

The background features a complex network graph with numerous nodes and edges. The nodes are represented by small circles in various colors including teal, orange, purple, and black. The edges are thin, dark grey lines connecting these nodes. A large, white, rounded rectangle with a slightly irregular border is superimposed on the left side of the image, framing the text.

# 深度学习

---

Lecture1 Introduction  
YMSC, Pang Tongyao

# About This Course

---

- It is
  - An introduction to deep learning methods, from basics to advanced topics.
  - Focused on algorithms and applications, emphasizing the development of strong coding skills. (Math is still important!)
- It isn't
  - Focused on rigorous mathematical definitions and derivations.
  - A comprehensive overview of all topics within deep learning.

# About This Course

---

- Grading Policy
  - Assignments: 3\*25%
  - Final Projects: 25%
  - Late policy: 10% off per day late; not accepted after 10 days
- Office hour
  - by appointment
  - email: [typang@tsinghua.edu.cn](mailto:typang@tsinghua.edu.cn)
  - office: 双清综合楼B533

# Reference

---

- Books

- Bishop, Christopher M., and Nasser M. Nasrabadi. **Pattern recognition and machine learning**, Springer, 2006.
- Bengio, Yoshua, Ian Goodfellow, and Aaron Courville. **Deep learning**, MIT press, 2017.
- Zhang, Aston and Lipton, Zachary C. and Li, Mu and Smola, Alexander J., **Dive into Deep Learning**, Cambridge University Press, 2023

- Related Courses:

- 6.8300/6.8301: Advances in Computer Vision, MIT
- 6.S978: Deep Generative Models, MIT
- CS231n: Deep Learning for Computer Vision, Stanford
- Reinforcement Learning, by David Silver, UCL

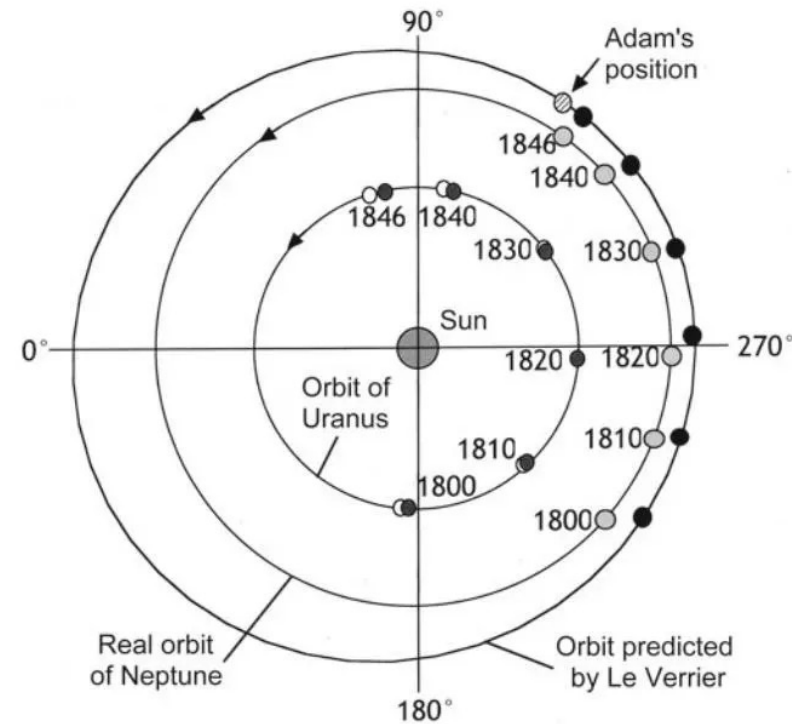
# What is Learning?



# Motivating Examples

Year	Housing Price
2021	2
2022	3
2023	4
2024	?

Predict the Housing Price of a City



## The Discovery of Neptune

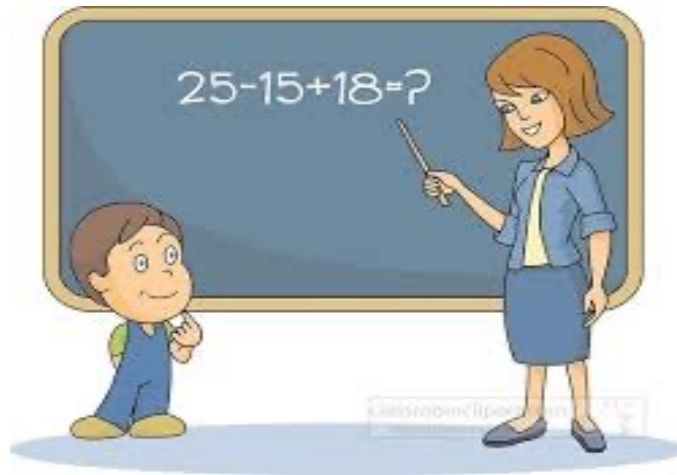
Astronomers noticed Uranus' motion differed from predictions and used its trajectory to predict the location and mass of an unknown celestial body.

# Motivating Examples

---



A doggie is learning



A kid is learning



Also learning.....

# Machine/Deep Learning

---

- **Machine learning** is the study of algorithms that can learn from experience.
- Experience: typically in the form of observational data or interactions with an environment.
- As a machine learning algorithm accumulates more experience, its performance improves.

## Artificial Intelligence



Any technique that enables computers to mimic human intelligence. It includes *machine learning*

## Machine Learning



A subset of AI that includes techniques that enable machines to improve at tasks with experience. It includes *deep learning*

## Deep Learning



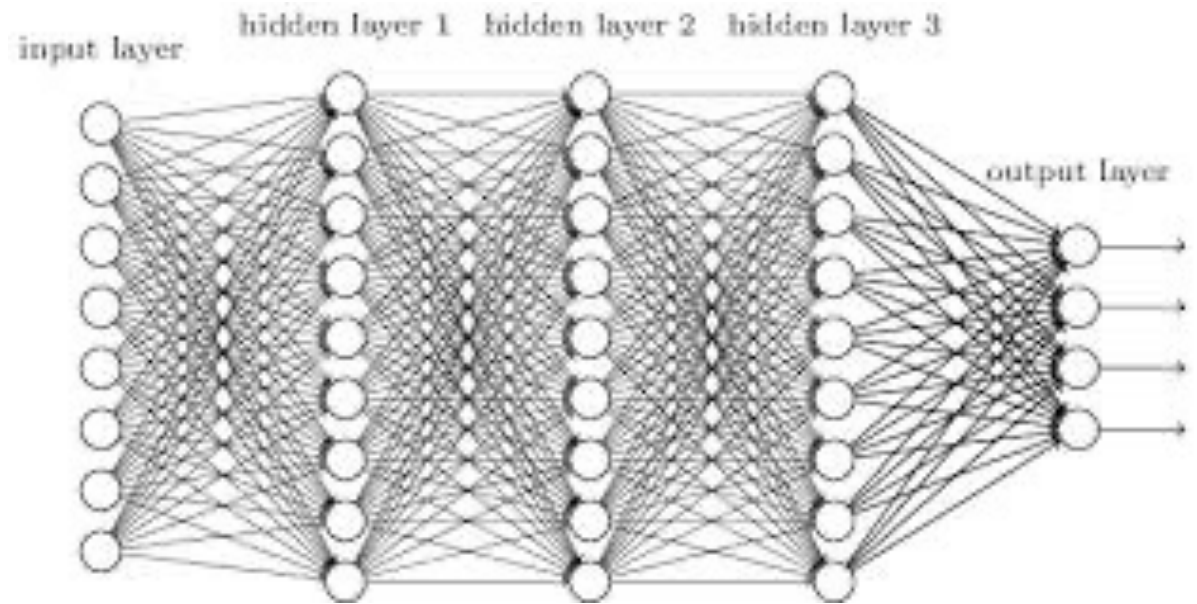
A subset of machine learning based on neural networks that permit a machine to train itself to perform a task.



