

Programming and processing of big data using python language in medicine

Otabek Ergashev^{1,*}, Nurillo Mamadaliev¹, Sardorbek Khonturaev¹, and Muzaffar Sobitov¹

¹Fergana branch of TUIT named after Muhammad Al-Khwarizmi, Fergana, Uzbekistan

Abstract. This article is devoted to the use and further application of Python libraries in the medical industry. These libraries include NumPy, Pandas, Scikit-learn, Keras and TensorFlow, Matplotlib, Seaborn, and Plotly. On the example of the Keras library, the problem associated with the use of medical data analysis was considered. Keywords: Neural network, Big data, data analysis, big data processing, Python library.

1 Introduction

In recent years, the use of big data in the medical industry has become increasingly prevalent. With the help of big data analytics, healthcare professionals can analyze large amounts of data to improve patient outcomes, identify trends, and develop new treatments. Python has emerged as one of the most popular programming languages for big data processing in medicine due to its versatility, ease of use, and powerful data analysis capabilities.

Python has a wide range of libraries and frameworks that are specifically designed for data analysis and machine learning. These libraries include NumPy, Pandas, Scikit-learn, Keras, and TensorFlow. These libraries provide healthcare professionals with the tools they need to analyze large amounts of data, including medical images and clinical trial results.

With the help of Python, healthcare professionals can use medical images to develop personalized treatments for patients. For example, medical images can be used to identify tumors, measure the size and shape of organs, and identify abnormalities in blood vessels. Python's ability to handle unstructured data makes it particularly useful for processing medical images, which are often large, complex, and unstructured [1-19].

Python is also useful in clinical trial data analysis. Clinical trials are a critical part of the drug development process, and they play a significant role in the approval of new drugs. Python can be used to analyze clinical trial data, including patient demographics, treatment efficacy, and adverse events. Python's powerful data analysis capabilities make it an ideal choice for clinical trial data analysis.

In addition to its data analysis capabilities, Python is also an excellent tool for data visualization. Data visualization is an essential part of big data processing, as it helps healthcare professionals to understand complex data sets and identify trends. Python has several libraries like Matplotlib, Seaborn, and Plotly that can be used to create visualizations of medical data.

* Corresponding author: ergashev1984@gmail.com

One of the key benefits of using Python for big data processing in medicine is its ability to handle unstructured data. Unstructured data is data that does not have a predefined data model or is not organized in a predefined manner. Examples of unstructured data in medicine include medical images, clinical notes, and electronic health records. Python has libraries like Natural Language Toolkit (NLTK) and OpenCV that can be used to process unstructured data.

Python is also well-suited for graph processing, which involves analyzing relationships between data points. Graph processing is particularly useful in medicine for analyzing social networks of patients, identifying outbreaks of infectious diseases, and analyzing the spread of diseases. Python has libraries like NetworkX and igraph that can be used for graph processing.

We will consider the problem associated with the application of medical data analysis, namely, diagnosing the risk of diabetes mellitus based on the patient's condition. To do this, we will write a simple neural network that solves an important practical problem.

2 Method of study

The experimental data, taken from the medical archive, is a comma-separated value file, in which each line corresponds to one questionnaire. There are nine parameters in the experimental data set consisting of anonymous records. The last one, the target one, shows whether the patient had diabetes mellitus or not (respectively, 1 or 0). Eight other parameters also have numerical values:

1. Number of pregnancies (all source patients are women at least 21 years of age).
2. Plasma glucose concentration 2 hours after administration in an oral glucose tolerance test.
3. Diastolic blood pressure (mm Hg).
4. The thickness of the skin fold in the triceps area (mm).
5. Serum insulin concentration ($\mu\text{U/ml}$).
6. Body mass index ($\text{weight in kg}/(\text{height in m})^2$).
7. A function that describes the genetic predisposition to diabetes (diabetes pedigree).
8. Age (years).

Loading required modules and data

The Keras library allows you to run neural networks with a minimum number of operations. Sequential Sequential from the `keras.models` module is used as a neural network model with `keras.layers` of the Dense type specified.

The Keras library is a high-level interface for creating neural networks. Keras is written in Python and runs on top of lower-level solutions like TensorFlow, CNTK, and Theano. Due to this, the program code is not only powerful, but also extremely compact (1).

```
from keras.models import Sequential  
from keras.layers import Dense  
import numpy
```

 (1)

For subsequent reproducibility of the results, we fix the random number generator using the `random.seed()` function from the `numpy` library. Reading data from the dataset (2):

```
numpy.random.seed(2)  
dataset = numpy.loadtxt("prima-indians-diabetes.csv", delimiter=",")
```

 (2)

Divide the data into feature matrix X and target variable vector Y (the last column of the dataset) (3):

```
X = dataset[:,0:8]  
Y = dataset[:,8] (3)
```

3 Result of study

Neural Network: Creating a Model

Create a neural network model (4):

```
model = Sequential() (4)
```

Let describe the structure of the neural network model. Let's define the input, output and hidden layers. Our neural network will have a dense (Dense) structure - each neuron is connected to all the neurons of the next layer. The output layer will consist of a single neuron that determines the likelihood of diabetes (Fig. 1).

