

Introduction to Machine Learning: Linear and Logistic Regression and Neural Networks using Python

1.1 What is Machine Learning?

Identifying Cats and Dogs



Pictures source: pixabay.com

How do we as humans learn to identify cats and dogs?

- At an early age, parents/siblings/teachers point to the animals or to pictures and say that this is a cat or this is a dog
- As this happens 10s or 100s of times during our early years, a model/rules to identify cats and dogs are encoded in our brains
- Using the model/rules, children become able to classify whether an animal is a cat or a dog without help from others
- This is referred to as Concept Learning in human psychology

Machine Learning – Definition

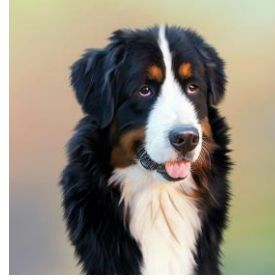
- Machine learning (ML) is the scientific study of algorithms and statistical models that computer systems use to perform a specific task without using explicit instructions, relying on patterns and inference instead
- Machine learning algorithms build a mathematical model based on sample data, known as "training data", in order to make predictions or decisions without being explicitly programmed to perform the task
- Source: https://en.wikipedia.org/wiki/Machine_learning
- How does relate to identifying cats and dogs?
 - The examples we are shown in early childhood constitute the training data
 - The information encoded in our brains is the “model” based on this data
 - In later years, we are able to perform the identification using the model without using any further explicit instructions

Machine Learning:

Three areas of application:

- Classification
- Clustering
- Game playing

Machine Learning: Classification



Pictures source: pixabay.com

For classification (cats and dogs) using machine learning:

- Training Data consists of 100s or 1000s of images each labelled as cat or dog
- Based on training data, Machine Learning system learns a model/rules to perform classification
- When new image (without label) is shown, system uses the model to make a prediction of cat or dog

Machine Learning: Clustering



Pictures source: pixabay.com

For clustering using machine learning:

- Input consists of images (cars and bicycles) without labels
- Machine Learning system is used to find patterns/“clusters” in the data
- In this case, expect 2 broad clusters (cars and bicycles)
- If sufficient data, system may also find sub-clusters (convertibles, 2-doors, 4-doors, road bikes, mountain bikes, blue objects, etc)

Machine Learning: Game Playing



By Goban1 - Own work, Public Domain,
<https://commons.wikimedia.org/w/index.php?curid=15223468>

For game playing (such as Go above), Machine Learning system:

- Is provided (encoded) with the rules of game (how to move pieces)
- Through self-play or provided with database of games, Machine Learning system learns what moves to make at each turn to ultimately win game
- AlphaZero program that plays Go, Chess, and Shogi was trained solely using self-play in a short period of time and was able to beat other systems convincingly
 - See <https://en.wikipedia.org/wiki/AlphaZero> for more details

Types of Machine Learning

Type	Description (Adapted from https://en.wikipedia.org/wiki/Machine_learning)
Supervised Learning	<p>Process of learning a function that maps input information to labelled output information. The input/output information is called the training data. The learned function is then used to predict output labels when new input information is provided.</p> <p>Applications: regression, image classification, language translation, spam classification, auto-completion, etc</p>
Unsupervised Learning	<p>Process of learning previously undetected patterns in a data set with no pre-existing labels and with a minimum of human supervision</p> <p>Applications: clustering, data mining</p>
Reinforcement Learning	<p>Process of learning what actions to take given the situation so as to maximize a reward</p> <p>Applications: Game playing (Atari, Chinese checkers, Chess, Go, etc), Industrial Control, etc</p>

1.2 Supervised Machine Learning: Linear Regression, Logistic Regression, and Neural Networks

Supervised Machine Learning

- Process of learning a function that maps input information to labelled output information. The input/output information is called the training data. The learned function is then used to predict output labels when new input information is provided.

