like to acknowledge the efforts put in by my project guide Ms. Ashweta Fondekar, who assisted me throughout the span of the project and helped me achieve my goal. I also would like to thank my parents for their support and others who have indirectly assisted me in this project.

Part II

Background

Data is everywhere. International initiatives coupled with global disruptive innovation are the leading causes for pushing datafication forward.

Datafication refers to the modern-day trend of digitalizing (or datafying) every aspect of life.

This data creation is enabling the transformation of data into new and potentially valuable forms. Entire municipalities are being incentivized to become smarter. In the not too distant future, our towns and cities will collect thousands of variables in real time to optimize, maintain, and enhance the quality of life for entire populations. One would reasonably expect that as well as managing traffic, traffic lights may also collect other data such as air quality, visibility, and speed of traffic. As a result of big data from connected devices, embedded sensors, and the IoT, there is a global need for the analysis, interpretation, and visualization of data.

Part III

The Pima Indians

1 Introduction

For centuries the Pima Indians of the southern Arizona have been targeted by a genetic quirk that has caused the people of the Pima Indian families to have long lived with kidney disease, blindness and amputations that attend diabetes. And, of course, death. For decades researchers have tried to understand why the Pima tribe faces such a high risk of developing genetic diabetes.

According to Greg Morago, a staff writer from The Courant,

”I have many inducements to be more careful. There’s my grandmother, who lost toes on both feet last year after an infection. There’s my mother, who pricks her fingers every day to test her blood sugar levels. There’s my younger sister, who last year got diabetes and whose teenage son only recently learned he’s borderline diabetic.”

Scientists have learned that diabetes develops when a person’s body doesn’t use insulin effectively. Volunteers from the Pima Indians tribe continue to help support research not unlike the recent clinical trials that linked lifestyle changes to preventing diabetes. These people have been living through, and dying from, diabetes have been placed in the forefront of diabetes prevention.

2 Diabetes Mellitus

More commonly referred to as "diabetes" – a chronic disease associated with abnormally high levels of the sugar glucose in the blood. Diabetes is due to one of two mechanisms:

Inadequate production of insulin (which is made by the pancreas and lowers blood glucose), or Inadequate sensitivity of cells to the action of insulin. The two main types of diabetes correspond to these two mechanisms and are called insulin dependent (type 1) and non-insulin dependent (type 2) diabetes. In type 1 diabetes there is no insulin or not enough of it. In type 2 diabetes, there is generally enough insulin but the cells upon which it should act are not normally sensitive to its action.

The signs and symptoms of both types of diabetes include increased urine output and decreased appetite as well as fatigue. Diabetes is diagnosed by blood glucose testing, the glucose tolerance test, and testing of the level of glycosylated hemoglobin (glycohemoglobin or hemoglobin A1C). The mode of treatment depends on the type of the diabetes.

The major complications of diabetes include dangerously elevated blood sugar, abnormally low blood sugar due to diabetes medications, and disease of the blood vessels which can damage the eyes, kidneys, nerves, and heart.

Part IV

Review of Literature

Part V

Data Preprocessing

1 Feature Extraction

1.1 Filters

1.2 Information Gain

1.3 Chi-square test

1.4 Fisher score

1.5 Correlation coefficient

1.6 Variance threshold

2 Wrappers

2.1 Recursive Feature Elimination

2.2 Sequential Feature Selection Algorithms

2.3 Genetic Algorithms

3 Embedded Methods

3.1 L1 (LASSO) Regularization

3.2 Decision Tree

Part VI

Machine Learning Algorithms

1 Linear Regression

1.1 Introduction to Linear Regression

Linear regression is a forecasting technique that can be use to predict the future of a number series based on the historic data given. The perks of using a linear regression model are as follows:

- produces decent and easy to interpret results.
- is computationally inexpensive.
- conversion of algorithm into code does not take much effort or time.
- numeric values as well as nominal values support is offered.

However, a major drawback of linear regression is that it **poorly models nonlinear data**.

