# Deep Learning and Applications

DSA 5204 • Lecture 1
Dr Low Yi Rui (Aaron)
Department of Mathematics



National University
of Singapore

# Information

**Instructor**
        Low Yi Rui (Aaron)
        Office: S17-05-19
        Email: matyrl@nus.edu.sg

**Slides and demo code will be available on Canvas**

**Reference Book**
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- https://www.deeplearningbook.org/

# Hybrid Classroom

## Lesson Plan

- One or two 10 min breaks
- Q&A sessions after the breaks
- Q&A at the end of class

## Zoom instructions

- Join the Zoom session on Canvas
- Type questions in the chat, or ask during the Q&A sessions
- Mute your microphones, but switch them on when you'd like to ask questions
- You are encouraged to turn on your webcam when asking questions

## PollEverywhere

- All polls can be accessed at

    https://pollev.com/aaronyrlow

# Attendance for SSG funded students

**Attendance will be taken via QR code during the first break**

# Assessment Policies

**Assessment**

- Test (30%)
- Graded Homework (10% × 2)
- Course Project (Report + Presentation: 50%)

**Late submission policy for homework and report**

- Please drop me an email if you are going to submit late
- Late submission penalty: -20% of the total grade per day late

  e.g. Actual score = 7/10. Submitted 1 day late.
  Adjusted score = 5/10.

# Consultation and Forum

**Consultation slots**

- available on Canvas (If filled, additional ones will be created)
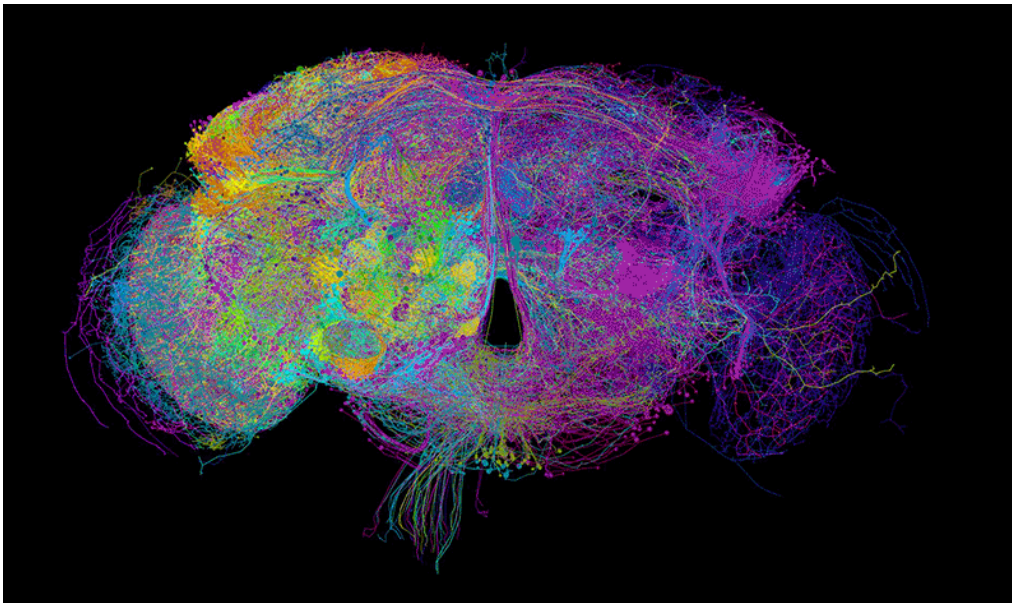- Currently set to 30min per session, Saturday 2-3pm