# Deep Learning Methods

---

- Supervised Learning
- Unsupervised Learning
  - Self-supervised Learning
  - **Generative Models**
- Reinforcement Learning
- Meta/Transfer Learning
- .....



Deep Neural Networks

# Applications

---

*"Swap sunflowers with roses"*



*"Add fireworks to the sky"*



*"Replace the fruits with cake"*



*"What would it look like if it were snowing?"*



*"Turn it into a still from a western"*



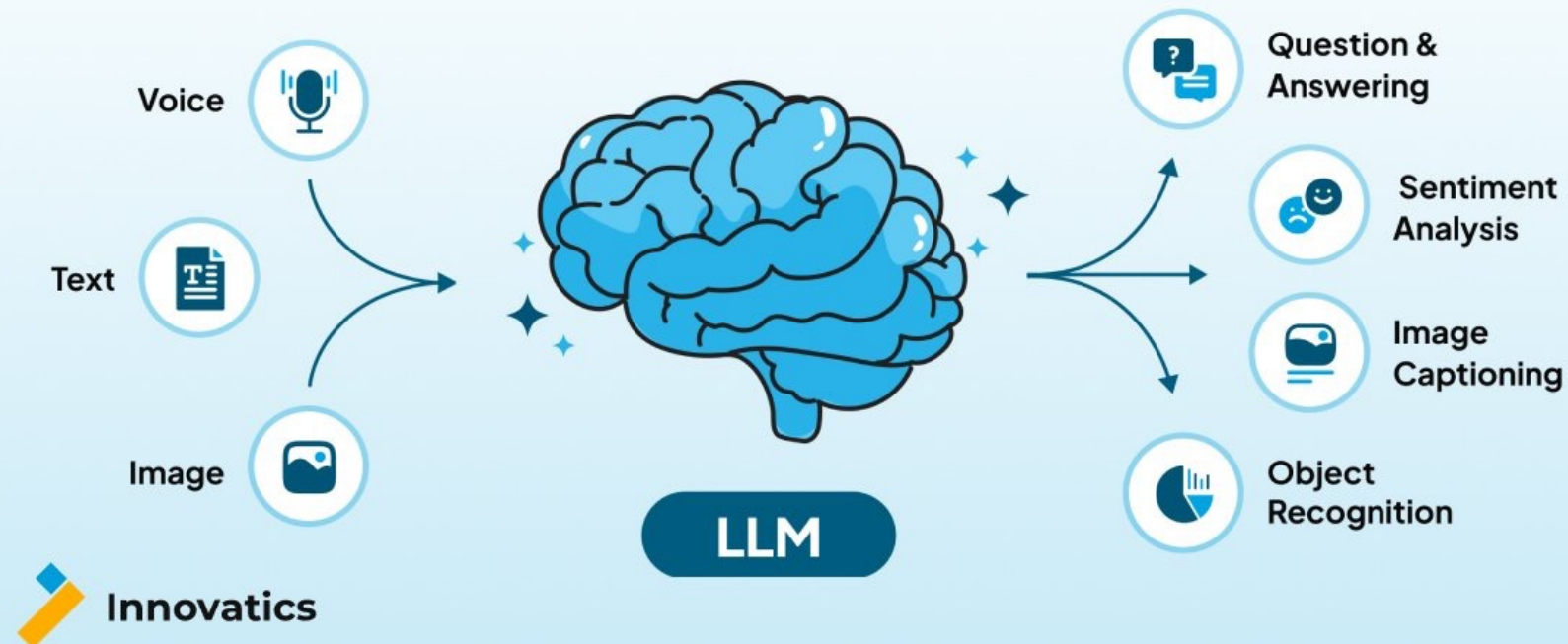
*"Make his jacket out of leather"*



Text-to-image Generation

# Applications

## Exploring Large Language Models (LLMs)



Gemini

ChatGPT

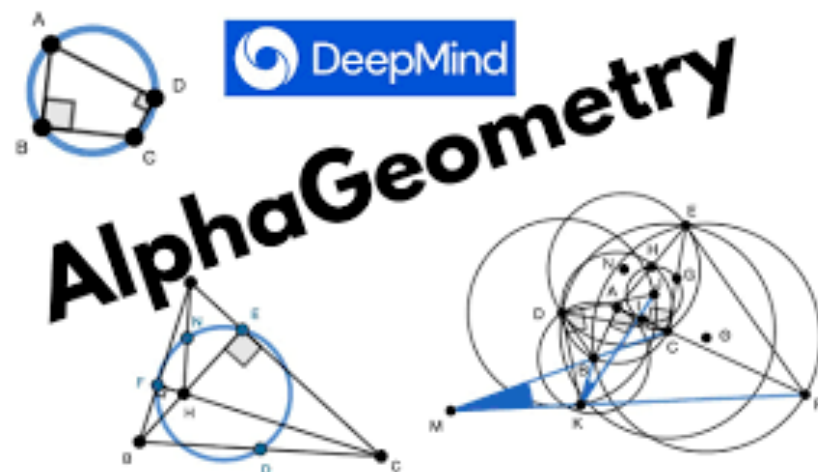
Claude

deepseek



# Applications

---



# Applications

---



Embodied Intelligence



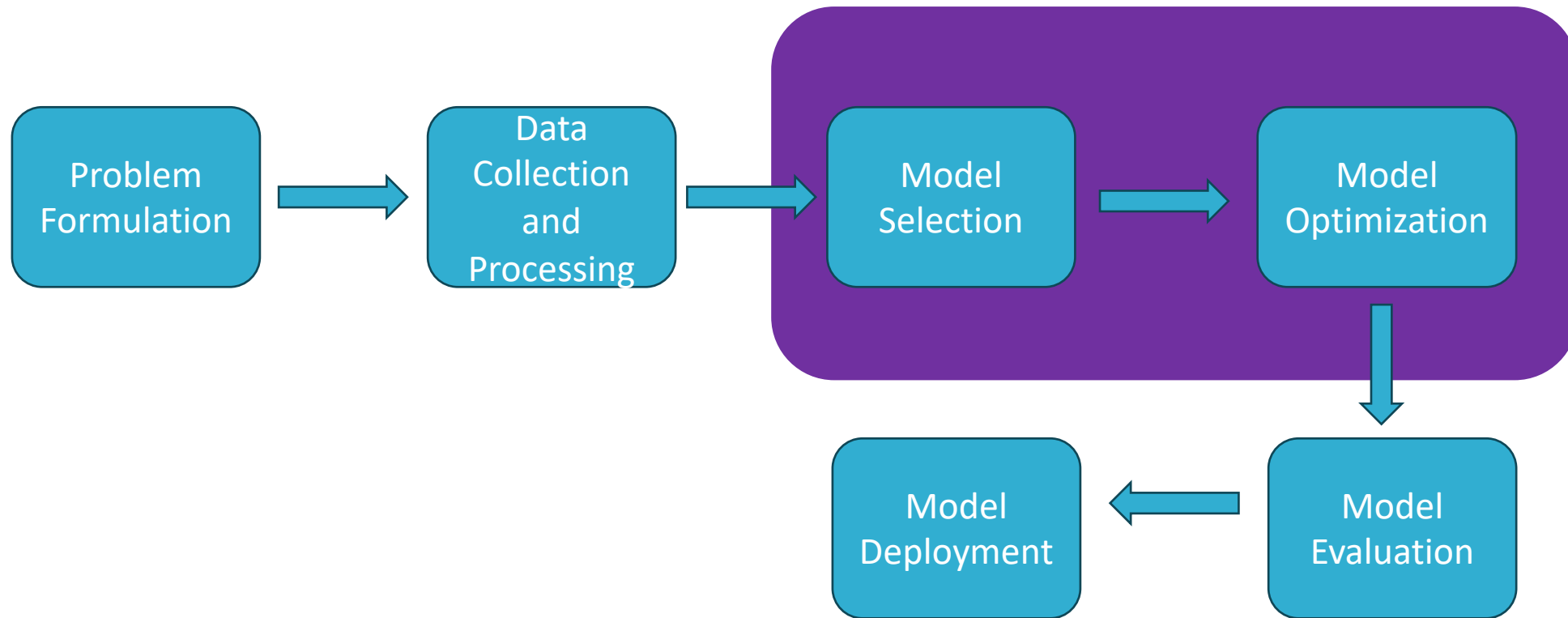