Considering a dataset that has values ranging from $X = \{x_1 + x_2 + x_3 + \dots + x_n\}$ where all the entries of the dataset are real numbers. Each x_i is associated with a corresponding value of y_i from the dataset $Y = \{y_1 + y_2 + y_3 + \dots + y_n\}$.

The most basic equation for linear regression can be expressed via this simple equation.

$$y = \beta_0 x + \beta_1 + \epsilon$$

So to minimize the error in the predictions, a way to calculate the error should be formulated. A loss function in machine learning is simply a measure of how different the predicted value is from the actual value. The Quadratic Loss Function to calculate the loss or error in our linear regression model. It can be defined as:

$$L(x) = \sum_{i=1}^n (y_i - p_i)^2$$

Therefore using the method of Least Squares, we can find the values of β_0 and β_1 .

$$\beta_0 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

The linear regression model with an error value close to 1.00 indicates a perfect model and those with values closer to 0.00 indicates a model that delivers poor performance.

2 Logistic Regression

Even if called regression, this is a classification method which is based on the probability for a sample to belong to a class. As our probabilities must be continuous in \mathbb{R} and bounded between $(0, 1)$, it's necessary to introduce a threshold function to filter the term z . The name logistic comes from the decision to use the sigmoid (or logistic) function:

3 K-Nearest Neighbours

The k-means algorithm is based on the (strong) initial condition to decide the number of clusters through the assignment of k initial centroids or means:

Then the distance between each sample and each centroid is computed and the sample is assigned to the cluster where the distance is minimum. This approach is often called minimizing the inertia of the clusters, which is defined as follows:

The process is iterative—once all the samples have been processed, a new set of centroids K (1) is computed (now considering the actual elements belonging to the cluster), and all the distances are recomputed. The algorithm stops when the desired tolerance is reached, or in other words, when the centroids become stable and, therefore, the inertia is minimized. Of course, this approach is quite sensitive to the initial conditions, and some methods have been studied to improve the convergence speed. One of them is called k-means++ (Karteeika Pavan K., Allam Appa Rao, Dattatreya Rao A. V., and Sridhar G.R., Robust Seed Selection Algorithm for K-Means Type Algorithms, International Journal of Computer Science and Information Technology 3, no. 5, October 30, 2011), which selects the initial centroids so that they are statistically close to the final ones. The mathematical explanation is quite difficult; however, this method is the default choice for scikit-learn, and it's normally the best choice for any clustering problem solvable with this algorithm.

3.1 Finding Optimum Number of Clusters

One of the most common disadvantages of k-means is related to the choice of the optimal number of clusters. An excessively small value will determine large groupings that contain heterogeneous elements, while a large number leads to a scenario where it can be difficult to identify the differences among clusters. Therefore, we're going to discuss some methods that can be employed to determine the appropriate number of splits and to evaluate the corresponding performance.

4 Decision Tree

4.1 Introduction

Decision trees are flowcharts that represent the decision-making process as rules for performing categorization. Decision trees start from a root and contain internal nodes that represent features and branches that represent outcomes. As such, decision trees are a representation of a classification problem. Decision trees can be exploited to make them easier to understand. Each decision tree is a disjunction of implications (i.e., if-then statements), and the implications are Horn clauses that are useful for logic programming. A Horn clause is a disjunction of literals. On the basis that there are no errors in the data in the form of inconsistencies, we can always construct a decision tree for training datasets with 100% accuracy. However, this may not roll out in the real world and may indicate overfitting, as we will discuss.

Classifying an example involves subjecting the data to an organized sequence of tests to determine the label. Trees are built and tested from top to bottom as so:

1. Start at the root of the model
2. Wait until all examples are in the same class
3. Test features to determine best split based on a cost function
4. Follow the branch value to the outcome
5. Repeat number 2
6. Leaf node output

The central question in decision tree learning is which nodes should be placed in which positions, including the root node and decision nodes. There are three main decision tree algorithms. The difference in each algorithm is the measure or cost function for which nodes, or features, are selected. The root is the top node. The tree is split into branches; evaluated through a cost function; and a branch that doesn't split is a terminal node, decision, or leaf.