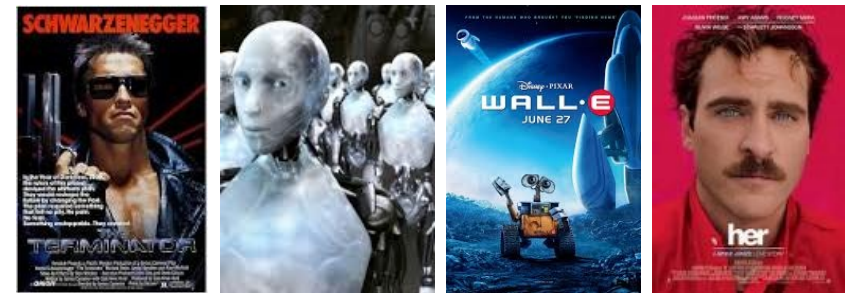**Discussion forum on Canvas**

# AI, Machine Learning & Deep Learning

# The goal of the study of AI

To **understand** and **replicate** intelligence



*Neurons in a fruitfly's brain*
*Zheng et al. Cell 2018, 174 (3), 730-743.e22.*



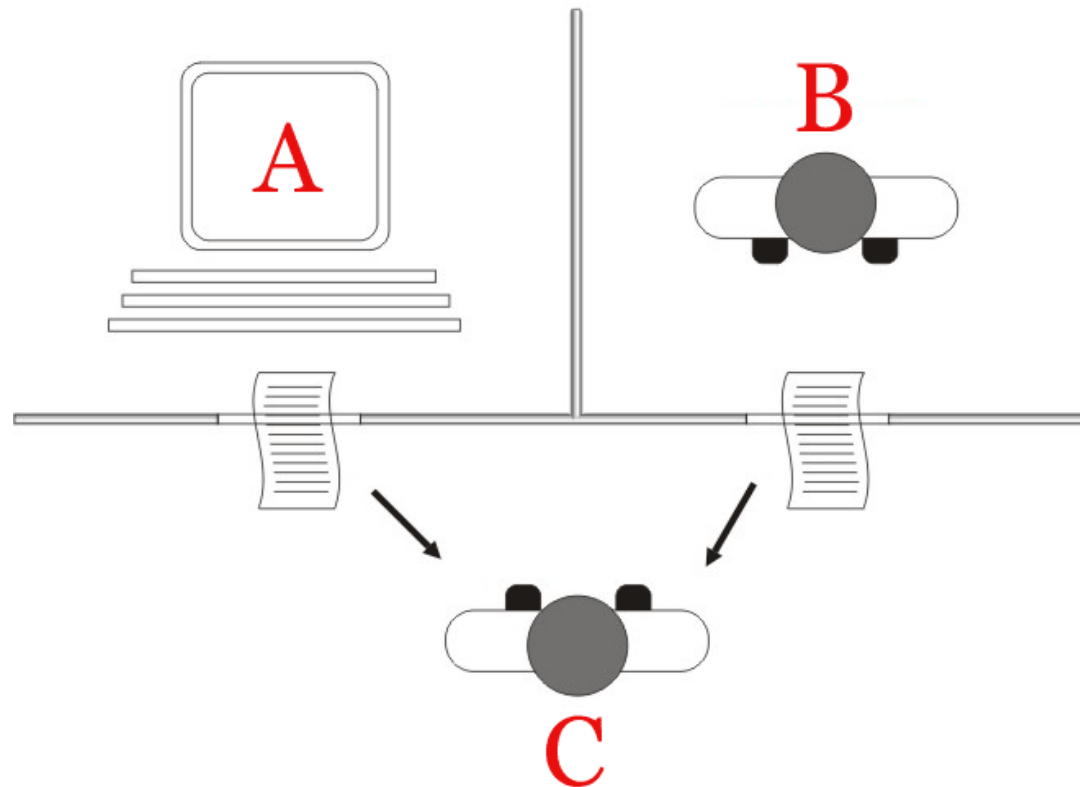Evolution of the concept of replication

# Can machines think?

*I propose to consider the question, "Can machines think?". This should begin with definitions of the meaning of the terms "machine" and "think". The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous...*

Alan Turing, 1950

Turing, A. M. "Computing machinery and intelligence". Mind 49 433-460.

# The Imitation Game



*I believe that in about fifty years' time it will be possible, to programme computers ... to make them play the imitation game so well that an average interrogator will not have more than 70 percent chance of making the right identification after five minutes of questioning.*
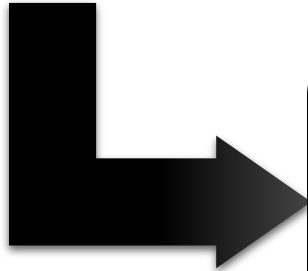
Alan M Turing, 1950

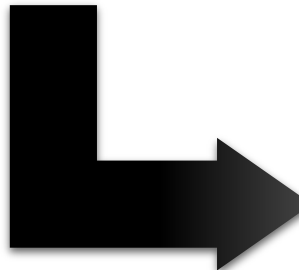https://en.wikipedia.org/wiki/Turing_test#/media/File:Turing_test_diagram.png

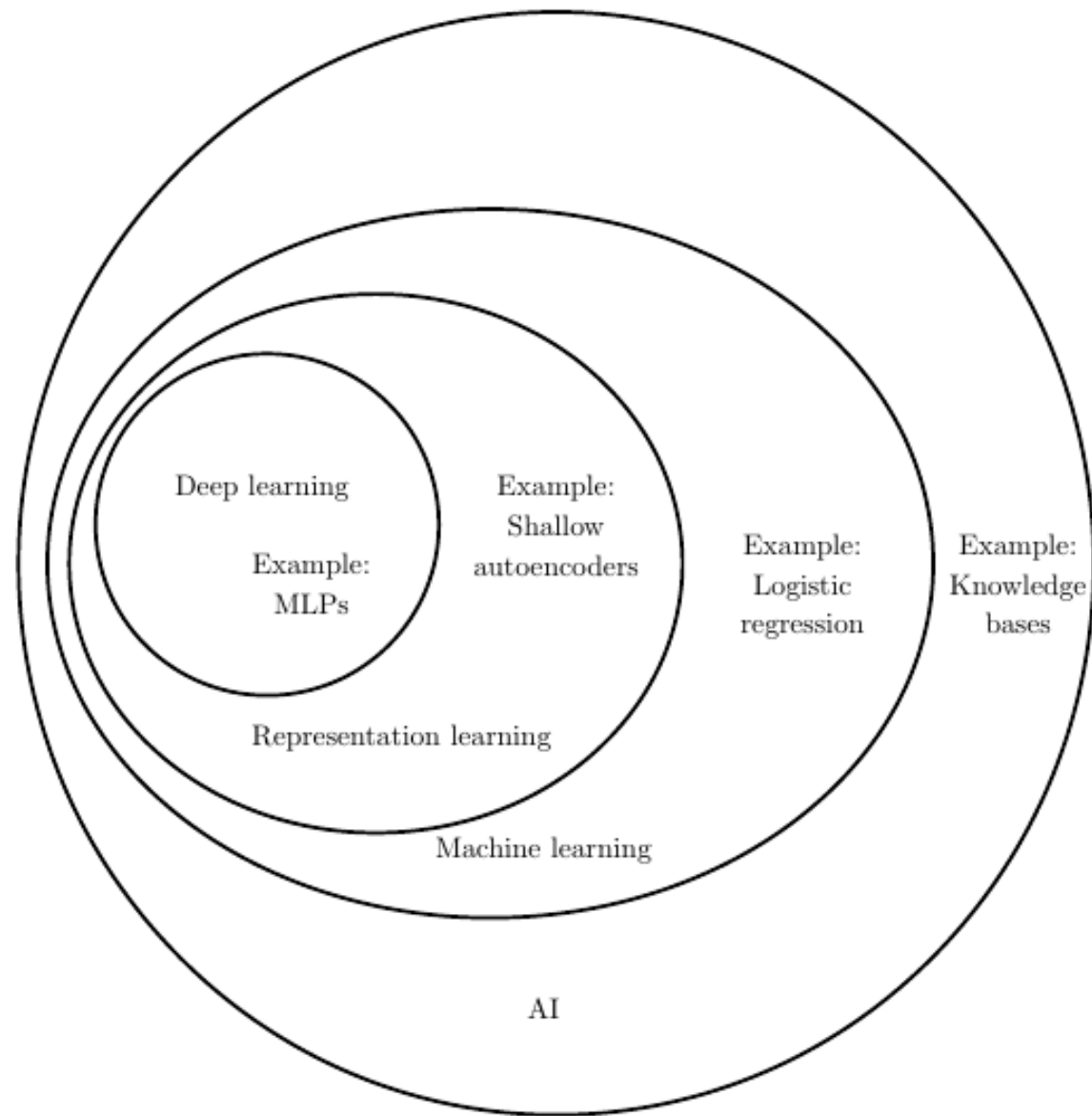# From artificial intelligence to machine learning

