



Summer Fellowship Report

On

Python for Machine learning

Submitted by

Anvita Thadavoose Manjummel

Vellore Institute of Technology, Chennai

Under the guidance of

Prof. Kannan M. Moudgalya

Chemical Engineering Department IIT Bombay

Mentors from IIT Bombay

Ms. Nirmala Venkat

Mentors from VIT Chennai

Dr. Subbulakshmi T Dr. Bhargavi Rentachintala

September 19, 2025

Acknowledgment

I would like to extend my sincere gratitude to my mentors at IIT Bombay - Ms. Nirmala for her constant support throughout the fellowship. Her guidance, suggestions, and encouragement at every stage of script and code development were invaluable in shaping this tutorial series.

I am also deepl grateful to the faculties at VIT Chennai, especially Dr. Subbulakshmi T and Dr. Bhargavi Rentachintala, for their technical expertise during the code development phase, and their valuable inputs on the content.

Last but not least, I thank my fellow intern T Harini, for her cooperation and support.

Contents

1	Intro	oduction	3
2	Tutorials on Python for Machine Learning		4
	2.1	Setup Python Environment for Machine Learning	4
	2.2	K Nearest Neighbor Classification	5
	2.3	K Nearest Neighbor Regression	
	2.4	Linear Regression	6
	2.5	Logistic Regression Binary Classification	7
	2.6	Logistic Regression MultiClass Classification	7
	2.7	Decision Tree	8
	2.8	Artificial Neural Networks	8
	2.9	Support Vector Machine	9
	2.10	K Means Clustering	10
	2.11	Random Forest	11
3	Code	e and Script Writing	12
4	Crea	ating slides	14
5	Phases of my Spoken Tutorial		15
	5.1	Outline	15
	5.2	Script	15
	5.3	Slides	15
	5.4	Novice Check	16
	5.5	Domain Check	16
	5.6	Admin Check	16
	5.7	Recording	16
6	Conclusion		17

Introduction

The Spoken Tutorial Project is an online learning initiative that provides high-quality video tutorials on Free and Open Source Software (FOSS). The project aims to make learning software skills simple, accessible, and affordable to all, regardless of location or background. The tutorials are designed to help learners study at their own pace, using only a computer and an internet connection. With minimal prerequisite knowledge required, the platform is especially suited for beginners as well as self-learners.

A key feature of Spoken Tutorials is the emphasis on practical, hands-on learning. The majority of each tutorial is devoted to live demonstration of the topic, ensuring that learners not only understand the concepts but also apply them in real time. To support this, all necessary resources—such as source code and data files used in the tutorials—are made available alongside the videos. To reinforce the concepts, practice problems are included at the end of each tutorial, encouraging learners to test their understanding independently.

The platform also provides an interactive support system through tutorial-specific forums. Here, learners can post queries and receive clarifications from domain experts, ensuring that doubts are resolved quickly and effectively. In addition, Spoken Tutorial offers a certification program, where learners can obtain certificates by successfully completing online assessments, thereby validating their skills. Over the years, the Spoken Tutorial project has grown into a national-level initiative with wide outreach across schools, colleges, and universities. By promoting FOSS tools and encouraging self-paced learning, it continues to bridge gaps in digital education and empower learners with valuable technical skills.

Tutorials on Python for Machine Learning

The Python for Machine Learning series in the Spoken Tutorial Project is designed to introduce learners to the fundamental concepts of machine learning using Python. Python is one of the most widely used programming languages in the field of Artificial Intelligence and Machine Learning due to its simplicity, readability, and the availability of powerful libraries such as NumPy, Pandas, Matplotlib, and Scikit-learn.

This tutorial series follows the step-by-step approach of the Spoken Tutorial methodology, ensuring that even learners with minimal programming background can follow along. The scripts and demonstrations are carefully designed to explain each concept in a clear and practical manner. The series not only covers the theoretical foundations of machine learning but also emphasizes hands-on implementation, allowing learners to write clean code, analyze datasets, and build simple ML models.

By the end of the series, learners will have gained a solid understanding of Python's role in machine learning and will be equipped with the basic skills needed to explore real-world applications of ML.

