## Machine Learning to Language Model **Topic 01 - Introduction to Machine Learning**

Jaihua Yen https://jaihuayen.github.io/

## Contents

- Al Models
- Gradient Descent
- Demo
- Wrap Up



A Models

## **Al Models**



**MACHINE LEARNING** Algorithms whose performance improve as they are exposed to more data over time.



https://www.globaltechcouncil.org/artificial-intelligence/clearing-the-confusion-ai-vs-machine-learning-vs-deep-learning/

#### **ARTIFICIAL INTELLIGENCE**

A program that can sense, reason, act, and adapt.

#### **DEEP LEARNING**

Subset of machine learning in which multilayered neural networks learn from vast amount of data.

### The Era of Al Al could finally be introduced into practice in general tasks!

- Tremendous improvement in computational resources.
- Enhancement of model architecture for model efficiency.
- Release open-source large pre-trained general models.  $\bullet$
- Development of novel model training algorithms.  $\bullet$









[MCUNet: Tiny Deep Learning on IoT Devices, Han et al. 2020]

### **Al Models** Al models could be used in these tasks:

	Forecasting	Clas
Tableau - Forecasting       ile     Data       Worksheet     Dashboard       Summarize       H	ap Format Server Window Help   Image: Apple and App	
H Average Line H Average Line Hedian with Quartiles	Order Pr     Order Pr       Critical     S       0K     OK       0K     Medium       Low	Incoming messages

#### ssification



.com/blog/classification-machine-learning

#### Generating

#### What is natural language processing

Natural Language Processing (NLP) is a branch of Artificial Intelligence (AI) that focuses on the interaction between computers and humans using natural language. The goal of NLP is to enable computers to understand, interpret, and generate human language.

NLP involves the use of techniques from computer science, linguistics, and machine learning to process, analyze, and generate natural language. Some common applications of NLP include sentiment analysis, language translation, text classification, chatbots, and speech recognition.

NLP is a rapidly evolving field with new developments and advancements being made regularly. As computers become better at processing language, the potential applications for NLP continue to expand, making it an important field of study in both industry and academia.

https://openai.com/blog/chatgpt



# Al models are functions!



#### Translate "你好" in English.



## **Al Models** Al models are functions built by <u>neural networks</u>



#### **Universal Approximation Theorem**

 $g(x) \approx f(x)$ 

Any function can be approximated by neural network!

https://cs231n.github.io/convolutional-networks/

#### The function we want to learn from data by training

# How to Train Al Models?

### **Machine Learning** Mathematical Fundations of Model Training



https://blogs.sap.com/2019/04/05/machine-learning-in-sap-strategy/

#### **Training Dataset** Given datasets and tags to train AI model.



X

https://www.cs.toronto.edu/~kriz/cifar.html



y
airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

#### **Model Training Cannot perfectly predict the true label in the first time**



We want to improve the model output in order to have better performance!

#### The difference between the true label and predicted label

#### **Loss Function Measurement of model performance**

The difference between the true label and predicted label



## Loss function (Cross Entropy Loss) $L = -\sum y_i \log \hat{y}_i$

The higher the predicted probability  $y_i = 1$ the lower the loss you will get

i=1

 $v_i = 0$ The whole term will be 0 so it doesn't matter



#### **Loss Function Measurement of model performance**

Loss function (Cross Entropy Loss) The difference between the true label and predicted label  $L = -\sum_{i} y_i \log \hat{y}_i$ V i=1 $= -\sum_{i} y_i \log g(x_i)$ ship automobile i=1 $L(w) \triangleq -\sum_{i} y_i \log(wx_i)$ i=1





#### **Loss Function Measurement of model performance**



### **Training Procedure Start for Random Initialization of Weight**





#### **Gradient Descent** A Approach to lower the function value



#### **Gradient Descent** A Approach to lower the function value



### **Gradient Descent A High-Dimension Overview**



https://blog.paperspace.com/intro-to-optimization-in-deep-learning-gradient-descent/

#### Update the three parameters at the same time!



### **Gradient Descent A High-Dimension Overview**



https://blog.paperspace.com/intro-to-optimization-in-deep-learning-gradient-descent/

$$w = [x, y, z]$$
$$w_t = w_{t-1} - i \nabla L(x, y, z)$$

This is called gradient, which makes this method called gradient descent!

$$\begin{bmatrix} x_t \\ y_t \\ z_t \end{bmatrix} = \begin{bmatrix} x_{t-1} \\ y_{t-1} \\ z_{t-1} \end{bmatrix} - r \begin{bmatrix} \frac{\partial L}{\partial x_{t-1}} \\ \frac{\partial L}{\partial y_{t-1}} \\ \frac{\partial L}{\partial z_{t-1}} \end{bmatrix}$$



### Backpropogation **A High-Dimension Overview**

Use loss to update the weights in backward order



https://cs231n.github.io/convolutional-networks/

## $\boldsymbol{w} = [x, y, z]$ $w_t = w_{t-1} - r\nabla L(x, y, z)$

$$\begin{bmatrix} x_t \\ y_t \\ z_t \end{bmatrix} = \begin{bmatrix} x_{t-1} \\ y_{t-1} \\ z_{t-1} \end{bmatrix} - r \begin{bmatrix} \frac{\partial L}{\partial x_{t-1}} \\ \frac{\partial L}{\partial y_{t-1}} \\ \frac{\partial L}{\partial z_{t-1}} \end{bmatrix}$$

Let's do this in <u>Colab</u>!



- We only predict the next word based on the previous word.
- Model the predicted probability of a certain word based on a given word.

 $C_i$  is the word in the *i* position.

• Here we give an example of a sentence: <s> Al could finally be introduced into practice in general tasks <e>

$$P(c_i \mid c_{i-1})$$

<s> AI could finally be introduced into practice in general tasks <e>

 $C_{i-1}$   $C_i$ AI <S>

AI

<s> AI could finally be introduced into practice in general tasks <e>

 $C_{i-1}$  $C_i$ 



<s> Al could finally be introduced into practice in general tasks <e>

 $C_{i-1}$  $C_i$ 

tasks

<e>

We want to train the following neural network







Let's also do this in <u>Colab</u>!

## Questions

- What if having too large/small learning rate (r)?
- What if we cannot cover all the cases so that some conditional probability is zero?
- What if we also consider the penalty of the loss of predicting the class that does not contain the ground truth label?



Wrap Up

## What We Have Gone Through

- Applications in machine learning
- Training a machine learning model
- Gradient descent
- Bi-gram language model



## What's Next

- Deep Neural Network
- Word Embedding
- Transfer Learning





