Introduction to ML & DL

Shan-Hung Wu shwu@cs.nthu.edu.tw

Department of Computer Science, National Tsing Hua University, Taiwan

Machine Learning

Shan-Hung Wu (CS, NTHU)

Introduction

Outline



- **1** What's Machine Learning?
- 2 What's Deep Learning?
- 3 About this Course...



Outline



2 What's Deep Learning?

3 About this Course...



Shan-Hung Wu (CS, NTHU)

• To solve a problem, we need an algorithm

• E.g., sorting

- To solve a problem, we need an algorithm
 - E.g., sorting
 - A priori knowledge is enough

- To solve a problem, we need an algorithm
 - E.g., sorting
 - A priori knowledge is enough
- For some problem, however, we do not have the a priori knowledge
 - E.g., to tell if an email is spam or not
 - The correct answer varies in time and from person to person

• To solve a problem, we need an algorithm

- E.g., sorting
- A priori knowledge is enough

• For some problem, however, we do not have the a priori knowledge

- E.g., to tell if an email is spam or not
- The correct answer varies in time and from person to person
- Machine learning algorithms use the *a posteriori knowledge* to solve problems

• To solve a problem, we need an algorithm

- E.g., sorting
- A priori knowledge is enough
- For some problem, however, we do not have the a priori knowledge
 - E.g., to tell if an email is spam or not
 - The correct answer varies in time and from person to person
- Machine learning algorithms use the *a posteriori knowledge* to solve problems
 - Learnt from *examples* (as extra input)

Example Data ${\mathbb X}$ as Extra Input

• Unsupervised:

$$\mathbb{X} = \{ \pmb{x}^{(i)} \}_{i=1}^{N}, \text{ where } \pmb{x}^{(i)} \in \mathbb{R}^{D}$$

• E.g., $\mathbf{x}^{(i)}$ an email

Example Data ${\mathbb X}$ as Extra Input

• Unsupervised:

$$\mathbb{X} = \{ \pmb{x}^{(i)} \}_{i=1}^N, \text{ where } \pmb{x}^{(i)} \in \mathbb{R}^D$$

• E.g., $\boldsymbol{x}^{(i)}$ an email

• Supervised:

$$\mathbb{X} = \{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=1}^{N}, \text{ where } \boldsymbol{x}^{(i)} \in \mathbb{R}^{D} \text{ and } \boldsymbol{y}^{(i)} \in \mathbb{R}^{K},$$

• E.g.,
$$y^{(i)} \in \{0,1\}$$
 a spam label

Shan-Hung Wu (CS, NTHU)

General Types of Learning (1/2)

• Supervised learning: learn to predict the labels of future data points

$$X \in \mathbb{R}^{N \times D}$$
:
 6 6 7 4 4 5
 $\mathbf{y} \in \mathbb{R}^{N \times K}$:
 $[\mathbf{e}^{(6)}, \mathbf{e}^{(1)}, \mathbf{e}^{(9)}, \mathbf{e}^{(4)}, \mathbf{e}^{(2)}]$
 $\mathbf{y}' \in \mathbb{R}^{K}$:
 7

General Types of Learning (1/2)

• **Supervised learning**: learn to predict the labels of future data points

$$\begin{array}{lll} X \in \mathbb{R}^{N \times D} : & \overbrace{e^{(6)}, e^{(1)}, e^{(9)}, e^{(4)}, e^{(2)}]}^{K' \in \mathbb{R}^{D}} & x' \in \mathbb{R}^{D} : & \overbrace{y' \in \mathbb{R}^{K} : ?}^{K' \in \mathbb{R}^{N \times K}} \\ \hline & Unsupervised \ learning: \ learn \ (latent) \ patterns \ in \ X, \ and \ optionally \end{array}$$

0 generate new x's



Retrieve

When did the US break away from England?

Model output The US took the first step towards gaining independence from GB when it declared independence on July 2, 1776 (although the event is now commemorated on July 4, 1776, the date when the Declaration of Independence was officially adopted by Congress) [1][2]. The Treaty of Paris was later signed on September 3, 1783. formally separating the United States from the British Empire [3].