Self-Driving

Start from the very  
beginning



# Typical Process

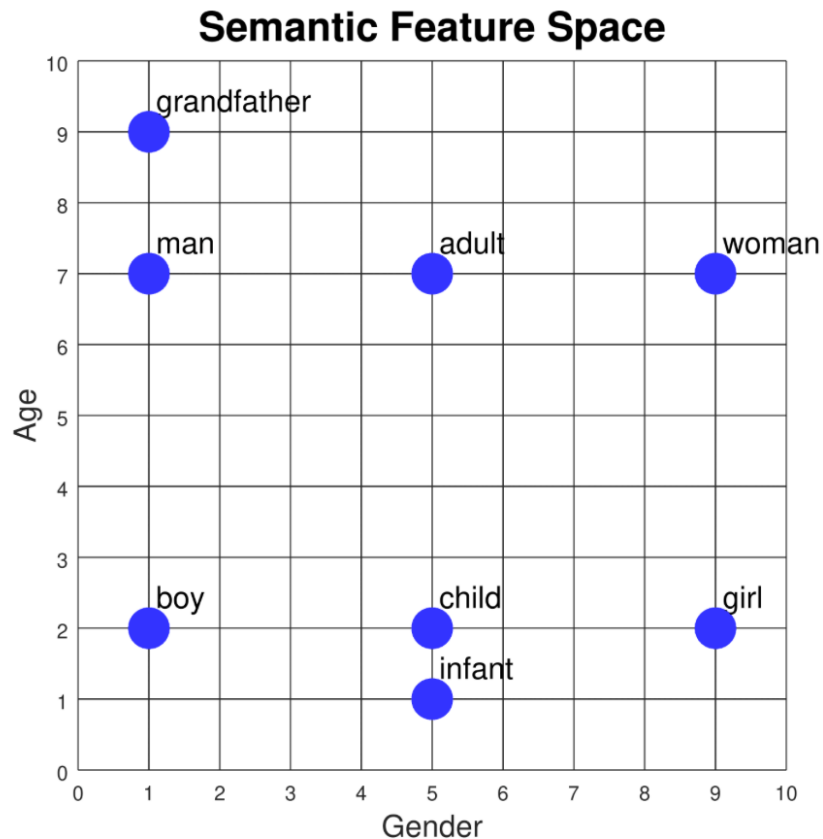
---



[illegible]



# Data Format



Word Coordinates		
	Gender	Age
grandfather	[ 1,	9 ]
man	[ 1,	7 ]
adult	[ 5,	7 ]
woman	[ 9,	7 ]
boy	[ 1,	2 ]
child	[ 5,	2 ]
girl	[ 9,	2 ]
infant	[ 5,	1 ]

--Dave Touretzky

# Data Format

---

## *One-hot Encoding*

apple = [ 1 0 0 0 0 ..... ]

bag = [ 0 1 0 0 0 ..... ]

cat = [ 0 0 1 0 0 ..... ]

dog = [ 0 0 0 1 0 ..... ]

elephant = [ 0 0 0 0 1 ..... ]

# Data Format

---

- In deep learning algorithms, most data is in the form of tensors, which can be better utilized for accelerated computation on GPUs.
- Tensor data can still have various structures, such as time series data, graph data, etc.
- Dataset Decomposition