2.1 Setup Python Environment for Machine Learning

This tutorial focuses on preparing the Python environment required for machine learning. It begins with an introduction to machine learning and the advantages of using Python, followed by an overview of essential libraries. Learners are guided through installing Miniconda on Ubuntu, creating and activating a dedicated conda environment, and installing all required libraries from a provided package file. The tutorial also covers setting up Jupyter Notebook, adding conda kernels, and introduces the basics of using Jupyter Notebook for implementing machine learning programs.

The tutorial covers the following:

- About Machine Learning.
- Why we use Python for Machine Learning.
- About Python libraries.
- Installation of Miniconda in Ubuntu OS.
- Creating a conda environment for Machine Learning.
- Activating conda environment for Machine Learning.
- Download MLpackage.txt file.
- Install all the libraries available in the txt file in Ubuntu OS.
- Installing Jupyter Notebook in the Machine Learning environment.
- Installing conda kernels in Machine Learning environment.
- About Jupyter Notebook and its basics.

2.2 K Nearest Neighbor Classification

This tutorial introduces the K-Nearest Neighbor (KNN) algorithm for classification tasks. It demonstrates how data points are assigned to the class most common among their nearest neighbors, covering distance metrics, training, and predicting labels for new data.

- Introduction to Nearest Neighbors and K Nearest Neighbor
- Introduction to KNN classification
- Explanation about Iris dataset
- KNN working example using one of the iris features
- Importing the necessary libraries
- Loading the Iris dataset
- Basic Data Exploration and Analysis
- Train and Test Split of dataset
- Choosing the K value using elbow method
- KNN classification model building
- Model prediction and outcome
- Evaluation metrics using classification report

2.3 K Nearest Neighbor Regression

The K-Nearest Neighbor Regression tutorial explains how KNN can be adapted to predict continuous values. It covers finding the average of the nearest neighbors to make predictions, implementing the algorithm in Python, and analyzing the results using performance metrics such as Mean Squared Error.

This tutorial covers the following:

- Introduction to K Nearest Neighbor Regression
- Various distance metrics used in KNN
- Importing the necessary libraries
- Loading the iris dataset
- Standard scaling of the dataset
- Train and Test Split of dataset
- Choosing the K value using elbow method
- KNN regression model building
- Model prediction and outcome
- Evaluation using MSE and Adjusted R squared score

2.4 Linear Regression

This tutorial teaches the concept of linear regression for predicting continuous outcomes. Learners explore the relationship between independent and dependent variables, fit a linear model using Python, visualize regression lines, and evaluate model performance using evaluation metrics

- About Linear Regression
- About Simple Linear Regression
- About Multiple Linear Regression
- About Evaluation Metrics
- Splitting the data into training and testing sets
- Implementing Simple Linear Regression model from scikit-learn

- Importing required Libraries
- Loading the dataset
- Evaluating the model's accuracy
- Implementing Multiple Linear Regression model from scikit-learn
- Evaluating the model's accuracy

2.5 Logistic Regression Binary Classification

The tutorial focuses on logistic regression for binary classification problems. It explains how probabilities are modeled using the sigmoid function, fitting the model in Python, and evaluating its performance.

This tutorial covers the following:

- Introduction to Logistic Regression
- Introduction to Binary classification
- Introduction to Multi class classification
- About Purchase prediction
- Implementing Binary classification
- Model Instantiation of Binary Classification and Model training
- Prediction for Train Data Verification for Binary Classification
- Predictions for Test Data for Binary Classification
- Calculate the ROC-AUC score on the training data
- Calculate the cross entropy loss for the training data

2.6 Logistic Regression MultiClass Classification

This tutorial extends logistic regression to handle multiple classes. Using Python libraries, learners train and test models on multiclass datasets.

- Implementing Multiclass classification
- Model Instantiation of Multiclass Classification and Model training

- Visualize this correlation using a heatmap
- Split the data into training and testing sets
- Build a multiclass classification model
- Prediction for Train Data Verification for Multiclass Classification
- Predictions for Test Data for Multiclass Classification
- Compare the predicted with the actual test class
- Visualize the confusion matrix of the model

2.7 Decision Tree

This tutorial provides a step-by-step introduction to building and evaluating a Decision Tree classifier. Learners also explore how tree visualization provides insights into the decision-making process, making it easier to understand how the model arrives at its outcomes.