Can machines think?

Can machines do what thinking beings do?

How can machines learn to do some things that thinking beings do?
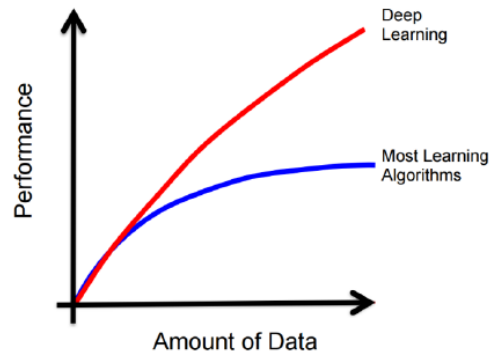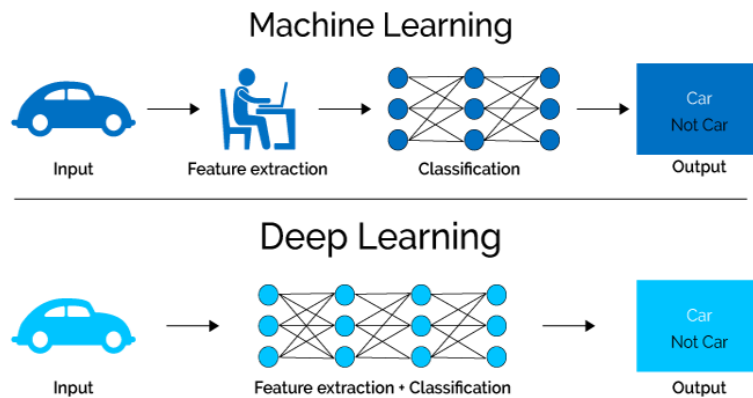
# A Brief History of Deep Learning

- 1940s: Artificial neural networks
- 1950s: Roots of the backpropagation algorithm
- 1950s-80s: Perceptron, RNNs
- 1980s-90s: CNN, LSTM, Bidirectional RNNs
- 2000s-Present: deep learning has become the state of the art for a variety of applications
  - AlexNet (2012)
  - DeepFace (2014)
  - Generative Adversarial networks (GANs) (2014)
  - AlphaGo (2016)
  - BERT (2018)
  - Molecular Dynamics, Protein folding, etc (2020-21)
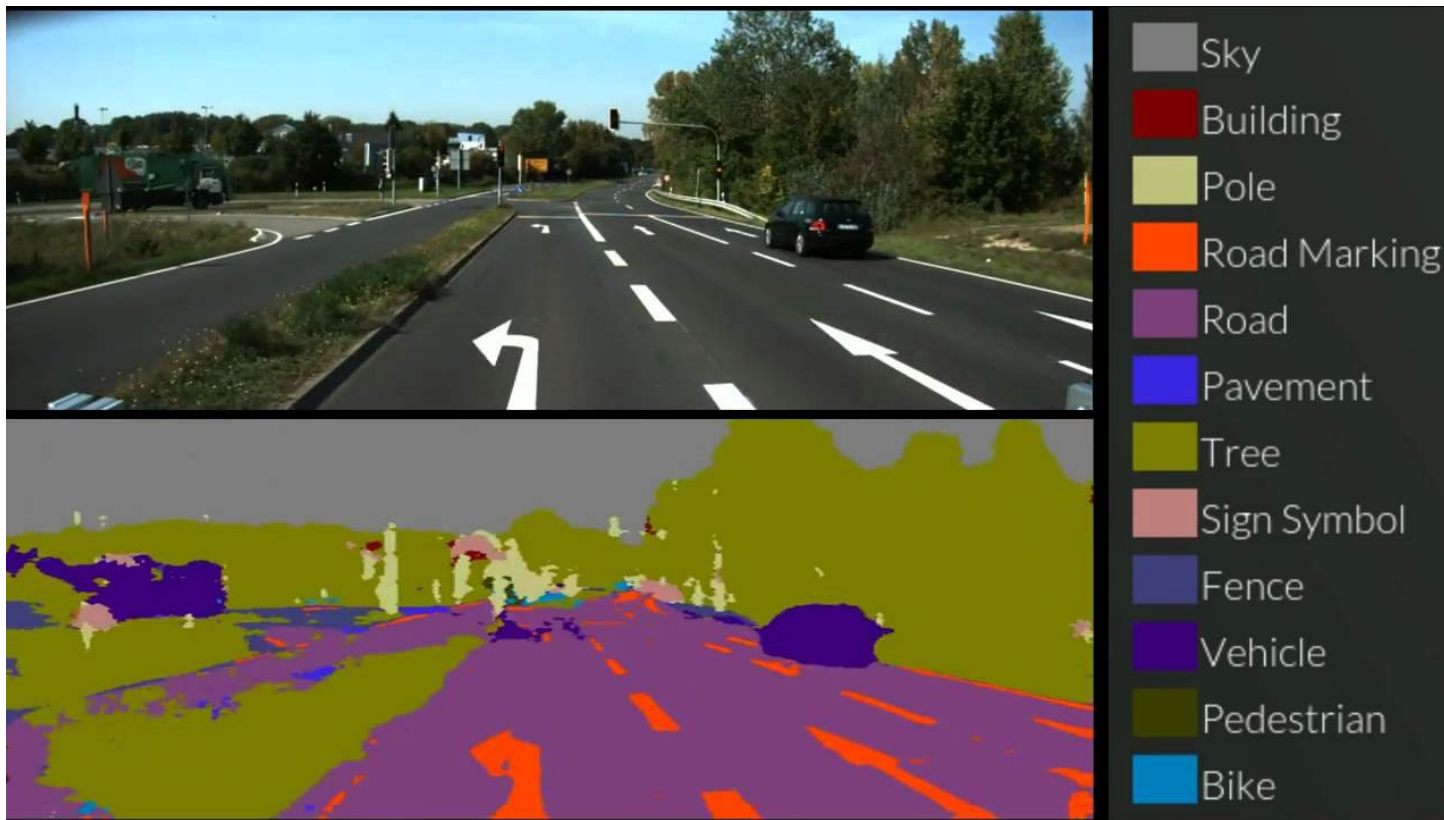
# Why deep learning now?