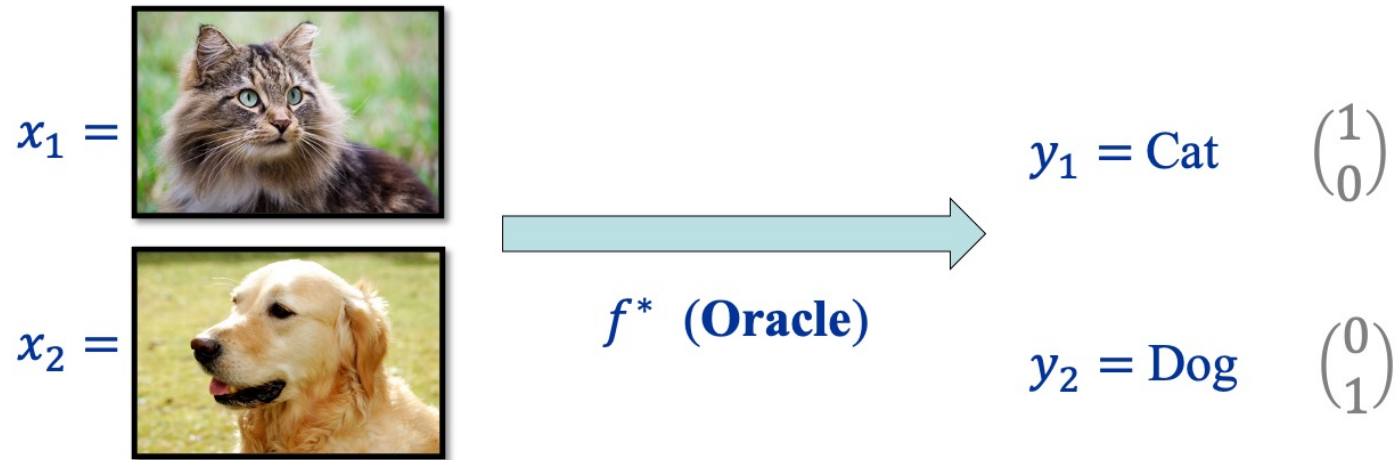
$$\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{test} \cup \mathcal{D}_{validation}$$

$\mathcal{D}_{train}$ : for training;  $\mathcal{D}_{test}$ : for evaluation;  $\mathcal{D}_{validation}$ : for selection

# Supervised Learning

---

- Training data with labels:  $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$ , input:  $x_i$ , output/label:  $y_i$ ,
- Goal: learn the mapping  $x_i \rightarrow y_i$



- The oracle  $f^*$  is unknown to us, except through the dataset

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

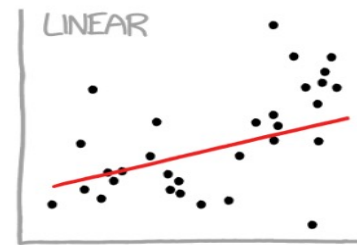
# Hypothesis Space

## Supervised Learning:

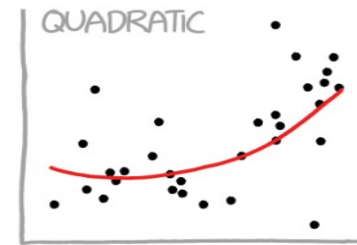
1. Define a **hypothesis space**  $\mathcal{H}$  consisting of a set of candidate functions, e.g.

$$\mathcal{H} = \{f: f(x) = w_0 + w_1x\}$$

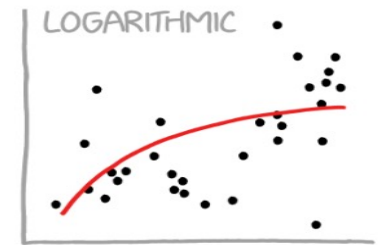
2. Find the “**best**” function  $\hat{f}$  in  $\mathcal{H}$  that **approximates**  $f^*$



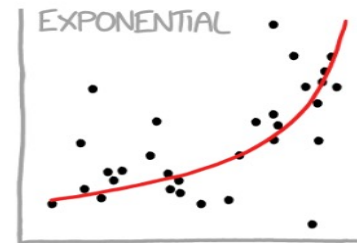
“HEY, I DID A REGRESSION!”



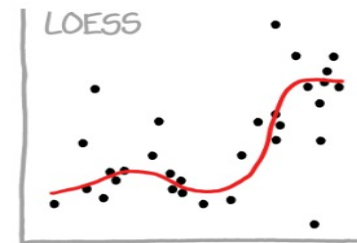
“I WANTED A CURVED LINE, SO I MADE ONE WITH MATH.”



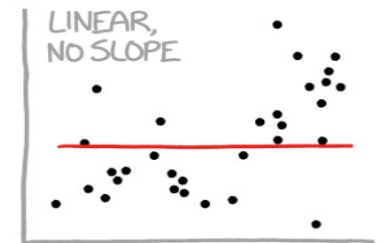
“LOOK, IT’S TAPERING OFF!”



“LOOK, IT’S GROWING UNCONTROLLABLY!”



“I’M SOPHISTICATED, NOT LIKE THOSE BUMBLING POLYNOMIAL PEOPLE.”



“I’M MAKING A SCATTER PLOT BUT I DON’T WANT TO.”

What you get depends on  $\mathcal{H}$ !

# Optimization

---

- Define a **loss function** to measure the distance between the model solution and the oracle solution. Then, we can find the best approximation by an optimization problem