Supervised Machine Learning: Applications

Application	Input Information	Output Information/Label
Regression	<ul style="list-style-type: none">House features: lot area, floor area, # of bathrooms, # of floors, size of garage, etc	<ul style="list-style-type: none">Price
Image Classification	<ul style="list-style-type: none">Images of cats and dogsBone x-raysImages of animals	<ul style="list-style-type: none">cat/dognormal/brokenAnimal name
Language Translation	<ul style="list-style-type: none">English words and phrases	<ul style="list-style-type: none">French translations
Spam Filter	<ul style="list-style-type: none">Messages	<ul style="list-style-type: none">spam/not spam
Auto-Completion	<ul style="list-style-type: none">(2,3,4,5 ...) word phrases	<ul style="list-style-type: none">Next word

Supervised Machine Learning: Three Approaches

This course presents three approaches for Supervised Learning:

- Linear Regression
 - Used for prediction of real values
 - Simple approach, which you have probably seen in your courses
 - Useful for introducing concepts
- Logistic Regression
 - Used for binary classification
 - Mathematics and coding are extensions of those for Linear Regression
- Neural Networks (Multi-Layer Perceptrons)
 - Used for binary and multi-class classification (or regression)
 - Mathematics and coding build on concepts for Linear and Logistic Regression

Linear Regression – Line Fitting

Training Data:

- Input information: X values
- Output information: Y values

Linear Regression Goal:

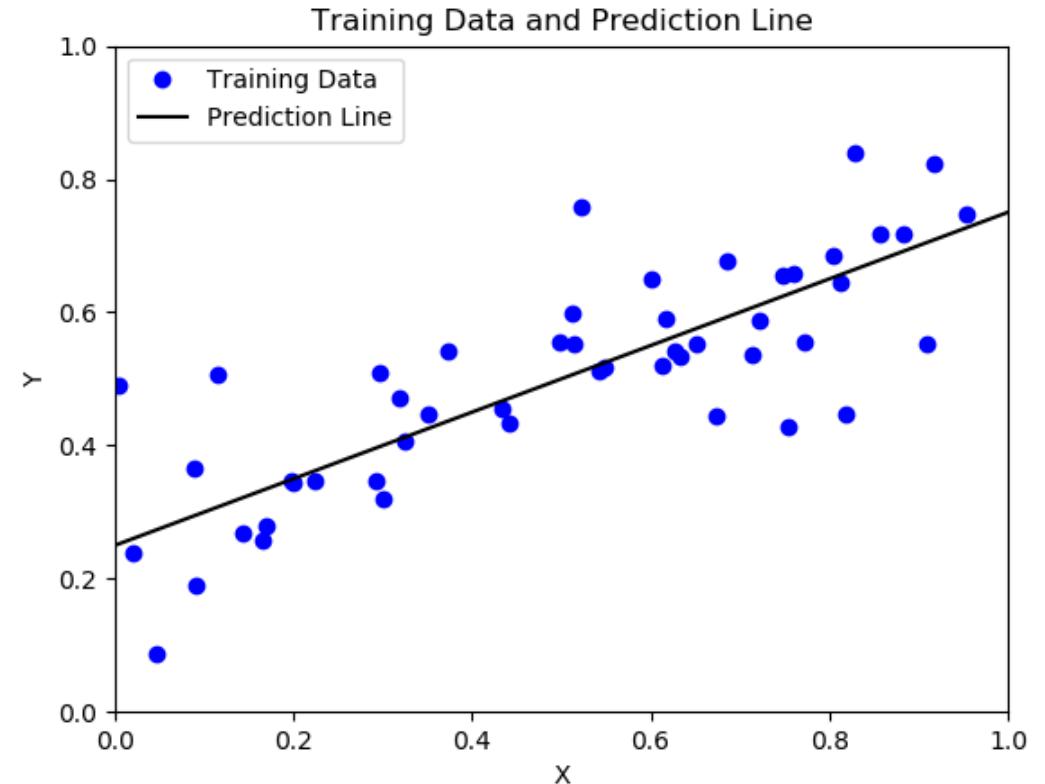
- Find straight line that best fits the training data

Prediction:

- Use line to predict Y values given new input X values

Why start with Linear Regression?