Decision trees are useful in the way that acquired knowledge can be expressed in an easy to read and understandable format (see Figure 4-1). It mimics human decision-making whereby priority—determined by feature importance, relationships, and decisions—is clear. They are simple in the way that outcomes can be expressed as a set of rules. Figure 4-1. Decision tree of $n = 2$ nodes Decision trees provide benefits in how they can represent big datasets and prioritize the most discriminatory features. If a decision tree depth is not set, it will eventually learn the data presented and overfit. It is recommended to set a small depth for decision tree modeling. Alternatively, the decision tree can be pruned, typically starting from the least important feature, or the incorporation of dimensionality reduction techniques.

4.2 Overfitting in Trees

Overfitting is a common machine learning obstacle, and not limited to decision trees. All algorithms are at risk of overfitting, and a variety of techniques exist to overcome this problem. Random forest or jungle decision trees can be extremely useful in this.

Pruning reduces the size of a decision tree by removing features that provide the least information. As a result, the final classification rules are less complicated and improve predictive accuracy. The accuracy of a model is calculated as the percentage of examples in the test dataset that is classified correctly.

- True Positive / TP: Where the actual class is yes, and the value of the predicted class is also yes.
- False Positive / FP: Actual class is no, and predicted class is yes
- True Negative / TN: The value of the actual class is no, and the value of the predicted class is no
- False Negative / FN: When the actual class value is yes, but predicted class is no
- Accuracy: $(\text{correctly predicted observation}) / (\text{total observation}) = (TP + TN) / (TP + TN + FP + FN)$
- Precision: $(\text{correctly predicted Positive}) / (\text{total predicted Positive}) = TP / (TP + FP)$
- Recall: $(\text{correctly predicted Positive}) / (\text{total correct Positive observation}) = TP / (TP + FN)$

Classification is a common method used in machine learning; and ID3 (Iterative Dichotomizer 3), C4.5 (Classification 4.5), and CART (Classification And Regression Trees) are common decision tree methods where the resulting tree can be used to classify future samples.

5 Random Forest

6 Gradient Boosting

7 Support Vector Machine

Support Vector Machines is an algorithm that is capable of handling linear as well as data that occurs non-linearly. For example, for a long time, SVMs were the best choice for MNIST dataset classification, thanks to the fact that they can capture very high non-linear dynamics using a mathematical trick, without complex modifications in the algorithm.

7.1 Linear Support Vector Machines

Let us consider a dataset of features we want to classify.

$$X = \{x_1, x_2, x_3, \dots, x_n\}$$

For the target variable, we will consider the dataset Y , with target outcomes as $\{0, 1\}$ indicating a true or false condition.

$$Y = \{y_1, y_2, y_3, \dots, y_n\}$$

8 Perceptron

The perceptron is the foundation of neural networks.