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

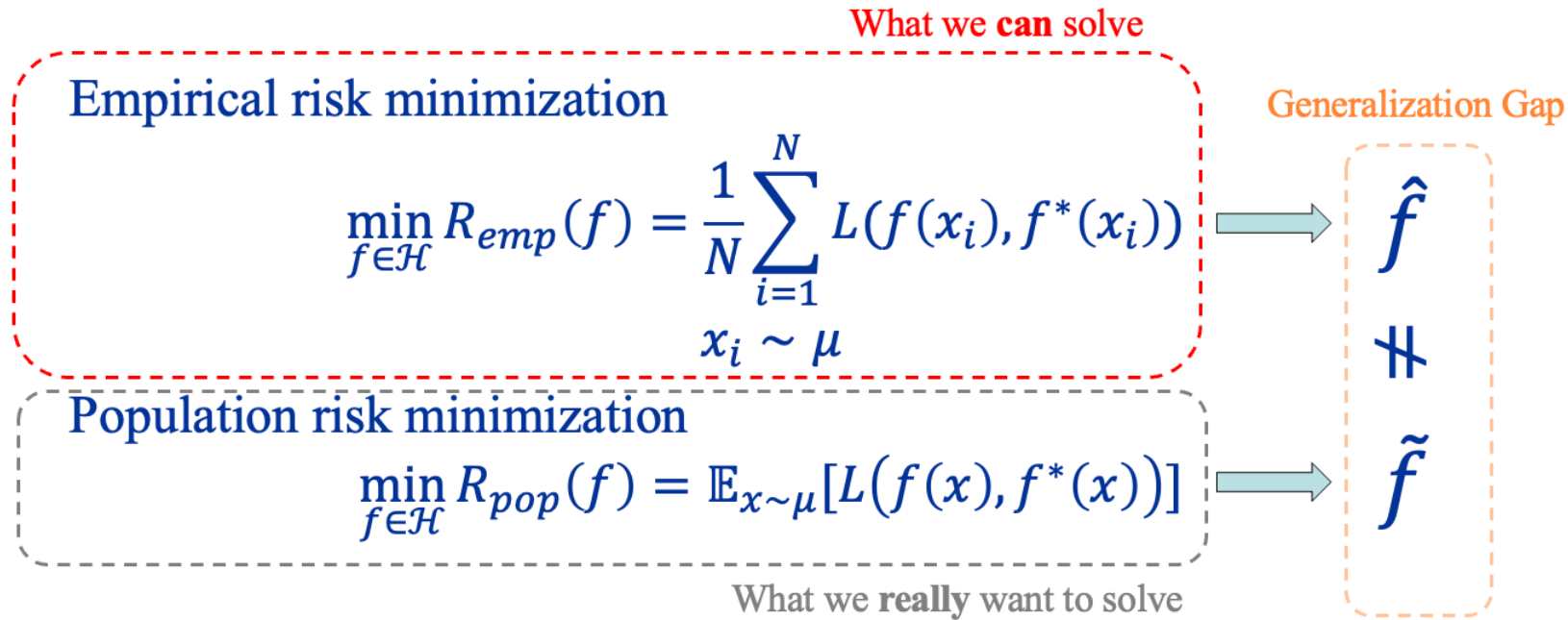
↑  
The empirical loss

- Optimization algorithms to be explored in the following lectures.....

# Generalization

---

- **Generalization ability:** to do well on new, unseen data.

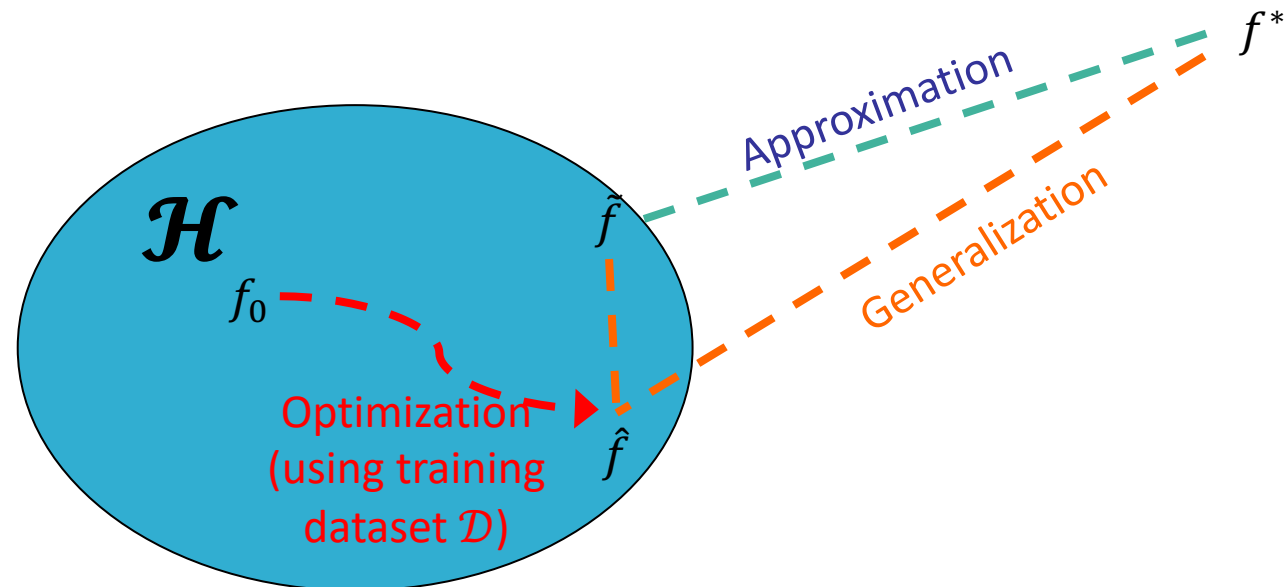


Generalization error depends on the training dataset, hypothesis spaces, and optimization algorithms etc.

# Paradigms of Supervised Learning

---

- **Approximation:** the distance between the hypothesis space and the oracle solution.
- **Optimization:** Seeks the best solution within the hypothesis space based on training data.
- **Generalization:** Examines the difference between the optimized solution and the true solution on unseen data.





# Linear Regression



# Linear Regression

---

- Linear Model with Continuous Output
- Step 1: Define the hypothesis space as the set of linear models.

$$\mathcal{H} = \{f: f(x) = w_0 + w_1x, w_0 \in \mathbb{R}, w_1 \in \mathbb{R}\}$$

- Step 2: Find the best approximation.

First, define the loss function, for example:

$$L(f_1, f_2) = \int p(x) \|f_1 - f_2\|_2^2 d\mu$$

Minimize the empirical loss over the training dataset  $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$

$$\min_{f \in \mathcal{H}} R_{emp}(f) = \min_{w_0, w_1} \frac{1}{2N} \sum_{i=1}^N (w_0 + w_1 x_i - y_i)^2$$

Solution:

$$\frac{\partial R_{emp}}{\partial w_0}(\hat{w}_0, \hat{w}_1) = 0 \text{ and } \frac{\partial R_{emp}}{\partial w_1}(\hat{w}_0, \hat{w}_1) = 0$$
$$\hat{w}_0 = \bar{y} - \hat{w}_1 \bar{x} \quad \hat{w}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} \quad \bar{x} = \frac{1}{N} \sum_i x_i \quad \bar{y} = \frac{1}{N} \sum_i y_i$$

Test process:  $y_{test} = \hat{w}_0 + \hat{w}_1 x_{test}$

# Linear Basis Models

---

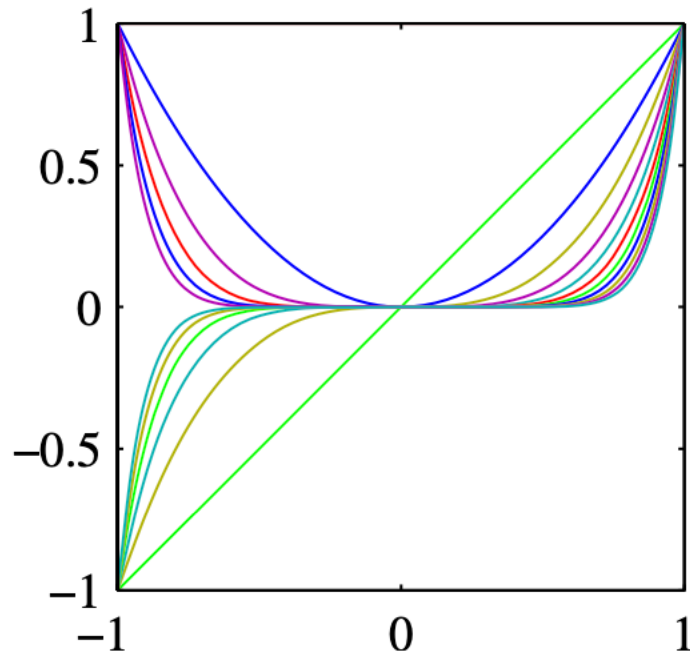
- Q: Can linear models only represent linear relationships between  $x$  and  $y$ ?
- **No!** Introduce basis functions or feature mappings can extend its ability to model more complex relationships.