## Data



## Computing

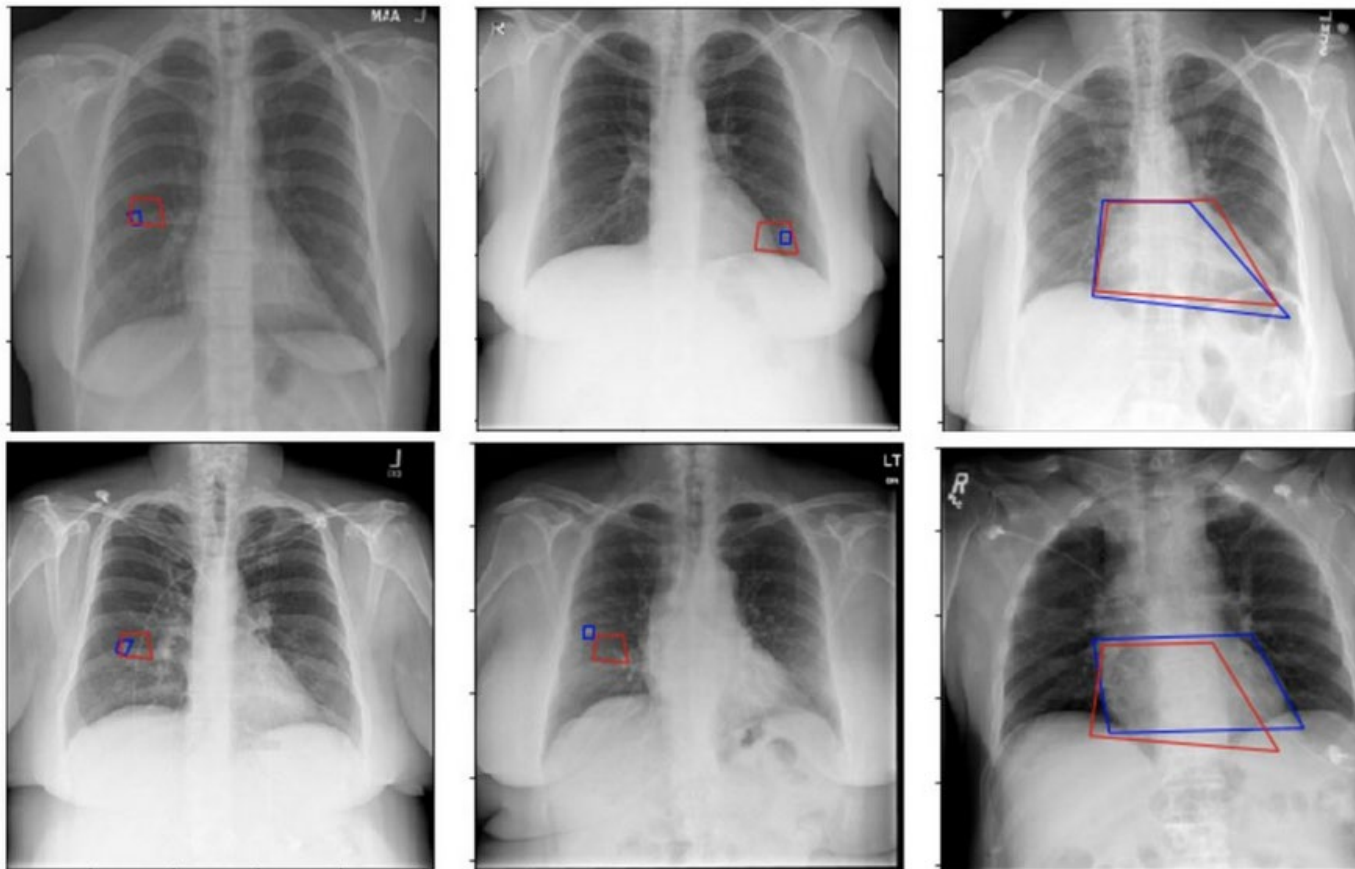# Deep Learning in Autonomous Driving



https://www.youtube.com/watch?v=kMMbW96nMW8
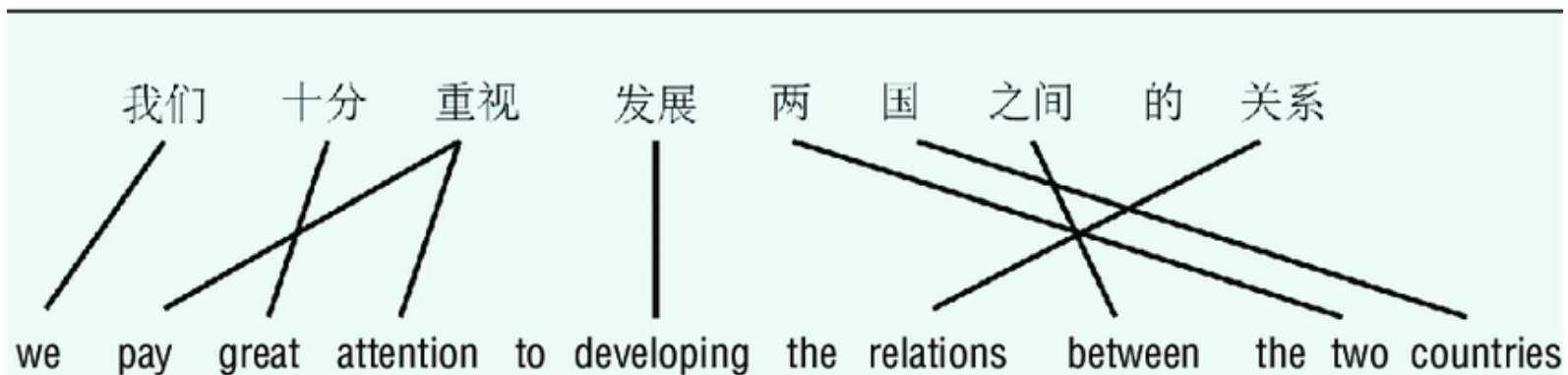
# Deep Learning in Medical Imaging



*Moradi, Mehdi, et al. International Conference on Medical Image Computing and Computer-Assisted Intervention. Springer, Cham, 2018.*

# Deep Learning in Machine Translation

| | |
|---|---|
| Word alignment | FNN, RecurrentNN |
| Translation rule selection | FNN, RAE, CNN |
| Reordering and structure prediction | RAE, RecurrentNN, RecursiveNN |
| Language model | FNN, RecurrentNN |
| Joint translation prediction | FNN, RecurrentNN, CNN |