- Simple problem with well known solution
- This course will present a general approach that can also be applied to Logistic Regression and Neural Networks



Logistic Regression – Binary Classification

Training Data:

- Input Information: points in (x_0, x_1) plane
- Output Information: label 0 (red) or 1 (blue) for each point

Logistic Regression Goal:

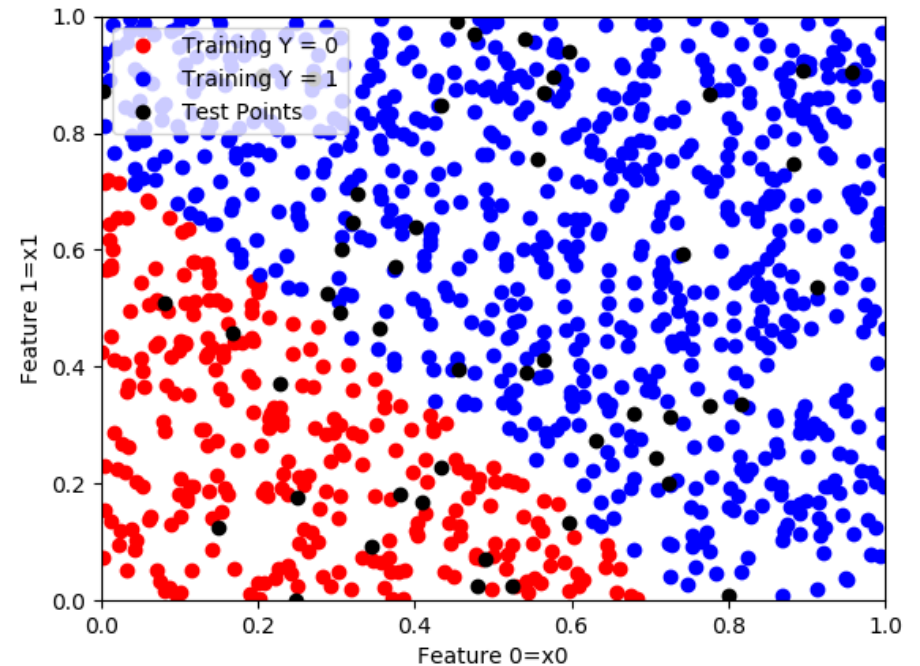
- Find function that best fits 0 and 1 labels in training data

Prediction:

- Using function, determine labels for new input test points (black points in picture)

Logistic Regression:

- Approach builds on that for Linear Regression
- Limited in variety of functions it can represent: not suitable when boundary between 0 and 1 regions is not straight line



Neural Networks – Binary Classification

Training Data:

- Input Information: points in (x_0, x_1) plane
- Output Information: label 0 (red) or 1 (blue)

Neural Networks Goal:

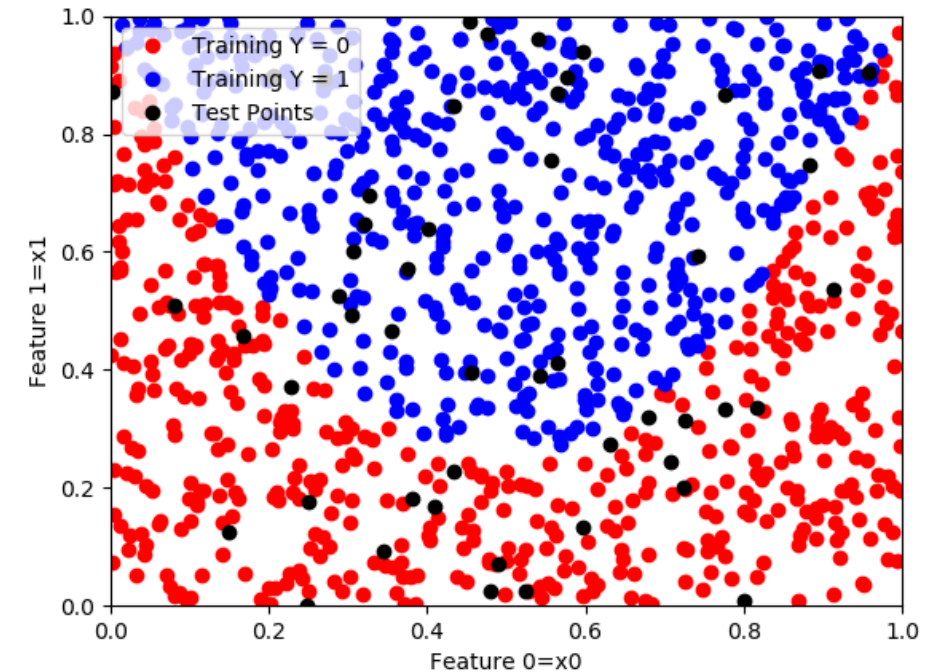
- Find function that best fits 0 and 1 labels in training data