$$\mathcal{H}_M = \left\{ f: f(x) = \sum_{j=0}^{M-1} w_j \phi_j(x) \right\}$$

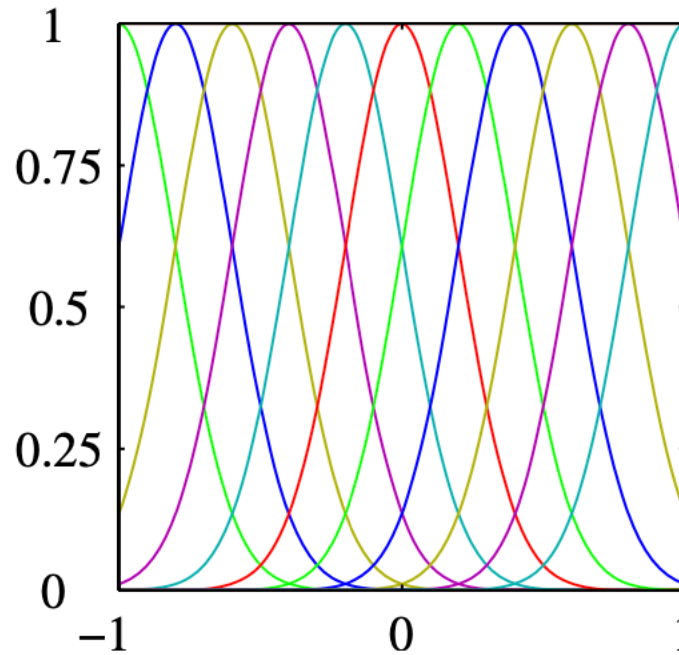
$\mathcal{H} = \{f: f(x) = w_0 + w_1 x, w_0 \in \mathbb{R}, w_1 \in \mathbb{R}\}$  corresponds to  
 $d = 1, M = 2, \phi_0(x) = 1, \phi_1(x) = x$

# Example of Basis Functions

---



Polynomials:  $\phi_j(x) = x^j$



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

-- Bishop

# Linear Basis Models

---

- Optimization:

$$\min_{f \in \mathcal{H}_M} R_{emp}(f) = \min_{w \in \mathbb{R}^M} R_{emp}(w)$$

$$= \min_{w \in \mathbb{R}^M} \frac{1}{2N} \sum_{i=1}^N (f(x_i) - y_i)^2$$

$$= \min_{w \in \mathbb{R}^M} \frac{1}{2N} \sum_{i=1}^N \left( \sum_{j=0}^{M-1} w_j \phi_j(x_i) - y_i \right)^2$$

where

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \cdots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \cdots & \phi_{M-1}(x_2) \\ \vdots & \ddots & \vdots \\ \phi_0(x_N) & \cdots & \phi_{M-1}(x_N) \end{pmatrix}$$

$$w = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_{M-1} \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

# Linear Regression

---

- To solve  $\min_{w \in \mathbb{R}^M} \frac{1}{2N} \|\Phi w - y\|^2$ ,

we can do this by setting  $\nabla R_{emp}(\hat{w}) = 0 \Leftrightarrow \Phi^T(\Phi \hat{w} - y) = 0$

Solve the linear equation systems:

- When  $\Phi^T \Phi$  is invertible, there is a unique solution  $\hat{w} = (\Phi^T \Phi)^{-1} \Phi^T y$
- When  $\Phi^T \Phi$  is not invertible, there are many solutions

$$\hat{w}(u) = \Phi^\dagger y + (I - \Phi^\dagger \Phi)u, \forall u \in \mathbb{R}^M, \Phi^\dagger = (\Phi^T \Phi)^{-1} \Phi^T$$

Q: how to select the desirable one?

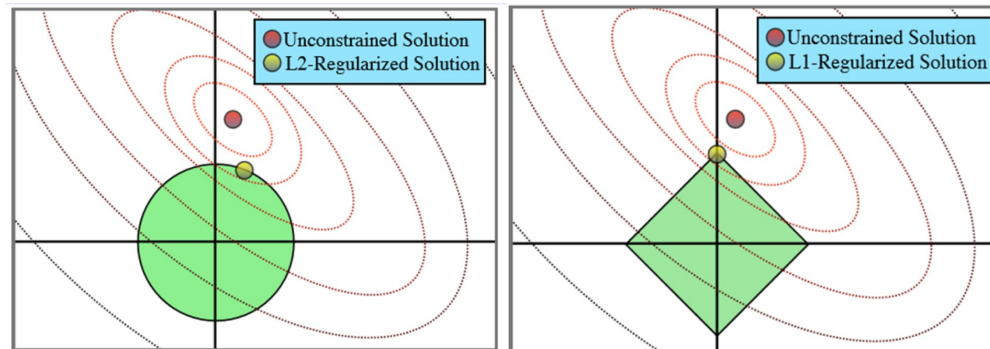
# Linear Regression

---

- Add regularization

$$\min_{w \in \mathbb{R}^M} \frac{1}{2N} \|\Phi w - y\|^2 + \underbrace{\lambda C(w)}_{\text{the regularization term}}$$

- $\ell_2$  regularization (ridge regression):  $C(w) = \|w\|_2^2$ .
- $\ell_1$  regularization (**lasso**):  $C(w) = \|w\|_1$ , pursuit sparsity
- .....