*Zhang, Jiajun, and Chengqing Zong. "Deep neural networks in machine translation: An overview." (2015).*
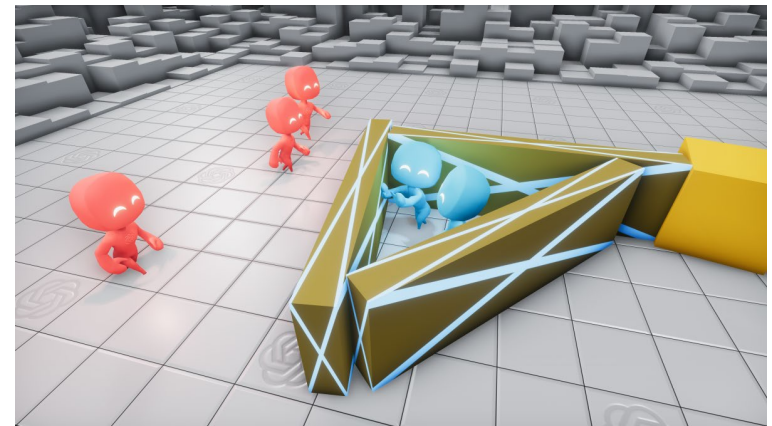
# Deep RL: towards AGI

## AlphaGo/AlphaZero



## AlphaStar



## OpenAI Five


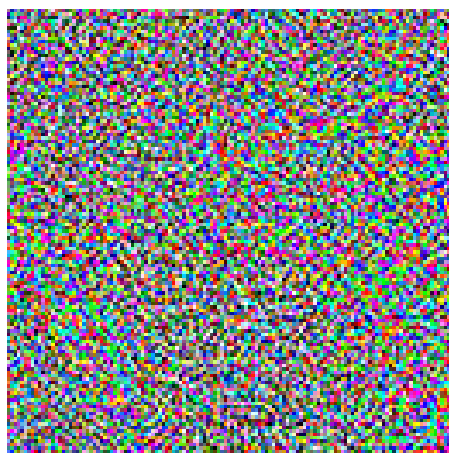
## OpenAI Hide-and-Seek

# Limitations of Deep Learning



"panda"
57.7% confidence
+ ε
"gibbon"
99.3% confidence

# This Class

Modern introduction to deep learning and its applications

- Architectures
- Learning algorithms
- Practical techniques
- Applications via demos/homework/projects
- An overview of active areas of DL research

# Deep Learning Research

- Deep learning research is very active!
- Learning how to learn is very important
- The course project serves this purpose (more info later)