Prediction:

- Using function, determine labels for new input test points (black points in picture)

Neural Network

- Can be used in case of more complex boundary between 0 and 1 regions
- Can be used for classification with more than 2 classes



Neural Networks – Multi-class Classification

Training Data:

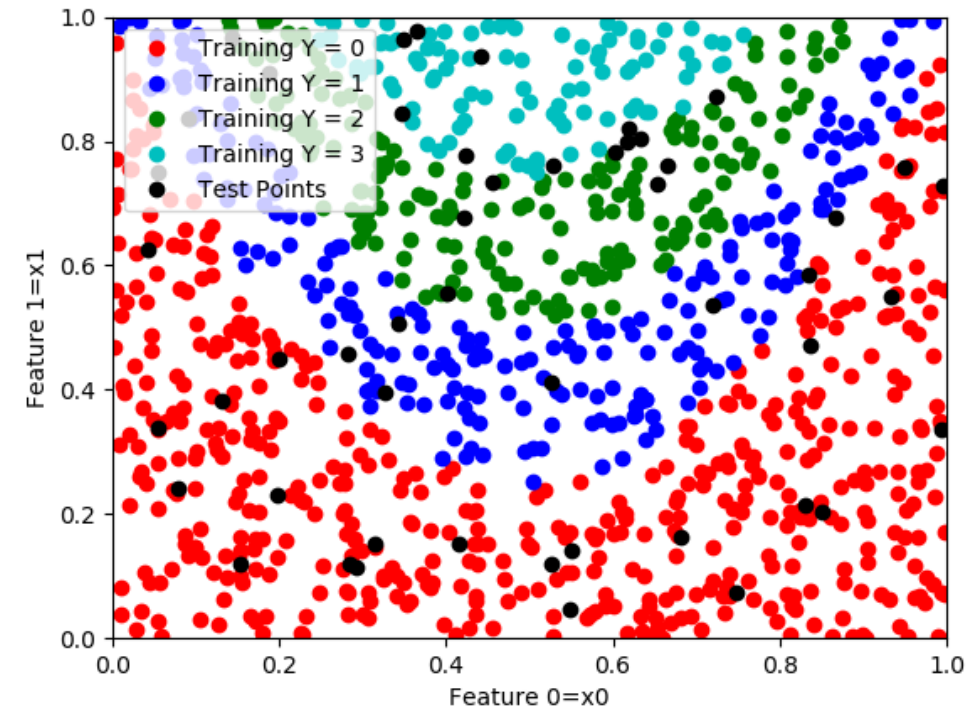
- Input Information: points in (x_0, x_1) plane
- Output Information: label 0 (red) or 1 (blue), 2 (green), 3 (cyan)

Neural Networks Goal:

- Find function that best fits 0,1,2,3 labels in training data

Prediction:

- Using function, determine labels for new input test points (black points in picture)