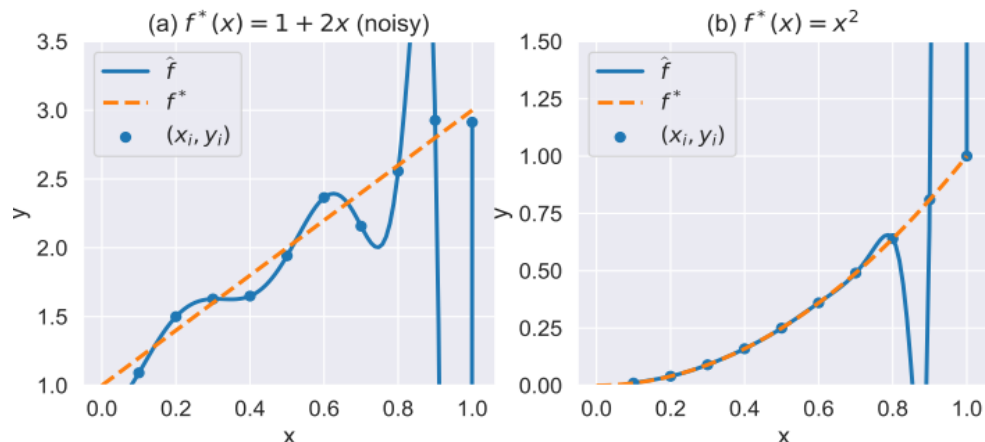


# Regularization and generalization

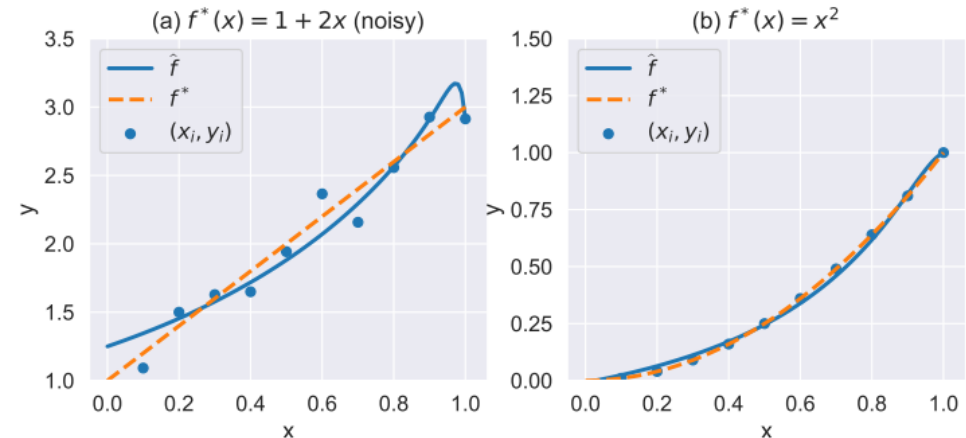
- We apply  $\ell^2$  regularization on the over-fitting examples
- Consider

$$\mathcal{H}_M = \{f: f(x) = \sum_{j=0}^{99} w_j x^j\} \text{ so } M = 100, \text{ but } N = 10$$

Without regularization



With regularization



# Linear Regression

---

- Ridge regression

$$\min_{w \in \mathbb{R}^M} \frac{1}{2N} \|\Phi w - y\|^2 + \lambda \|w\|_2^2$$

$$\hat{w} = (\Phi^T \Phi + 2N\lambda I)^{-1} \Phi^T y$$

- Lasso: no analytical solution

$$\min_{w \in \mathbb{R}^M} \frac{1}{2N} \|\Phi w - y\|^2 + \lambda \|w\|_1$$

# Optimization Algorithms

---

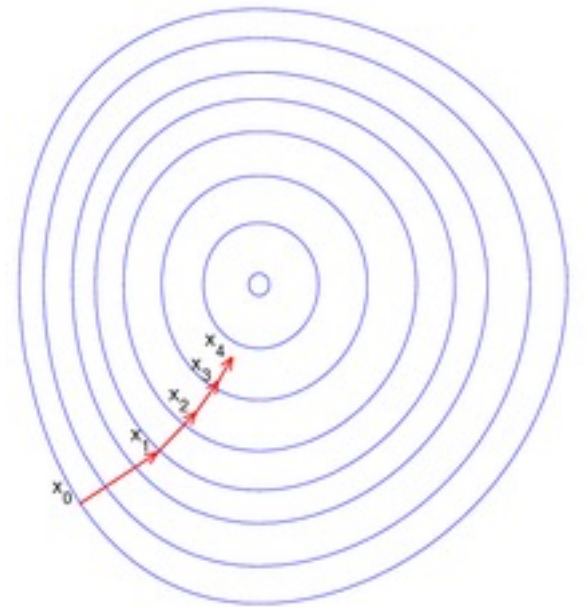
- A necessary first-order optimality condition for  $\min f(x)$

$$\nabla f(x^*) = 0$$

- Solve it iteratively:

$$\text{Gradient Descent : } x_{n+1} = x_n - \gamma \nabla f(x_n)$$

- Provided  $\gamma < \|\nabla^2 f\|$ , it can be shown that  $\|\nabla f(x_n)\| \rightarrow 0$



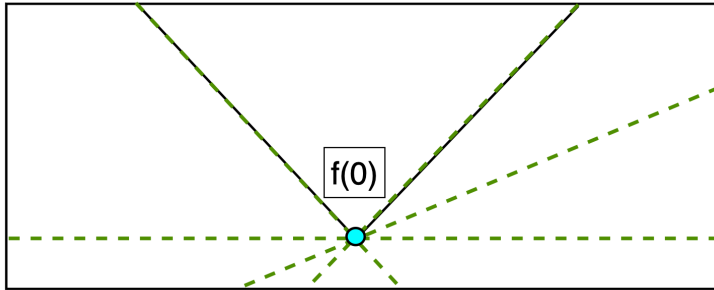
# Optimization Algorithms

---

- $\|\cdot\|_1$  is not differentiable at 0. Consider subgradient

$$\partial f(x_0) = \{g | f(x) \geq f(x_0) + g^T(x - x_0)\}$$

$$\partial|x| = \begin{cases} 1, & x > 0 \\ -1, & x < 0 \\ [-1,1], & x = 0 \end{cases}$$



- Subgradient methods:  $w_{k+1} = w_k - \gamma_k g_k$ ,  $g_k \in \partial f(w_k)$
- Due to the slow convergence rate of subgradient methods, many advanced algorithms have been proposed to solve Lasso, e.g. ADMM, proximal gradient descent.

# Classification

---



This image by Nikita is  
licensed under [CC-BY 2.0](#)

(assume given a set of possible labels)  
{dog, cat, truck, plane, ...}



cat

(Discrete output)

# Logistic Regression

---

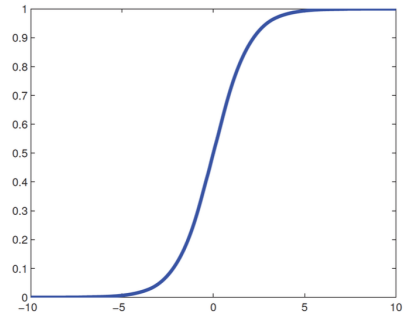
- Binary classification problem:

one-hot encoding for the output  $\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}$

This two-dimensional vector can be understood as the probability for each class and can take continuous values.