General Types of Learning (2/2)

 Reinforcement learning: learn from "good"/"bad" feedback of actions (instead of correct labels) to maximize the goal





General Types of Learning (2/2)

 Reinforcement learning: learn from "good"/"bad" feedback of actions (instead of correct labels) to maximize the goal





- AlphaGo [1] is a hybrid of reinforcement learning and supervised learning
 - Supervised learning from the game records
 - Then, reinforcement learning from self-play

Shan-Hung Wu (CS, NTHU)

Introduction

Data collection, preprocessing (e.g., integration, cleaning, etc.), and exploration

- Data collection, preprocessing (e.g., integration, cleaning, etc.), and exploration
 - Split a dataset into the training and testing datasets

- Data collection, preprocessing (e.g., integration, cleaning, etc.), and exploration
 - Split a dataset into the training and testing datasets
- Model development
 - **1** Assume a *model* $\{f(\cdot; w)\}$ that is a collection of candidate functions f's (representing posteriori knowledge) we want to discover
 - (1) f is assumed to be parametrized by w

- Data collection, preprocessing (e.g., integration, cleaning, etc.), and exploration
 - Split a dataset into the training and testing datasets
- Model development
 - **1** Assume a **model** $\{f(\cdot; w)\}$ that is a collection of candidate functions f's (representing posteriori knowledge) we want to discover
 - 1 f is assumed to be parametrized by w
 - Define a cost function C(w; X) (or functional C[f; X]) that measures "how good a particular f can explain the training data"

- Data collection, preprocessing (e.g., integration, cleaning, etc.), and exploration
 - Split a dataset into the training and testing datasets
- Model development
 - Assume a *model* {f(·;w)} that is a collection of candidate functions f's (representing posteriori knowledge) we want to discover
 f is assumed to be parametrized by w
 - Define a *cost function* $C(w; \mathbb{X})$ (or functional $C[f; \mathbb{X}]$)
 - Define a cost function C(w;X) (or functional C[f;X]) that measures "how good a particular f can explain the training data"
- 3 Training: employ an algorithm that finds the best (or good enough) function f*(·; w*) in the model that minimizes the cost function

$$w^* = \arg\min_w C(w; \mathbb{X})$$

Shan-Hung Wu (CS, NTHU)

- Data collection, preprocessing (e.g., integration, cleaning, etc.), and exploration
 - Split a dataset into the training and testing datasets
- Model development
 - Assume a *model* {f(·;w)} that is a collection of candidate functions f's (representing posteriori knowledge) we want to discover
 f is assumed to be parametrized by w
 - Define a cost function C(w; X) (or functional C[f; X]) that measures "how good a particular f can explain the training data"
- 3 Training: employ an algorithm that finds the best (or good enough) function f*(·; w*) in the model that minimizes the cost function

$$w^* = \arg\min_{w} C(w; \mathbb{X})$$

Testing: evaluate the performance of the learned f* using the testing dataset

Shan-Hung Wu (CS, NTHU)

- Data collection, preprocessing (e.g., integration, cleaning, etc.), and exploration
 - Split a dataset into the training and testing datasets
- Model development
 - Assume a *model* {f(·;w)} that is a collection of candidate functions f's (representing posteriori knowledge) we want to discover
 f is assumed to be parametrized by w
 - Define a cost function C(w; X) (or functional C[f; X]) that measures "how good a particular f can explain the training data"
- 3 Training: employ an algorithm that finds the best (or good enough) function f*(·; w*) in the model that minimizes the cost function

$$w^* = \arg\min_w C(w; \mathbb{X})$$

- Testing: evaluate the performance of the learned f* using the testing dataset
- 5 Apply the model in the real world

Shan-Hung Wu (CS, NTHU)

Introduction

- **1** Training dataset: $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i$
- 2 Testing dataset: $\mathbb{X}' = \{(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)})\}_i$

- **1** Training dataset: $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i$
- 2 Testing dataset: $\mathbb{X}' = \{(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)})\}_i$
- Model development

1 Model:
$$\{f : f(x; w) = w^{\top}x\}$$

- **1** Training dataset: $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i$
- 2 Testing dataset: $\mathbb{X}' = \{(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)})\}_i$
- Model development
 - **1 Model**: $\{f : f(x; w) = w^{\top}x\}$
 - **2** Cost function: $C(w; \mathbb{X}) = \Sigma_i 1(w; f(\mathbf{x}^{(i)}; w) \neq y^{(i)})$

- **1** Training dataset: $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i$
- 2 Testing dataset: $\mathbb{X}' = \{(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)})\}_i$
- Model development
 - **1 Model**: $\{f : f(x; w) = w^{\top}x\}$
 - **2 Cost function**: $C(w; \mathbb{X}) = \sum_i 1(w; f(\mathbf{x}^{(i)}; w) \neq y^{(i)})$
- **3** Training: to solve $w^* = \arg\min_{w} \Sigma_i 1(w; f(\mathbf{x}^{(i)}; w) \neq y^{(i)})$

- **1** Training dataset: $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i$
- 2 Testing dataset: $\mathbb{X}' = \{(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)})\}_i$
- Model development
 - **1 Model**: $\{f : f(x; w) = w^{\top}x\}$
 - **2 Cost function**: $C(w; \mathbb{X}) = \sum_i 1(w; f(\mathbf{x}^{(i)}; w) \neq y^{(i)})$
- **3** Training: to solve $w^* = \arg\min_{w} \Sigma_i 1(w; f(x^{(i)}; w) \neq y^{(i)})$
- **• Testing**: accuracy $\frac{1}{|\mathbb{X}'|} \Sigma_i 1(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)}; f(\mathbf{x}'^{(i)}; \mathbf{w}^*) = \mathbf{y}'^{(i)})$