Approach for Supervised Machine Learning

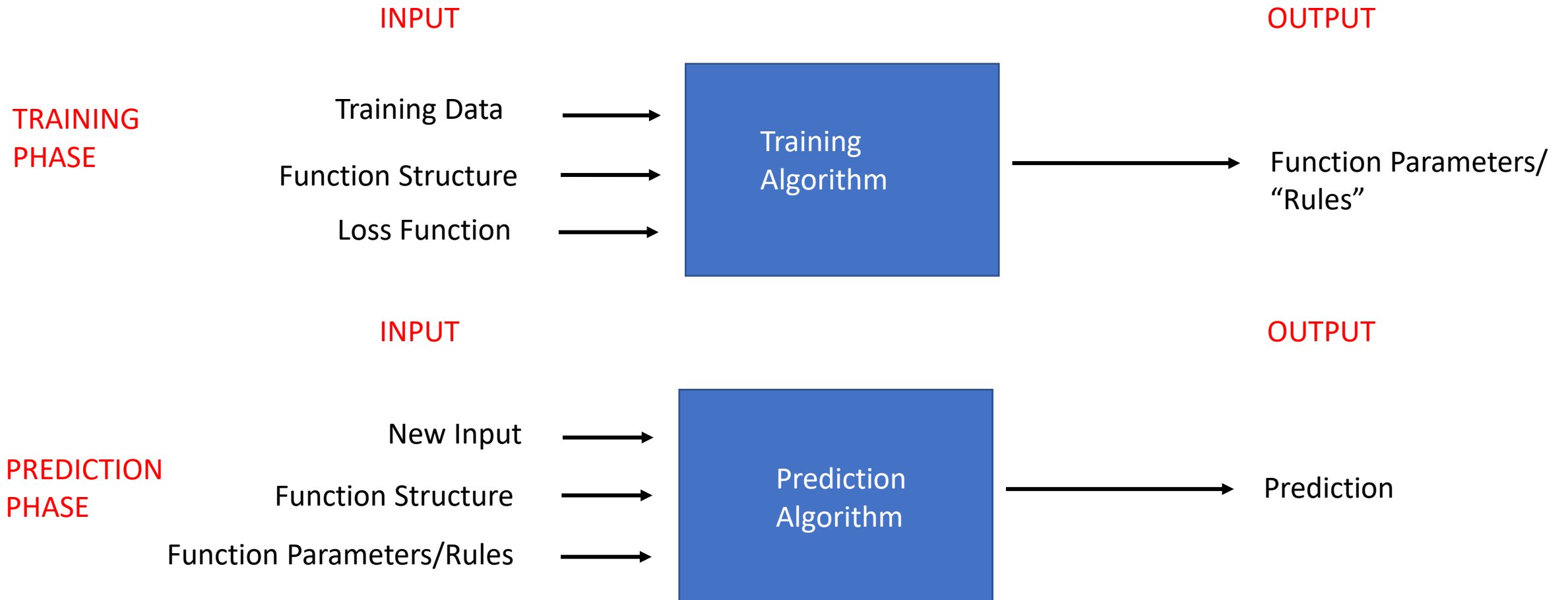
General approach for Supervised Machine Learning:

- Training Data: Start with input/output information
- Function Structure: assume a functional form with unknown function parameters to map training input info to output info
- Define Loss Function: to measure how well function maps input info to output info
- Training Algorithm: find function parameters to minimize Loss function for training data
- Prediction Algorithm: use function structure and learned parameters to compute output info when new input info is provided

For Linear Regression (in 1 dimension):

- Training Data: points in X-Y plane
- Function Structure: assume $Y = W \cdot X + b$ (slope W , intercept b are “Function Parameters”/Rules)
- Loss function: squared error over training data
- Training Algorithm: find W and b to make Loss function as small as possible for training data
- Prediction Algorithm: use function structure and “learned” W , b to compute output Y when new input X values are provided

Supervised Machine Learning: Training and Prediction



What this Course Covers – Math and Algorithms

Course covers mathematics and algorithms for:

Linear Regression, Logistic Regression and Neural Networks

Mathematics and Algorithms:

- Format of the training data
- Function structures
- Loss functions
- Training algorithms for determining function parameters
- Prediction algorithms
- Techniques for improving performance of machine learning systems

What this Course Covers – ML Framework

- A machine learning framework consists of building blocks for designing, training and validating deep neural networks, through a high level programming interface
- This course presents a code walkthrough/development of a basic machine learning framework (in Python) for setting up and training Linear Regression, Logistic Regression, and Neural Networks using an object oriented approach
- Many examples are presented to show how to use the framework

What this Course Covers – Case Studies

Course shows how to use framework to address 3 case studies:

- House Price Prediction (Regression)
 - Use Linear Regression to predict house prices
- Spam Classification (Binary Classification)
 - Use a Neural Network to classify text messages as spam or not spam
- Digits Classification (Multi-class Classification)
 - Use a Neural Network to identify digits (0,1,...,9) from images

1.3 About this Course

About this Course

Students completing this course will:

1. Gain an understanding of the mathematical foundations and algorithms for Linear Regression, Logistic Regression, and Neural Networks approaches for supervised machine learning
2. Gain an understanding of computer codes (written in Python) for Linear Regression, Logistic Regression, and Neural Networks. The course will walk through development of a basic machine learning framework.
3. Gain an understanding how to measure and improve performance of supervised machine learning systems
4. Gain a foundation for further study/practice of machine learning

Course Prerequisites

All courses have a starting point – an assumption about the base of students' knowledge. Courses cannot assume no knowledge and teach all the prerequisites. Students who don't have the prerequisite knowledge/skills listed below must acquire them before (or while) taking this course.

