

EE2211 Tutorial 1

Question 1:

What is the difference between ML (Machine Learning) and AI (Artificial Intelligence)?

Question 2:

Which of the following is the most reasonable definition of machine learning?

- (a) Machine learning is the field of allowing robots to act intelligently.
- (b) Machine learning is the science of programming computers.
- (c) Machine learning only learn from unlabeled data.
- (d) Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed.

Question 3:

A computer program is said to *learn* from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E . Suppose we feed a learning algorithm a lot of historical weather data, and have it learn to predict weather. In this setting what is T ?

- (a) The historical weather data.
- (b) The probability of it correctly predicting a future data's weather.
- (c) The weather prediction task.
- (d) None of these.

Question 4:

Suppose you are working on weather prediction and use a learning algorithm to predict tomorrow's temperature (in degrees Centigrade/Fahrenheit).

(i) Would you treat this as a classification or a regression problem?

- (a) Regression.
- (b) Classification.
- (c) Clustering.
- (d) None of these.

(ii) What kind of data should you gather? **temp and date**

Question 5:

You want to develop learning algorithms to address each of the following two problems.

P1: You'd like the software to examine your email accounts, and decide whether each email is a spam or not.

P2: You have a large quantity of green tea (e.g., 1000kg) with a record of previous sales. You want to predict how much of it will sell over the next 6 months.

Should you treat these as classification or as regression problems?

- (a) Treat both P1, P2 \rightarrow regression problems.
- (b) Treat both P1, P2 \rightarrow classification problems.
- (c) Treat P1 \rightarrow regression problem, P2 \rightarrow classification problem.
- (d) Treat P1 \rightarrow classification problem, P2 \rightarrow regression problem.

Question 6:

Suppose you are working on stock market prediction. Typically tens of millions of shares of a company's stock are traded each day. You would like to predict the number of shares that will be traded tomorrow.

(i) Would you treat this as a classification or a regression problem?

- (a) Regression.
- (b) Classification.
- (c) Clustering.
- (d) None of these.

(ii) If the data you have collected involved millions of attributes, what would you do?

categorised then apply regression separately

dimensional reduction, reduce size of data by removing unnesscasiy data, speed up learning and prediction

if need exact answer like yes/no, supervised
if discovering category, unsupervised

Question 7:

Some of the problems below are best addressed using a supervised learning algorithm, and the others with an unsupervised learning algorithm. Which of the following would you apply supervised learning to? (Select all that apply) Assume some appropriate dataset is available for your algorithm to learn from.

- (a) Determine whether there are vocals (i.e., a human voice singing) in each audio clip extracted from a piece of music, or it is a clip of only musical instruments and no vocals.
- (b) Given data on how 1000 medical patients respond to an experimental drug (such as effectiveness of the treatment, side effects, etc.), discover whether there are different categories or “types” of patients in terms of how they respond to the drug, and if so what these categories are.
- (c) Given a large dataset of medical records of patients suffering from heart disease, try to learn whether there might be different clusters of such patients for which we might tailor separate treatments.
- (d) Given a set of data which contains the diet and the occurrence of diabetes from a population over a 10-year period. Predict the odds of a person developing diabetes over the next 10 years.

T = predict if a patient is COVID-19 infected

P = prediction accuracy

E = labelled patients and their age and health conditions, symptomatic data

if odd of every single year, regression
if yes no get diabetes, classification

Question 8:

Suppose you are working on a machine learning algorithm to predict if a patient is COVID-19 infected according to the patient’s particulars such as age and health conditions, symptomatic data, such as fever, dry cough, tiredness, aches and pains, sore throat, diarrhoea, conjunctivitis, and headache etc. What are the Task, Performance, and Experience involved according to the definition of machine learning?

Question 9:

We use labelled data for supervised learning, where the labels are used as the desired target of prediction for classifiers. Which of the next data are the useful labelled data?

- (a) To build an image object classifier to discriminate between apple and orange, we have many fruit images labelled with the country of origin.
- (b) To build a system to predict the number of COVID cases for tomorrow given the past daily record, we have a collection of daily data for a period of 12 months.
- (c) To build a classifier to automatically evaluate student essays, we have collected a set of student essays that have not been graded by teachers.

Question 10:

Determine whether each of the following is “inductive” or “deductive” reasoning?

- (a) The first coin I pulled from the bag is a penny. The second and the third coins from the bag are also pennies. Therefore, all the coins in the bag are pennies. inductive guessing population from sample, induce to big, inductive
- (b) All men are mortal. Harold is a man. Therefore, Harold is mortal. deductive

Question 11:

Find a problem of your interest and formulate it as a machine learning problem. List out the input features and output response and provide your choice regarding the types of learning [such as supervised or unsupervised learning.

T: identify bucket from lidar map

supervised

P: classification accuracy

E: lidar map with bucket bounding boxes

From Tutorial 2 onwards, we shall use Python for some computation. Here are some Python Resources:

Installing scikit-learn (Ref: [Book2] Andreas C. Muller and Sarah Guido, “Introduction to Machine Learning with Python: A Guide for Data Scientists”, O’Reilly Media, Inc., 2017)

scikit-learn depends on two other Python packages, NumPy and SciPy. For plotting and interactive development, you should also install matplotlib, IPython, and the Jupyter Notebook. We recommend using the following prepackaged Python distribution, which provides the necessary packages:

Anaconda

A Python distribution made for large-scale data processing, predictive analytics, and scientific computing. Anaconda comes with NumPy, SciPy, matplotlib, pandas, IPython, Jupyter Notebook, and scikit-learn. Available on Mac OS, Windows, and Linux, it is a very convenient solution and is the one we suggest for people without an existing installation of the scientific Python packages. Anaconda now also includes the commercial Intel MKL library for free. Using MKL (which is done automatically when Anaconda is installed) can give significant speed improvements for many algorithms in scikit-learn.

Some tutorials that might be useful:

A quickstart tutorial on NumPy: <https://numpy.org/devdocs/user/quickstart.html>

Some community tutorials on Pandas: https://pandas.pydata.org/pandas-docs/stable/getting_started/tutorials.html

Scikit-learn tutorials: <https://scikit-learn.org/stable/tutorial/index.html>

Python Numpy Tutorial (with Jupyter and Colab):

<https://cs231n.github.io/python-numpy-tutorial/#jupyter-and-colab-notebooks>