- **1** Training dataset: $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i$
- 2 Testing dataset: $\mathbb{X}' = \{(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)})\}_i$
- Model development
 - **1** *Model*: $\{f : f(x; w) = w^{\top}x\}$
 - **2** Cost function: $C(w; \mathbb{X}) = \Sigma_i 1(w; f(\mathbf{x}^{(i)}; w) \neq y^{(i)})$
- **3** Training: to solve $w^* = \arg\min_{w} \Sigma_i 1(w; f(x^{(i)}; w) \neq y^{(i)})$
- **④ Testing**: accuracy $\frac{1}{|\mathbb{X}|} \Sigma_i 1(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)}; f(\mathbf{x}'^{(i)}; \mathbf{w}^*) = \mathbf{y}'^{(i)})$
- **(5)** Use f^* to predict the labels of your future emails

- **1** Training dataset: $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i$
- 2 Testing dataset: $\mathbb{X}' = \{(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)})\}_i$
- Model development
 - **1** *Model*: $\{f : f(x; w) = w^{\top}x\}$
 - **2** Cost function: $C(w; \mathbb{X}) = \Sigma_i 1(w; f(\mathbf{x}^{(i)}; \mathbf{w}) \neq y^{(i)})$
- **3** Training: to solve $w^* = \arg\min_{w} \Sigma_i 1(w; f(x^{(i)}; w) \neq y^{(i)})$
- **④ Testing**: accuracy $\frac{1}{|\mathbb{X}|} \Sigma_i 1(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)}; f(\mathbf{x}'^{(i)}; \mathbf{w}^*) = \mathbf{y}'^{(i)})$
- **(5)** Use f^* to predict the labels of your future emails
 - See Notation

Outline

1 What's Machine Learning?

2 What's Deep Learning?

3 About this Course...



Shan-Hung Wu (CS, NTHU)

Introduction

Deep Learning

• ML where an $f(\cdot; \mathbf{w})$ has many (deep) layers

$$\hat{\mathbf{y}} = f^{(L)}(\cdots f^{(2)}(f^{(1)}(\mathbf{x}; \mathbf{w}^{(1)}); \mathbf{w}^{(2)}) \cdots; \mathbf{w}^{(L)})$$
$$\mathbf{x} \rightarrow f^{(1)}(\cdot; \mathbf{w}^{(1)}) \rightarrow f^{(2)}(\cdot; \mathbf{w}^{(2)}) \rightarrow \cdots \qquad f^{(L)}(\cdot; \mathbf{w}^{(L)}) \rightarrow \hat{\mathbf{y}}$$

Deep Learning

• ML where an $f(\cdot; w)$ has many (deep) layers

$$\hat{\mathbf{y}} = f^{(L)}(\cdots f^{(2)}(f^{(1)}(\mathbf{x}; \mathbf{w}^{(1)}); \mathbf{w}^{(2)}) \cdots; \mathbf{w}^{(L)})$$
$$\mathbf{x} \rightarrow f^{(1)}(\cdot; \mathbf{w}^{(1)}) \rightarrow f^{(2)}(\cdot; \mathbf{w}^{(2)}) \rightarrow \cdots \qquad f^{(L)}(\cdot; \mathbf{w}^{(L)}) \rightarrow \hat{\mathbf{y}}$$

Pros:

- Learns to pre-process data automatically
- Learns a complex function (e.g., visual objects to labels)

Deep Learning

• ML where an $f(\cdot; w)$ has many (deep) layers

$$\hat{\mathbf{y}} = f^{(L)}(\cdots f^{(2)}(f^{(1)}(\mathbf{x}; \mathbf{w}^{(1)}); \mathbf{w}^{(2)}) \cdots; \mathbf{w}^{(L)})$$
$$\mathbf{x} \rightarrow f^{(1)}(\cdot; \mathbf{w}^{(1)}) \rightarrow f^{(2)}(\cdot; \mathbf{w}^{(2)}) \rightarrow \cdots \qquad f^{(L)}(\cdot; \mathbf{w}^{(L)}) \rightarrow \hat{\mathbf{y}}$$

Pros:

- Learns to pre-process data automatically
- Learns a complex function (e.g., visual objects to labels)

Cons:

- Usually needs large data to train a model well
- Higher computation costs (for both training and testing)

Shan-Hung Wu (CS, NTHU)

Outline

1 What's Machine Learning?

2 What's Deep Learning?



3 About this Course...



Shan-Hung Wu (CS, NTHU)

Introduction

Target Audience

- Graduate (and senior undergraduate) CS students
 - Easy-to-moderate level of theory
 - Coding and engineering (in Python)
 - Clean datasets (small & large)

Target Audience

- Graduate (and senior undergraduate) CS students
 - Easy-to-moderate level of theory
 - Coding and engineering (in Python)
 - Clean datasets (small & large)
- No prior knowledge about ML is needed

Topics Covered

• Supervised, unsupervised learning, and reinforcement learning

Topics Covered

Supervised, unsupervised learning, and reinforcement learning
with *structural* output:



A man holding a tennis racquet on a tennis court.