Statistics of acceptance rate NeurIPS

# Mathematics and Programming

Mathematics
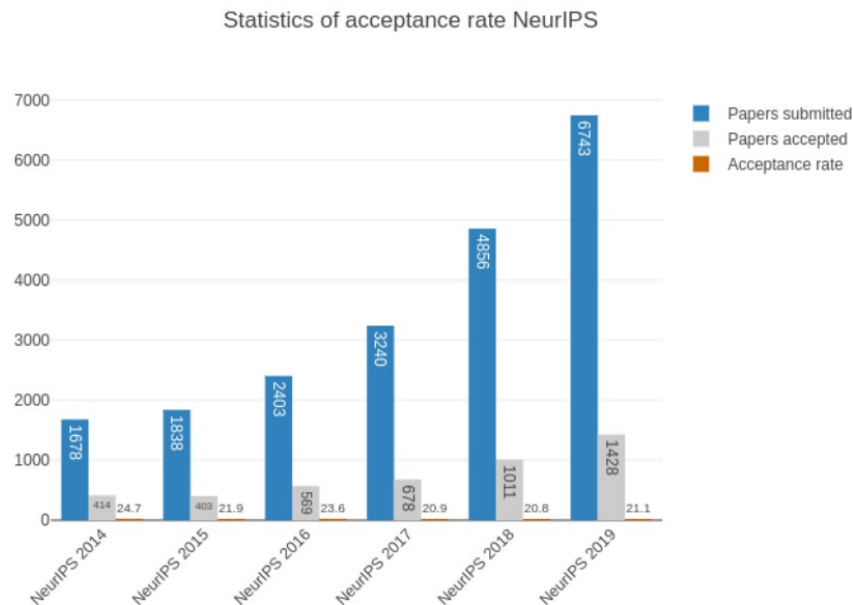
- Linear algebra
- Probability
- *Deep learning Book* (Chapters 2-3)

Programming

- Basic Python
- Useful libraries to know: Numpy, Scipy, Sklearn, Pandas, Seaborn/Matplotlib
- Deep learning libraries: Tensorflow(Keras)/Pytorch
- Tutorials are widely available online

Bishop, *Pattern Recognition and Machine Learning.*

DSA5102/
DSA5105

DSA5204

When poll is active respond at

**PollEv.com/aaronyrlow**



Have you taken DSA5102 or DSA5105 (or equivalent?)

# 0 surveys completed

0 surveys underway

# Machine Learning Basics

# A concrete definition of learning

*A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E.*

**Tom Mitchell**

**Programming**

| | | |
|---|---|---|
| $8 \div 2(2 + 2)$ | Inputs → | |
| | | Computer |
| calc(*args) | Program → | |

Outputs → 16

**Learning**

| | | |
|---|---|---|
| $8 \div 2(2 + 2)$ | Inputs → | |
| | | Computer |
| 16 | Outputs → | |

Program → calc(*args)

# Tasks

Prediction



Transcription



following

Translation

# Tasks

Imputation



Generation

# Experience

Experience = Data

Usually represented by vectors/matrices/tensors

# Performance

Performance measures how well a machine learning model does in a given task.

Examples:

- Regression: mean squared error of predictions
- Classification: accuracy of classification
- Image generation?

# Other issues

- Learning algorithm: how do we improve $P$ on task $T$ with experience $E$?

- What is an objective measure of performance $P$?

# Linear Regression as a Canonical Example

# The T, E, P of Regression

(T) Task:

- Inputs: $x \in \mathbb{R}^d$    Outputs: $y \in \mathbb{R}$
- Underlying relationship: $y = f^*(x)$
- We want to predict outputs given inputs

(E) Experience:
$$\mathcal{D} = \{x_i, y_i = f^*(x)\}_{i=1}^{N}$$

or a noisy version
$$\mathcal{D} = \{x_i, y_i = f^*(x) + \epsilon_i\}_{i=1}^{N}$$



$f^*$

Outputs ($y_i$)

Inputs ($x_i$)

# The T, E, P of Regression

Before discussing performance, we need to define a **hypothesis space**

$$\mathcal{H} = \{\text{collection of candidates } f\}$$

(P) Performance of a $f \in \mathcal{H}$:

$$\text{Distance}(f, f^*)$$

Learning:

Find a $\hat{f}$ that is closest to $f^*$

$$\hat{f} = \arg\min_{f \in \mathcal{H}} \text{Distance}(f^*, f)$$

# How is distance measured?

A notion of distance can be defined by a **loss function** $L$ evaluated over the dataset

$$\text{Distance} = \frac{1}{N} \sum_{i=1}^{N} L(f(\boldsymbol{x}_i), y_i)$$

Example: square loss

$$L(y', y) = \frac{1}{2}(y' - y)^2$$

We will encounter many other loss functions throughout this class!

This distance is called the **empirical risk**, and we denote it by $R_{\text{emp}}(f, \mathcal{D})$

# Empirical Risk Minimization

With the empirical risk at hand, we can concretely formulate the sentence

*Increase performance P with experience E on task T*

as an optimization problem

$$\hat{f} = \mathrm{argmin}_{f \in \mathcal{H}}\, R_{\mathrm{emp}}(f, \mathcal{D})$$

This is known as **empirical risk minimization** (ERM)

# Example: Linear Regression

Linear regression is a particular choice of $\mathcal{H}$ -- linear functions

$$\mathcal{H} = \{f(x) = w^T x \; : \;\; w \in \mathbb{R}^d\}$$

$$\mathcal{H}$$

# Example: Linear Regression

Empirical risk minimization under square loss:

$$\min_{w} R_{\text{emp}}(\boldsymbol{w}) = \frac{1}{2N} \sum_{i=1}^{N} (\boldsymbol{w}^T \boldsymbol{x}_i - y_i)^2$$

How do we solve this?