- A linear hypothesis space:  $\{u(x): u = w^T x, x \in \mathbb{R}^n, w \in \mathbb{R}^n\}$ .
- Softmax: Map the extracted feature  $u$  to the space of one-hot codes

$$\mu = \frac{1}{1 + e^{-u}}, 1 - \mu = \frac{e^{-u}}{1 + e^{-u}}$$



# Logistic Regression

---

- Loss function
  - $\ell_2$ -loss
  - Maximum likelihood (KL divergence)

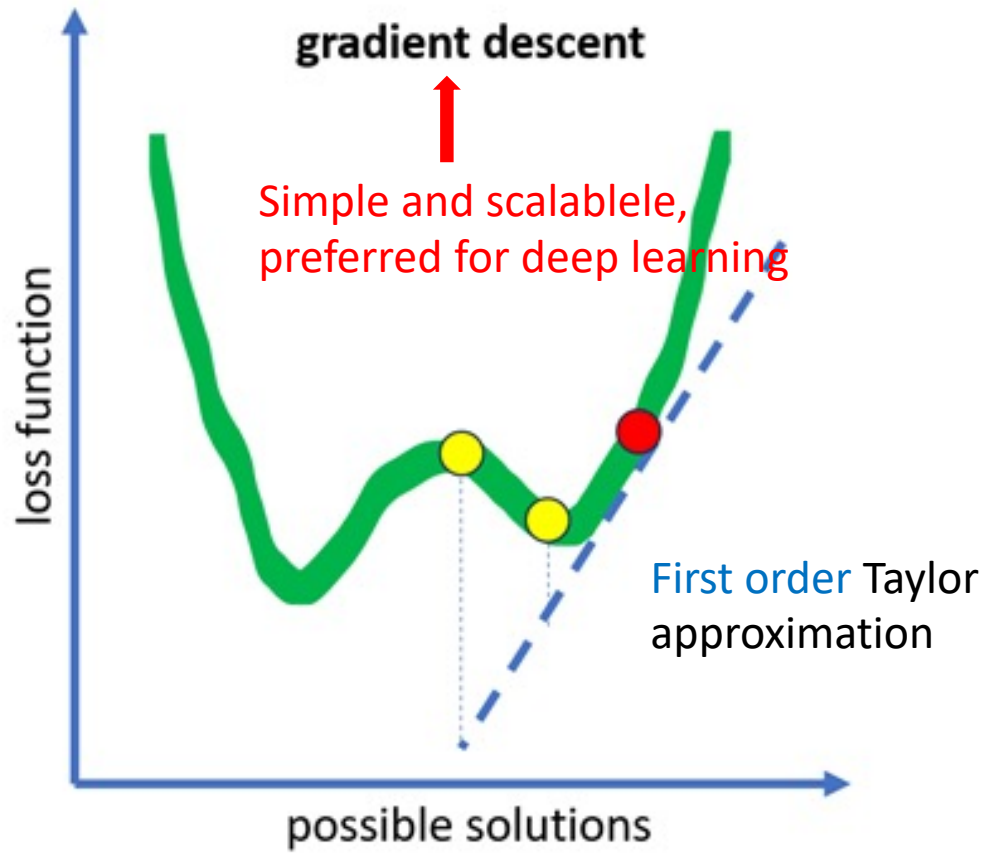
$$\max \prod_{i=1}^N p(y_i|x_i, w) \Leftrightarrow \min - \sum_{i=1}^N \log p(y_i|x_i, w) \text{ (cross entropy)}$$

$$- \sum_{i=1}^N \log p(y_i|x_i, w) = \sum_{i=1}^N -y_i \log \mu_i - (1 - y_i) \log(1 - \mu_i), \mu_i = \frac{1}{1 + e^{-w^T x}}$$

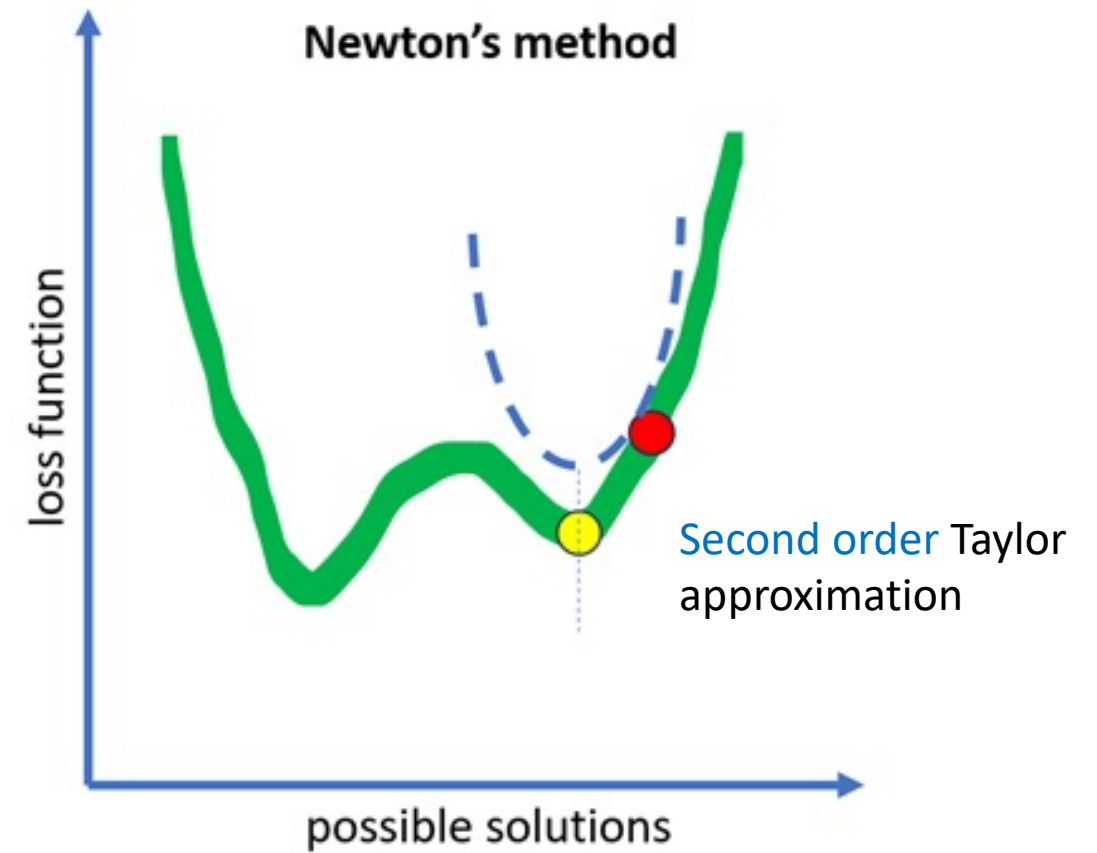
- Optimization: smooth and convex loss function; gradient descent, newton methods...

# Optimization

---



lower convergence, low memory  
requirements and computational cost



faster convergence, high memory  
requirements and computational cost



# Summary

---

- Machine Learning: learning from experience (or data).
- Supervised learning is to learn the mapping from  $x$  to  $y$  given paired training dataset  $\{(x_i, y_i)\}$ . Key steps include
  1. Define the hypothesis space
  2. Define the empirical loss
  3. Find the best approximation under the given empirical loss