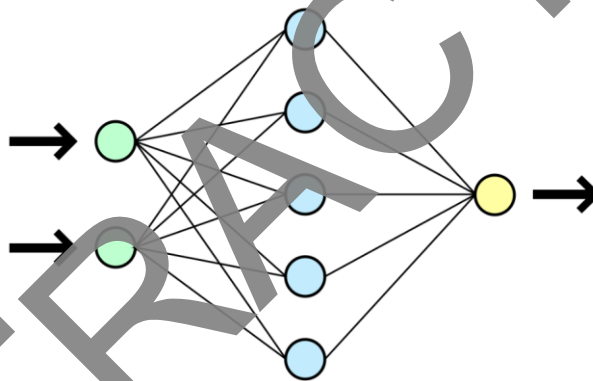


Fig. 1. Structure of the neural network of the output layer that determines the probability of diabetes.

The layer is added to the model using the add() method. For the input layer, you must specify the number of input_dim features, which in our case is 8 (5):

```
model.add(Dense(12, input_dim=8, activation='relu')) (5)
```

If the feature sets form a multidimensional table, then instead of the input_dim parameter, you can use the input_shape parameter, which takes a tuple with the number of elements in each of the dimensions.

We will use the ReLU function as an activation function for all layers, except for the output one. For the output layer, we use the sigmoid function to determine the final probability of disease risk (Fig. 2).

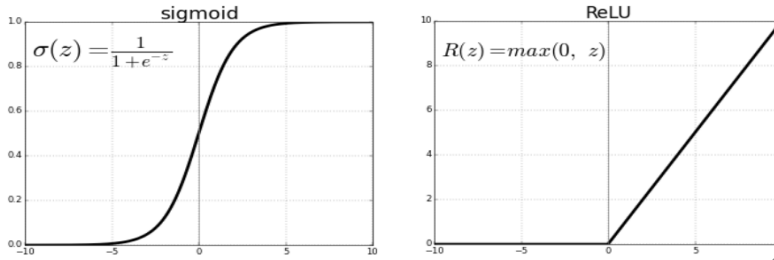


Fig. 2. Graph of the sigmoid function for determining the final probability of disease risk.

Let's create three hidden layers and one output layer of our neural network (6):

```
model.add(Dense(15, activation='relu'))  
model.add(Dense(8, activation='relu'))  
model.add(Dense(10, activation='relu'))  
model.add(Dense(1, activation='sigmoid'))
```

 (6)

The first numbers passed to Dense are the numbers of neurons experimentally optimized as a result of varying the structure of the neural network. You can change the number of hidden layers and the neurons they contain to improve the predictive quality of the model.

Before you start training the model, you need to compile it using the compile() method (7):

```
model.compile(loss="binary_crossentropy", optimizer="adam",  
metrics=['accuracy'])
```

 (7)

The method is passed three parameters:

loss - loss function - an object that the model seeks to minimize;

optimizer - optimizer, we use the built-in adam stochastic optimization method described in the publication by Diederik Kingma and Jimmy Ba;

metrics - list of optimization metrics, 'accuracy' metric is used for classification tasks.

Neural network: training and evaluation of the result

To train the neural network, we use the fit() method (8):

```
model.fit(X, Y, epochs = 1000, batch_size=10)
```

 (8)

The epochs parameter - "epochs" - the number of neural network passes over all dataset records (selected based on how quickly the model approaches the desired predictive accuracy with each new pass), batch_size - the number of sample objects taken in one step. During the training process, the API will output the corresponding lines with the values of the loss function and the metric for each of the epochs.

Let's evaluate the result of neural network training. The evaluate() method returns the loss function and metric values for the trained model (9):

```
scores = model.evaluate(X, Y)  
print("\n%s: %.2f%%", % (model.metrics_names[1], scores[1]*100)
```

 (9)

The last line in the formatted form displays the accuracy of the forecast for our model for the given accuracy metric(10):

```
acc: 87.89%
```

 (10)

4 Conclusions

Python is mainstream for processing large amounts of data in medicine, so big data processing using Python has the potential to revolutionize the medical industry. By the reason it provides a wide range of tools and libraries for working with data in various formats, including images and experimental results and usage of Python, healthcare professionals can analyze large amounts of data to improve patient outcomes, identify trends, and develop new treatments. Python's versatility, ease of use, and powerful data analysis capabilities make it an ideal choice for big data processing in medicine.