# The Least Squares Formula

Define the design matrix and output vector

$$X = \begin{pmatrix} - & \boldsymbol{x}_1 & - \\ - & \boldsymbol{x}_2 & - \\ \vdots & \vdots & \vdots \\ - & \boldsymbol{x}_N & - \end{pmatrix} \qquad \boldsymbol{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

Then, the ERM can be rewritten as

$$\min_{\boldsymbol{w}} R_{\text{emp}}(\boldsymbol{w}) = \frac{1}{2N} \|X\boldsymbol{w} - \boldsymbol{y}\|^2$$

# The Least Squares Formula

The minimum of $R_{\text{emp}}(\boldsymbol{w})$ can be found by

$$\nabla_{\boldsymbol{w}} R_{\text{emp}}(\boldsymbol{w}) = 0$$

This gives the equation

$$X^T(X\boldsymbol{w} - \boldsymbol{y}) = 0$$

which yields

$$\boldsymbol{w} = \widehat{\boldsymbol{w}} = (X^TX)^{-1}X^T\boldsymbol{y}$$
$$\hat{f}(\boldsymbol{x}) = \widehat{\boldsymbol{w}}^T\boldsymbol{x}$$

Ordinary Least Squares Formula

# Extension to Affine Functions

Very often, linear functions are not rich enough to model our data.

In the simplest case, we can consider **affine functions**

$$\mathcal{H} = \{f(x) = \boldsymbol{w}^T x + b \ : \ \boldsymbol{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$$

Linear  $\mathcal{H}$

Affine  $\mathcal{H}$

# Extension to Affine Functions

The least squares formula is the same, except we change the definition of the design matrix

$$X = \begin{pmatrix} - & \boldsymbol{x}_1 & - \\ - & \boldsymbol{x}_2 & - \\ \vdots & \vdots & \vdots \\ - & \boldsymbol{x}_N & - \end{pmatrix} \quad \Rightarrow \quad X = \begin{pmatrix} - & \boldsymbol{x}_1 & - & 1 \\ - & \boldsymbol{x}_2 & - & 1 \\ \vdots & \vdots & \vdots & \vdots \\ - & \boldsymbol{x}_N & - & 1 \end{pmatrix}$$

$$w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix} \quad \Rightarrow \quad w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \\ b \end{pmatrix}$$

# Model Capacity and Underfitting

With every $\mathcal{H}$ comes the concept of **capacity**:

*How large is the hypothesis space $\mathcal{H}$?*



(a) $f^*(x) = 1 + 2x$    (b) $f^*(x) = x^2$

under-fitting

# Extension to Linear Basis Models

Clearly, linear models cannot fit functions like

$$f^*(x) = x^2 \qquad f^*(x) = x^3 \qquad f^*(x) = \cos x$$

This motivates **linear basis models** in the form

$$\mathcal{H} = \left\{ f(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}) = \sum_{i=1}^{m} w_i \phi_i(\boldsymbol{x}) \; : \; \boldsymbol{w} \in \mathbb{R}^m \right\}$$

The functions $\boldsymbol{\phi} = (\phi_1, \dots, \phi_m)$ are called **basis functions** and are chosen *a priori.* They are also called **feature maps.**

# Examples of basis functions

Some choices of basis functions in 1D

- Polynomial basis: $\phi_j(x) = x^j$

- Gaussian basis: $\phi_j(x) = \exp\left(-\frac{(x-m_j)^2}{2s^2}\right)$

- Sigmoid basis: $\phi_j(x) = \sigma\left(\frac{x-m_j}{s}\right)$ with $\sigma(b) = \frac{1}{1+e^{-b}}$



Polynomial          Gaussian          Sigmoid

# Capacity of Linear Basis Models

Note that linear basis models strictly generalizes linear/affine models!

- Set $\phi_i(x) = x_i$ for $i = 1, \ldots, d$ and $\phi_{d+1}(x) = 1$

Using the basis

$\phi_i(x) = x^{i-1}$ for $i = 1,2,3$

We can easily fit this with a linear model

$\Longrightarrow$



(b) $f^*(x) = x^2$

Legend: $\hat{f}$ (line), $(x_i, y_i)$ (points)

# Universality

So, how big is the capacity of linear basis models?

With appropriate choices of basis functions and $m \to \infty$, they are **universal**, in the sense that they can approximate* any* function/relationship!
Examples: Fourier series, Weierstrass approximation

# Least Squares Formula for Linear Basis Models

Least square formula obtained by changing the design matrix

$$X = \begin{pmatrix} - & \boldsymbol{x}_1 & - \\ - & \boldsymbol{x}_2 & - \\ \vdots & \vdots & \vdots \\ - & \boldsymbol{x}_N & - \end{pmatrix} \Rightarrow \Phi = \begin{pmatrix} - & \boldsymbol{\phi}(\boldsymbol{x}_1) & - \\ - & \boldsymbol{\phi}(\boldsymbol{x}_2) & - \\ \vdots & \vdots & \vdots \\ - & \boldsymbol{\phi}(\boldsymbol{x}_N) & - \end{pmatrix}$$

The new empirical risk is

$$R_{\text{emp}}(\boldsymbol{w}) = \frac{1}{2N} \|\Phi\boldsymbol{w} - \boldsymbol{y}\|^2$$

Solution: $\widehat{\boldsymbol{w}} = (\Phi^T\Phi)^{-1}\Phi^T\boldsymbol{y}$

# So, why not use many $\phi_i$'s?

Let us use

$$\mathcal{H} = \left\{ f(x) = \sum_{j=0}^{99} w_j x^j : \boldsymbol{w} \in \mathbb{R}^{100} \right\}$$

to fit our linear and quadratic examples



**This is called over-fitting. Why are these bad?**

# Learning ≠ Optimization/ERM

We want to do well on **unseen** data! In other words, our model must **generalize.**

Empirical risk minimization

Generalization Gap

$$\min_{f \in \mathcal{H}} R_{\text{emp}}(f) = \frac{1}{N} \sum_{i=1}^{N} L(f(\boldsymbol{x}_i), f^*(\boldsymbol{x}_i))$$

$$\boldsymbol{x}_i \sim \mu$$

$$\hat{f}$$

$$\nVdash$$

Population risk minimization

$$\min_{f \in \mathcal{H}} R_{\text{pop}}(f) = \mathbb{E}_{\boldsymbol{x} \sim \mu}[L(f(\boldsymbol{x}), f^*(\boldsymbol{x}))]$$

$$\tilde{f}$$

What we **really** want to solve

# Hypothesis Space and Generalization



*Curve fitting methods and the message they send. https://xkcd.com/2048/*

# How do we estimate the population risk?

We do not know the sample distribution $\mu$, but we already have samples from it in our dataset $\mathcal{D}$.

