

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. $\rightarrow E$
- (b) The probability of it correctly predicting a future data's weather.
- ✓(c) The weather prediction task. $\rightarrow T$
- (d) None of these.

$\frac{P}{\rightarrow}$ prediction error (predicted vs actual temp)
 \rightarrow accuracy between predicted / actual

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. \rightarrow output continuous
- (b) Classification. \rightarrow output categorical
- (c) Clustering.
- (d) None of these.

(ii) What kind of data should you gather?

\rightarrow current weather, temperature, cloud cover, humidity

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. \rightarrow classify if spam or not (classification)

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. \rightarrow regression, sales over time

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.

\downarrow

curve fitting

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? \rightarrow supervision

- ✓(a) Regression. \rightarrow predict num of shares
- (b) Classification.
- (c) Clustering. \rightarrow unsupervised
- (d) None of these.

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

\rightarrow take only the relevant features (useful)

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.

no label
group based
on similarity
of features
→ clustering

- (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. (Supervised) → give data of human voice singing
- (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. → clustering, find out types of data (unsupervised)
- (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. → clustering (unsupervised)
- (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. (Supervised)
data → output (labels)

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?

→ Task: Classify patients into “infected” and “non-infected”
Performance: Accuracy of classification
Experience: patient’s symptoms and actual diagnosis

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. → not good, maybe shape, colour
- ✓ (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.

→ essay should be graded

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
- (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.

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>

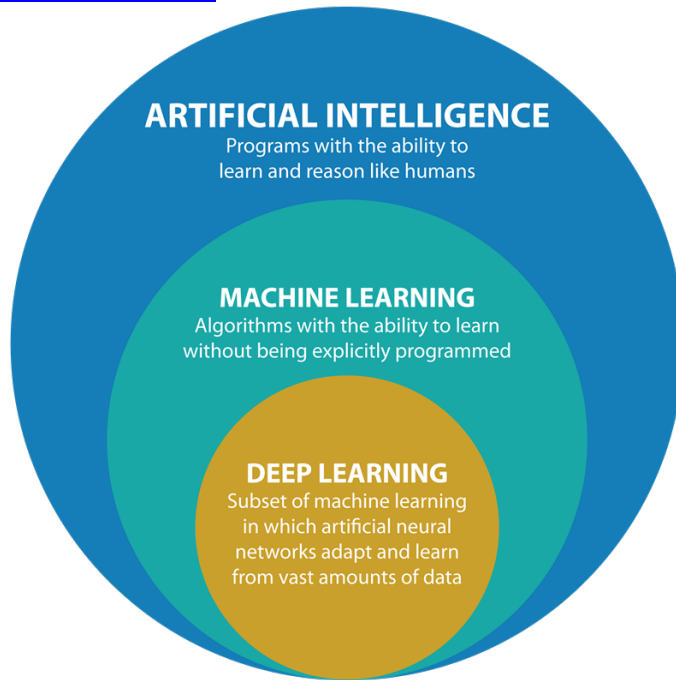
EE2211 Tutorial 1

Question 1:

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

Suggested discussion: Artificial Intelligence is the broader concept of machines being able to carry out tasks in a way that we would consider "smart". And, Machine Learning is a current application of AI based around the idea that we should really just be able to give machines access to data and let them learn for themselves.

<https://www.forbes.com/sites/bernardmarr/2016/12/06/what-is-the-difference-between-artificial-intelligence-and-machine-learning/#741adc6b2742>



Ref: <https://www.vsinghbisen.com/technology/ai/difference-between-artificial-intelligence-and-machine-learning/>

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.

Ans: (d)

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.

Ans: (c)

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?

Ans: (i)(a)

(ii) Weather forecasts are made by collecting quantitative [data](#) (e.g., changes in [barometric pressure](#), current weather conditions, and sky condition or [cloud](#) cover) about the current state of the atmosphere at a given place and using [meteorology](#) to project how the atmosphere will change.

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 → regression problems.
- (b) Treat both P1, P2 → classification problems.
- (c) Treat P1 → regression problem, P2 → classification problem.
- (d) **Treat P1 → classification problem, P2 → regression problem.**

Ans: (d)

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?

Ans: (i)(a), (ii)(**extract relevant features**)

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.**

Ans: (a), (d)

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 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?

Ans: (please refer to the definition of Task, Performance, and Experience in the lecture notes)

Task: patient classification into 'infected' or 'uninfected'

Performance: accuracy of classification

Experience: patient's symptomatic data with actual diagnosis

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.

Ans:

- (a) The useful fruit images should be labelled with apple or orange. Country of origin doesn't tell apple or orange. Therefore, the data is not useful
- (b) We can use n days of historical data as the input, and n+1th day's data as the target. This dataset is useful;
- (c) The useful dataset should include student essays and the grades. Student essays are the input, and the grades are the desired target of prediction. This dataset is not useful.

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.
- (b) All men are mortal. Harold is a man. Therefore, Harold is mortal.

Ans: (a) inductive, (b) 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.

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 distributions, 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>