9 Multilayered Perceptron

Part VII

Real World Application

Part VIII

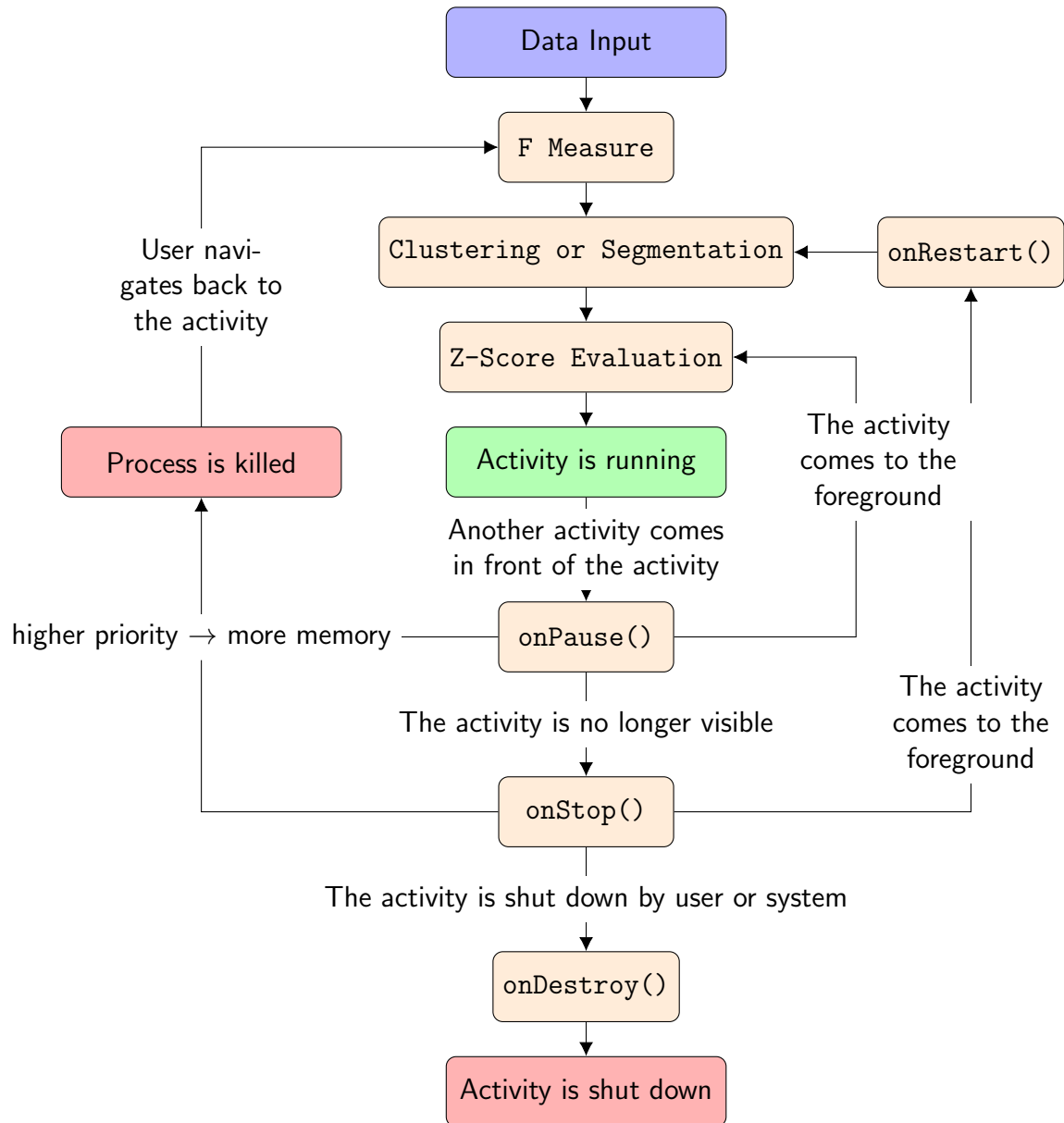
Building Information Systems for Prediction

1 Introduction

In today's world we have many algorithms and multiple datasets along with sufficient data available for testing and training the algorithms.

2 Feature Selection

3 Models



4 Conclusion

Part IX

Conclusion

Part X

Bibliography

References

- [1] Wei M, Gibbons LW, Mitchell TL *et al.* (1999) The Association between cardiorespiratory fitness and impaired fasting glucose and type 2 diabetes mellitus in men. *Ann Intern Med* **130**, 427-34.
- [2] Jr., W. C. S. (2017, January 26). Definition of Diabetes mellitus. Retrieved November 17, 2019, from <https://www.rxlist.com/script/main/art.asp?articlekey=2974>.
- [3] Zou, Q., Qu, K., Luo, Y., Yin, D., Ju, Y.,& Tang, H. (2018, November 6). Predicting Diabetes Mellitus With Machine Learning Techniques. Retrieved November 17, 2019, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6232260/>.