Therefore, we can use a train-test split
$$\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{test}}$$

- The training set $\mathcal{D}_{\text{train}}$ is used for empirical risk minimization
- The testing set $\mathcal{D}_{\text{test}}$ is used as an ultimate evaluation of performance
- Important: our learning algorithm should never peek at $\mathcal{D}_{\text{test}}$!

# Generalization

The problem of generalization is a central one in machine learning.

In this course, we will study a variety of ways to improve generalization, including

- Choice of architecture
- Regularization
- Data augmentation
- Optimization methods

# Classification vs Regression

**Let us consider a $K$-class classification problem**



What should we
understand as a
linear model in
this case?

# The Argmax Way

**Linear models for regression**

$$y = w^T \phi(x)$$

**Linear models for classification**

$$y_i = w_i^T \phi(x) \quad i = 1, 2, \ldots, K$$
$$\text{Prediction: } \arg \max_{i=1,\ldots,K} y_i$$

**Is this the same as SVM (Binary classification by hyperplanes)?**

$$\text{Sign}(w^T x + b)$$

# Limitations of This Approach

**Let's write a linear $K$-class classifier in a more compact way**

- Define a matrix $W \in \mathbb{R}^{K \times m}$
- Define an argmax function $g: \mathbb{R}^K \to \{1, \dots, K\}$ by
$$g(\mathbf{z}) = i \ \text{ if } \ z_i > z_j \text{ for all } j \neq i$$
  To break ties, we can alternatively define
$$g(\mathbf{z}) = \arg\min_{i=1,\dots,K} \left\{ z_i : z_i = \max_{j=1,\dots,K} z_j \right\}$$
- A linear classification model is thus
$$f: \mathbb{R}^d \to \{1, 2, \dots, K\}$$
$$f(\mathbf{x}) = g\big(W\boldsymbol{\phi}(\mathbf{x})\big)$$

# Limitations of This Approach

We can train these networks by defining a loss, say square error and minimize

$$L(W) = \frac{1}{2}\left(y_i - g\big(W\boldsymbol{\phi}(\boldsymbol{x_i})\big)\right)^2$$

No exact solution, so we perform gradient descent

$$W_{t+1} = W_t - \eta\nabla_W L(W)$$

What's wrong?

# Alternatively: One-hot Encoding and Softmax Activation

We can use the one-hot encoding to represent each label $y_i = k$ (class $k$) as a vector in $\mathbb{R}^K$

$$\boldsymbol{y}_i = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^K, \qquad i = 1, \dots, N$$

$k^{\text{th}}$ Position

**Then, our classification model can return a row of probabilities of a sample belong to each class, e.g.**
$$f(\boldsymbol{x}_i) = (a_1, \dots, a_k, \dots, a_K) \in \mathbb{R}^K, \qquad i = 1, \dots, N$$

Probability of belong to class $k$

This can be done by using the **softmax** function
$s: \mathbb{R}^K \to \mathbb{R}^K$

$$s(\mathbf{z})_k = \frac{\exp(z_k)}{\sum_j \exp(z_j)}$$

What does it do?

This gives the following hypothesis space
$$\mathcal{H} = \{f(\mathbf{x}) = s(W\boldsymbol{\phi}(\mathbf{x})): W \in \mathbb{R}^{K \times m}\}$$

Notice that $f \in \mathcal{H}$ always outputs a vector which can be interpreted as **probabilities** over $K$ classes

# Loss Functions for Classification

By right, classification should be done with zero-one loss: 1 if incorrect, 0 if correct. This is just accuracy!

Why do we not use this?

Instead, we consider loss surrogate
$$L: \mathbb{R}^K \times \mathbb{R}^K \to \mathbb{R}_+$$
where
- If $y = y'$, then $L(y, y')$=0
- If $y \approx y'$, then $L(y, y')$ is small
- $L$ is differentiable, with non (almost) everywhere zero gradients

# Examples

We can still use the square loss
$$L(\boldsymbol{y}, \boldsymbol{y}') = \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{y}'\|^2$$
What could be some issues with this?

The cross-entropy loss is defined as
$$L(\boldsymbol{y}, \boldsymbol{y}') = -\sum_{k=1}^{K} y'_k \log y_k$$
Theses are not the only choices! Any differentiable "distance" between two probability distributions suffices.

# Summary

- The T, E, P formulation of machine learning
- Linear regression as an illustration
    - Hypothesis space and Capacity
    - Learning algorithm
    - Underfitting and Overfitting
    - Performance metrics and test set
- Classification: use one-hot encoding, compare probabilities

# Useful tools for programming

Version control with **Git**

- https://www.freecodecamp.org/news/what-is-git-and-how-to-use-it-c341b049ae61/

Interactive python with **Jupyter notebooks**

- https://www.datacamp.com/community/tutorials/tutorial-jupyter-notebook

Data visualization using **Seaborn** and **Pandas**

- https://jakevdp.github.io/PythonDataScienceHandbook/04.14-visualization-with-seaborn.html

# Demo (Linear Regression)

**See uploaded under folder "Demos"**