This tutorial covers the following:

- Introduction to Decision Tree
- Describing the dataset
- Importing required Libraries
- Loading the dataset
- Encoding Categorical Features
- Splitting the dataset into Training and Testing sets
- Training Decision Tree Classifier
- Evaluating the model's accuracy
- Plotting Confusion matrix
- Visualizing Decision Tree

2.8 Artificial Neural Networks

This tutorial introduces Artificial Neural Networks (ANNs) and their structure through the Multi-Layer Perceptron (MLP). It explains the architecture of neurons and layers, and applies the concepts to the Breast Cancer dataset. The tutorial demonstrates training an MLP classifier, making predictions, and evaluating performance using common metrics.

This tutorial covers the following:

- Introduction to Artificial Neural Networks
- Introduction to Multi-Layer Perceptron
- About ANN Architecture
- Explanation of Neuron Structure
- Importing necessary libraries
- Loading Breast Cancer dataset
- Basic Data Exploration and Analysis
- Train and Test split of dataset
- MLP Classification model building
- Model prediction and outcome
- Evaluation of model's performance

2.9 Support Vector Machine

This tutorial introduces Support Vector Machines (SVMs), covering both Linear and Non-Linear (RBF) approaches. Using the California Housing dataset, it explains preprocessing steps such as encoding and train-test splitting. The tutorial demonstrates building classification models with Linear SVM, followed by Non-Linear SVM, and compares their predictions and evaluation results.

- About Support Vector Machine
- Introduction to Linear SVM
- Introduction to Non-Linear SVM
- Explanation of the California Housing dataset
- Importing necessary libraries
- Loading the dataset
- Label Encoding

- Train and Test Split of dataset
- Linear SVM classification model building
- Model prediction and outcome
- Evaluation for Linear SVM classification
- Non-Linear (RBF) SVM classification model building
- Model prediction and outcome
- Evaluation for Non-Linear (RBF) SVM classification

2.10 K Means Clustering

This tutorial introduces K-Means Clustering and explains its working for grouping similar data points. It covers evaluation using the Silhouette Score and uses a customers dataset to demonstrate the process. It then shows how to instantiate and fit the K-Means model, assign cluster labels, and visualize the resulting clusters to interpret patterns within the data.

- Introduction to K-means Clustering
- Working of K-means Clustering
- Description about Silhouette Score
- Description about the customers dataset
- Importing required Libraries
- Loading the dataset
- Data Exploration
- Finding optimal number of clusters
- Instantiating K-means Clustering model
- Clustering the data
- Visualizing the Clusters for the Data

2.11 Random Forest

This tutorial introduces Ensemble Learning and focuses on the Random Forest algorithm. It explains how Random Forest combines multiple decision trees to improve predictive performance.

- Introduction to Ensemble Learning
- Introduction to Random Forest
- Importing Libraries
- Loading the dataset
- Data Preprocessing
- Train and Test Split
- Model Instantiation of Random Forest and Model training
- Prediction for Train Data Verification for Random Forest
- Predictions for Test Data for Random Forest
- MSE and Adjusted R square score for Random Forest

Code and Script Writing

For the tutorials described in the previous chapter, I prepared the source code along with the corresponding scripts, including visual cues and narration, and recorded the tutorial videos. The scripts were carefully written in accordance with the Spoken Tutorial guidelines. Each segment of the code was explained in detail, and additional references were provided whenever prerequisite concepts could not be fully addressed within the tutorial.

The source code and the corresponding scripts can be found in the following links:

• Setup Python environment for Machine Learning

Code: Sample.ipynb

Script: Tutorial-1 Setup Python Environment for Machine Learning

 $\bullet\,$ K Nearest Neighbor Classification

Code: KNN classification.ipynb

Script: Tutorial-2 KNN Classification Script

• K Nearest Neighbor Regression

Code: KNNregression.ipynb

Script: Tutorial-3 KNNregression

• Linear Regression

Code: LinearRegression.ipynb

Script: Tutorial-4 Linear Regression

• Logistic Regression Binary Classification

Code: LR_Binary.ipynb

Script: Tutorial-5 LR_Binary

• Logistic Regression MultiClass Classification

Code: LR_Multiclass.ipynb Script: Tutorial-6 LR_Multiclass

• Decision Tree

Code: DecisionTree.ipynb Script: Tutorial-7 Decision Tree

• Artificial Neural Networks

Code: Artificial Neural Networks.ipynb

Script: Tutorial-8 Artificial Neural Networks Script

 $\bullet\,$ Support Vector Machine

Code: SVM.ipynb

Script: Tutorial-9 SVM $\,$

• K Means Clustering

Code: K means clustering.ipynb

Script: Tutorial-10 K-means Clustering Script

• Random Forest

Code: RandomForest.ipynb

Script: Tutorial-11 Random Forest

Creating slides

Following the development of the code and script, I also prepared slides for each tutorial in the series. These slides typically begin with the introduction, learning objectives, system requirements, and prerequisites, ensuring that learners know what to expect and what background knowledge is required. For each tutorial, the slides provide brief theoretical explanations of the topic along with references to the datasets or files used. The concluding slides present a summary of the concepts covered, followed by assignments or practice problems to reinforce learning. Additional slides also include information about the Spoken Tutorial project, forums for learner support, and acknowledgments.

The slides for the tutorials can be found at the following links:

- Setup Python Environment for Machine Learning
- K Nearest Neighbor Classification
- K Nearest Neighbor Regression
- Linear Regression
- Logistic Regression Binary Classification
- Logistic Regression MultiClass Classification
- Decision Tree
- Artificial Neural Networks
- Support Vector Machine
- K Means Clustering
- Random Forest

Phases of my Spoken Tutorial

5.1 Outline

The outline defines the structure of the Spoken Tutorial, listing all the topics to be covered in a sequential manner. It serves as a roadmap, ensuring that the content maintains a logical flow, beginning with basic concepts and gradually progressing to advanced topics.

5.2 Script

The script contains detailed, step-by-step instructions, explanations, and examples that form the core of the tutorial. It is written in a simple and precise manner, enabling learners, especially beginners, to follow the process without difficulty. Key points for demonstration are also included to maintain clarity. The script follows the Spoken Tutorial guidelines and is written in a clear and concise manner for easy understanding.

5.3 Slides

Slides provide supporting visual content for the tutorial. They include an introduction, prerequisites, system requirements, and key concepts, making the learning process more structured and engaging. The slides are created using LaTeX to ensure a neat layout and professional presentation. The color scheme and formatting remain consistent across all tutorials, in line with the Spoken Tutorial guidelines.

5.4 Novice Check

At this stage, the script and slides undergo a novice check to ensure that the material can be easily understood by learners with no prior experience in the subject. Feedback is used to simplify explanations, add missing details, and eliminate ambiguities. The Novice check for my tutorials were performed by Ms. Rashmi Patankar from the Spoken Tutorial team. She performed the check multiple times to ensure that the content is well explained and that the source code is bug-free.

5.5 Domain Check

The domain check is conducted by subject experts to verify technical accuracy and correctness of the explanations. The domain check was performed by Dr. Subbulakshmi T from VIT Chennai and Ms. Nirmala Venkat from IIT Bombay.

5.6 Admin Check

The admin check ensures compliance with Spoken Tutorial guidelines. This includes reviewing the format, structure, consistency of language, and slide presentation standards. Approval at this stage confirms that the material is ready for recording. The admin check was done by Ms. Nirmala from the Spoken Tutorial team

5.7 Recording

In the final phase, the tutorial is recorded using the approved script and slides. The recording combines screen demonstrations with clear narration, ensuring that learners can follow along smoothly. The recording is produced in high quality, adhering strictly to the approved script and maintaining focus on the relevant content without deviations.

Conclusion

The FOSSEE fellowship has been an enriching experience that allowed me to contribute to the Spoken Tutorial learning platform while strengthening my programming skills. Beyond technical growth, it also helped me develop key professional skills such as project management and teamwork.

Through the process of script writing, I learned the value of writing clean code and maintaining thorough documentation to ensure a seamless learning experience for the users. Additionally, while preparing the tutorial slides, I gained hands-on exposure to the versatile formatting features of LaTeX.

I am deeply grateful to the Spoken Tutorial team and my mentors at VIT Chennai for providing me with this valuable opportunity and their constant support throughout the fellowship.

References

- $\bullet \ \, \rm https://spoken-tutorial.org/$
- $\bullet \ \, \rm https://scikit-learn.org/stable/$