Prerequisites:

- Linear Algebra
 - Students should be familiar with vectors, matrices, transpose, matrix multiplication, least squares
- Multivariable Calculus
 - Students should be familiar with computing partial derivatives and employing the chain rule
- Python Programming
 - Students should be able to write and run Python 3 programs
 - It is preferable, but not required, that students be familiar with the NumPy package for scientific computing. Key features and functions will be explained in the course

Audience for this Course

This course is suitable for:

1. Students without any previous experience with machine learning
2. Students who have some knowledge of the subject and would like a refresher and/or gain a more detailed understanding of the mathematical foundations, algorithms, and development of a machine learning framework in Python

Course Approach

1. The underlying mathematics is explained in detail
 - The course will provide motivation and derivations for all the underlying mathematics
 - Examples provided in the presentations are complemented by Jupyter notebook demos
2. Algorithms are described in detail
 - The course will show how the underlying mathematics connects to the algorithms used in machine learning
3. Code development walkthrough is performed
 - This course walk through development of a basic machine learning framework that can be used for Linear Regression, Logistic Regression, and Neural Networks
 - Walkthrough will emphasize how algorithms are translated into Python code and how a machine learning framework can be created using an object oriented approach

How to Get the Most from this Course

One cannot become proficient in a subject by simply watching someone else!

How to make the most from this course:

1. Take notes as you go through the material and work out the mathematical derivations by yourself
2. Ask questions on the forums
3. Do the programming
 - Programming videos will start with objectives and high-level summary of the code structure, followed by video of code implementation
 - You may want to just review the objectives and then design and implement the code by yourself.
 - Alternatively, review objectives and high-level summaries, and then code for yourself. Can always check videos of code walkthroughs
 - All codes are made available

Why Code from Scratch?

- A machine learning framework consists of building blocks for designing, training and validating deep neural networks, through a high level programming interface.
- Following link has a listing of more than a dozen machine learning frameworks:

https://en.wikipedia.org/wiki/Comparison_of_deep-learning_software

Examples: Tensorflow, Keras, Pytorch, Caffe, Theano, etc

- Ultimately, going forward for machine learning projects, research, and production development, you should probably use one of these frameworks, that has been tested and optimized. Some of these frameworks have been also developed for GPU (which are typically faster than CPU).
- Why spend the time learning the math and algorithms and coding from scratch when one can use these frameworks “off the shelf”?
- My fundamental belief (and a reason for creating this course) is that to truly understand what is going on in a frameworks one must work through the mathematics and algorithms and write codes from scratch

About the Instructor

- Holds PhD in Applied Mathematics
- Worked for 10 years in university settings as post-doc and professor conducting research in applied math and numerical analysis and teaching undergraduate- and graduate-level courses
- Worked for 17+ years in financial risk management at a software start-up, financial information services company, and a large international bank

Course Outline:

- Chapter 1: Introduction
- Chapter 2: Python Demos
- Chapter 3: Mathematical Concepts
- Chapter 4: Linear and Logistic Regression
- Chapter 5: Neural Networks
- Chapter 6: Optimization, Validation, Regularization, and Performance
- Chapter 7: Case Studies
- Chapter 8: Introduction to Tensorflow
- Chapter 9: Summary and Thank You

1.4 Resources and Set Up

Resources and Set Up

- Instructor will use Windows 10 machine
 - Course material not specific to Windows – can use MacOS or Linux
- All code examples written in Python
- Course will run programs using
 - Jupyter Notebook
 - For demos
 - Command Window
 - For the machine learning framework
- Should have text editor compatible with Python for writing and editing programs
 - Examples: Atom, Sublime, Notepad+, etc
 - Instructor will use Sublime, but you can use your favourite editor

Course Material

- Course material located at:

<https://github.com/satishchandrareddy/IntroML/>

- Folder: Code
 - Contains course machine learning framework and drivers
 - Series of folders(versions) mirror the lectures. Idea is to discuss underlying mathematics and algorithms and then build the basic machine learning framework in manageable chunks, as opposed to presenting all the theory and then do the coding.
- Folder: Examples
 - Jupyter notebooks show python code examples that complement the examples presented in the lectures
- Folder: Presentations
 - PDF files of presentations

Summary of Course Code Versions

Version	Description
Version 1.1	Initial functionality for Linear Regression and Derivative Testing
Version 1.2	Add Training and Prediction functionality for Linear Regression
Version 1.3	Add functionality for Logistic Regression
Version 2.1	Add functionality for Neural Networks for binary classification
Version 2.2	Add functionality for Neural Networks for multi-class classification
Version 3.1	Add functionality for Momentum, RmsProp, and Adam optimizers and mini-batch optimization
Version 3.2	Add functionality for validation and additional accuracy calculations
Version 3.3	Add functionality for hyperparameter searches and regularization
Version 4.1	Add functionality for house price, Spam, and MNIST case studies
Version 5.1	Add Tensorflow examples

Packages used in Course

Component	Version	Description/Comments
Python	3.7.1	It is assumed that students have Python on their machine and have ability to run Jupyter notebooks. If you don't have Python on your machine, it is probably best to install Python using the Anaconda package, which aims to simplify package management. See: https://www.anaconda.com/ for download and installation of Anaconda. See: https://www.python.org for information about Python. I have found that Python 3.7.1 is consistent with the following package versions.
NumPy	1.18.1	Package for scientific computing. See https://numpy.org/ for details.
Matplotlib	3.0.3	Package for plotting. See https://matplotlib.org/ for details
pandas	0.23.4	Package containing data structures and data analysis tools – will use to load data from csv file. See https://pandas.pydata.org/
scikit-learn (sklearn)	0.20.1	Package for machine learning. Will use its text processing functions for spam classification. See https://scikit-learn.org/stable/index.html
tensorflow	2.1.0	Open source platform/framework for machine learning. See http://tensorflow.org
copy, csv, pathlib, time & unittest		These packages are part of the python release

Virtual Environment

- Should set up a “Virtual Environment” for this course
- Virtual environment will allow user to install specific versions of packages within environment without conflicting with different versions of packages used outside of that environment
 - Example: course uses Numpy version 1.18.1
 - You may want to use another version of Numpy for other purposes
 - Virtual environments will allow you to use both versions on your machine without conflicts

Package Installation Using Requirements File

- Can install all packages using requirements file
- Course requirements_IntroML.txt file located in folder IntroML/Code

- Contents of file:

numpy==1.18.1

matplotlib==3.0.3

pandas==0.23.4

scikit-learn==0.20.1

tensorflow==2.1.0

Note: all supporting packages will also be installed!

Virtual Environment and Package Installation

Task	Commands (Issue in Anaconda Prompt window)
Create virtual environment using python 3.7.1 (NAME is user specified name of environment)	conda create -n NAME python=3.7.1
List all virtual environments	conda env list
Activate virtual environment (NAME is user specified name of environment)	conda activate NAME
Deactivate current virtual environment	conda deactivate
Installing packages (using requirements_IntroML.txt) - Make sure desired environment is activated (use double dash - - before "file")	conda install --file requirements_IntroML.txt OR if there is an issue with the above, try pip install -r requirements_IntroML.txt
List packages in environment - Make sure desired environment is activated	conda list
Connect virtual environment to Jupyter Notebook - Make sure desired environment is activated (NAME is user specified name of environment) (Use double dash - - before "user" and "name")	pip install --user ipykernel python -m ipykernel install --user --name=NAME

References

- Following page has useful information about conda commands
<https://docs.conda.io/projects/conda/en/latest/commands.html>
- For details on setting up the virtual environment and connecting it to the Jupyter notebook: <https://janakiev.com/blog/jupyter-virtual-envs/>
- For details on installing packages in the virtual environment:
<https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-pkgs.html>

1.4 Resources and Set Up DEMO

- Download resources from Github:

<https://github.com/satishchandrareddy/IntroML/>

- Environment Set Up

- Create environment for the course
- Install relevant packages
- Connect environment to Jupyter notebook
- Test environment set up on Jupyter notebook