A group of young people playing a game of Frisbee



Two pizzas sitting on top of a stove top oven



A man flying through the air while riding a snowboard

- Part 1: math review (2 weeks)
 - Linear algebra
 - Probability & information theory
 - Numerical optimization

- Part 1: math review (2 weeks)
 - Linear algebra
 - Probability & information theory
 - Numerical optimization
- Part 2: machine learning basics (3 weeks)
 - Learning theory
 - Parametric/non-parametric models
 - Experiment design

- Part 1: math review (2 weeks)
 - Linear algebra
 - Probability & information theory
 - Numerical optimization
- Part 2: machine learning basics (3 weeks)
 - Learning theory
 - Parametric/non-parametric models
 - Experiment design
- Part 3: deep supervised learning (6 weeks)
 - Neural Networks (NNs), CNNs, RNNs, Transformers

- Part 1: math review (2 weeks)
 - Linear algebra
 - Probability & information theory
 - Numerical optimization
- Part 2: machine learning basics (3 weeks)
 - Learning theory
 - Parametric/non-parametric models
 - Experiment design
- Part 3: deep supervised learning (6 weeks)
 - Neural Networks (NNs), CNNs, RNNs, Transformers
- Part 4: unsupervised learning & generative AI (2 weeks)
 - · Autoencoders, manifold learning, generative models

- Part 1: math review (2 weeks)
 - Linear algebra
 - Probability & information theory
 - Numerical optimization
- Part 2: machine learning basics (3 weeks)
 - Learning theory
 - Parametric/non-parametric models
 - Experiment design
- Part 3: deep supervised learning (6 weeks)
 - Neural Networks (NNs), CNNs, RNNs, Transformers
- Part 4: unsupervised learning & generative AI (2 weeks)
 - · Autoencoders, manifold learning, generative models
- Part 5: reinforcement learning (2 weeks)
 - Value/gradient policies, action/critics, reinforce RNNs

Shan-Hung Wu (CS, NTHU)

Introduction

Grading (Tentative)

Prerequisite quiz: 10%

- On next Thu (9/12)
- You can take this course if you are within top-90
- Contests (x 3): 45%
 - At the end of each part
- Assignments: 20%
 - Come with the labs
- Final exam: 25%
- Bonus: up to 10%
 - Math labs (x 4)
 - Optional ML topics (x 2)
 - Course participations

Classes Info

- Lectures on Tue (2 hours)
 - Concepts & theories
 - with companion videos
- Labs on Thu (1 hour)
 - Implementation (in Python) & engineering topics
- TA time: 4:20pm–5:30pm on Thu at Delta 729
- More info can be found in the course website

Outline

1 What's Machine Learning?

2 What's Deep Learning?

3 About this Course...



Q: Should we team up for the contests?

A: Yes, 2~4 students per team

- **Q**: Should we team up for the contests?
- A: Yes, 2~4 students per team
- **Q:** Which GPU card should I buy?
- A: Nvidia GTX 1060 or above; 1050 Ti (4G RAM) minimal

- Q: Should we team up for the contests?
- A: Yes, 2~4 students per team
- **Q:** Which GPU card should I buy?
- A: Nvidia GTX 1060 or above; 1050 Ti (4G RAM) minimal
- **Q:** Do we need to attend the classes?
- A: No, as long as you can pass. But we offer attendance bonus...

- Q: Should we team up for the contests?
- A: Yes, 2~4 students per team
- **Q:** Which GPU card should I buy?
- A: Nvidia GTX 1060 or above; 1050 Ti (4G RAM) minimal
- **Q:** Do we need to attend the classes?
- A: No, as long as you can pass. But we offer attendance bonus...
- Q: Is this a light-loading course or heavy-loading one?
- A: Should be very heavy to most students. Please reserve your time

FAQ (2/2)

- **Q:** What's the textbook?
- A: No formal textbook. But if you need one, read the Deep Learning book

FAQ (2/2)

- **Q**: What's the textbook?
- A: No formal textbook. But if you need one, read the Deep Learning book
- Q: Why some sections are marked with "*" or "**" in the slides?
- A: The mark "*" means "can be skipped for the first time reader," and "**" means "materials for reference only"

TODO

- Assigned reading:
 - Calculus
 - Get your feet wet with Python

Reference I

 D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al.

Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.