References

1. W. Raghupathi, *Big data analytics in healthcare: Promise and potential*. Health Information Science and Systems (2014)
2. A. Dash et al, *Big data in healthcare: management, analysis and future prospects*. Journal of Big Data (2019)
3. C. Ma, *Integrated Pharmacy Research and Practice* **4**. 91-99 (2015)
4. B. Ekmekci, *An Introduction to Programming for Bioscientists: A Python-Based Primer*. PLOS Computational Biology (2016)
5. DP. Nazareth, *First application of quantum annealing to IMRT beamlet intensity optimization*. Spaans JD Phys Med Biol. (2015)
6. Rebentrost P. *Quantum support vector machine for big data classification*. Phys Rev Lett. (2014)
7. I. Siddikov, O. Porubay, E3S Web of Conferences **304** (01001) (2021), <https://doi.org/10.1051/e3sconf/202130401001>
8. I. Siddikov, O. Porubay, O. Minalilov, Journal of Physics: Conference Series **2373**. 8, 082014 (2022), doi: [10.1088/1742-6596/2373/8/082014](https://doi.org/10.1088/1742-6596/2373/8/082014)
9. O. Porubay, I. Siddikov, K. Madina, *Algorithm for optimizing the mode of electric power systems by active power* 2022 International Conference on Information Science and Communications Technologies (ICISCT). – IEEE, 2022. – pp. 1-4, doi: [10.1109/ICISCT55600.2022.10146996](https://doi.org/10.1109/ICISCT55600.2022.10146996)
10. I.K. Siddikov O.V. Pombay, AIP Conference Proceedings **2432**. 1 (2022), <https://doi.org/10.1063/5.0089473>
11. I. Siddikov, et.al, International Journal of Electrical and Computer Engineering **14**. 1. 184-191 (2014), doi: <http://doi.org/10.11591/ijece.v14i1.pp184-191>
12. D.M. Umurzakova, IIUM Engineering Journal **22**, 1. 287-297 (2021). <https://doi.org/10.31436/iiumej.v22i1.1415>.
13. D.M. Umurzakova, *Mathematical Modeling of Transient Processes of a Three-pulse System of Automatic Control of Water Supply to the Steam Generator When the Load Changes* // 14th International IEEE Scientific and Technical Conference Dynamics of Systems, Mechanisms and Machines, Dynamics. – Russia 2020. DOI: [10.1109/Dynamics50954.2020.9306117](https://doi.org/10.1109/Dynamics50954.2020.9306117)
14. Umurzakova D.M., *Neuro-fuzzy Control Algorithm of Dynamic Objects with Uncertainty of a Priori Information* / International conference on information science and communications technologies applications, trends and opportunities (ICISCT 2020). Tashkent University of information technologies named after Muhammad al-Khwarizmi. –Tashkent. 4-6 November, 2020. DOI: [10.1109/ICISCT50599.2020.9351462](https://doi.org/10.1109/ICISCT50599.2020.9351462)

15. Umurzakova D.M., *Synthesis of an Automatic Control System of a Multidimensional Dynamic Object Under Information Uncertainty* / International conference on information science and communications technologies applications, trends and opportunities (ICISCT 2020). Tashkent University of information technologies named after Muhammad al-Khwarizmi. –Tashkent. 4-6 November, 2020.
DOI: 10.1109/ICISCT50599.2020.9351451.
16. D.M. Umurzakova, IIUM Engineering Journal **22**, 1, 287-297 (2021).
<https://doi.org/10.31436/iiumej.v22i1.1415>.
17. Usarov, M., Usarov, D., Usanov, F., Askarhodjaev, S., Shamsiyev, D. (2023) PIJournal of Construction and Engineering Technology, **1(2)**, 1-6.
18. Ismoilov, M. M., Obidova, G., Juraeva, M., Meliqo'ziev, A., Maxkamova, D., Toshpo'latova, M. (2024) E3S Web of Conferences **508**, 06002.
<https://doi.org/10.1051/e3sconf/202450806002>
19. Rayimdjanova, O., Mukhtarov, F., Ismoilov, M. M., Abdusamatov, H., Abdullaeva, M., Madaminov, M. (2024) E3S Web of Conferences **508**, 01004.
<https://doi.org/10.1051/e3sconf/202450801004>

RETRACTED