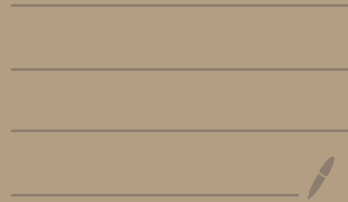


Lecture 1

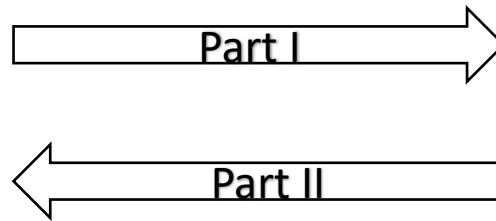
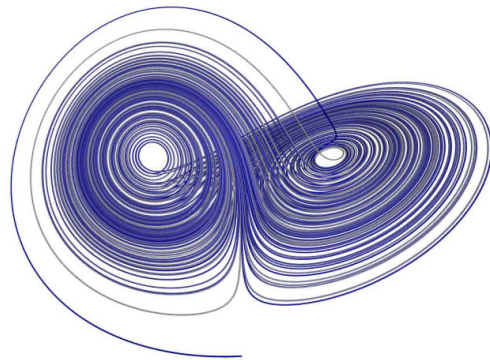


Dynamical Systems and Machine learning

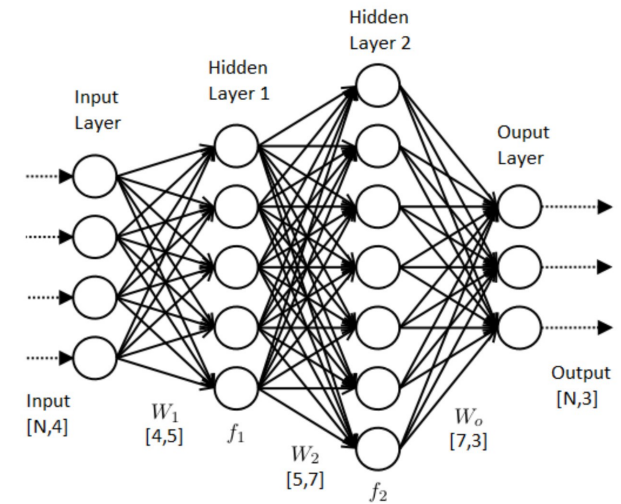
Qianxiao Li
(NUS, A*STAR)

Overview

Dynamical
Systems



Machine
Learning



Syllabus

- Introduction
- Part I: Dynamical Systems Approach to Deep Learning
 - Optimal control theory
 - Optimal control formulation of deep learning
 - Dynamics/Control-inspired training algorithms
 - Dynamics/Control-inspired network architectures
 - Mathematical results
- Part II: Data-Driven Dynamical Systems
 - PCA/SVD based model reduction methods (DMD, POD)
 - Koopman operator methods
 - Data driven control

Other Matters

- Class participation is encouraged!
 - Use the chat
 - TA will moderate
- You can email me at qianxiao@nus.edu.sg
 - Questions about the material
 - Mistakes in lecture notes
 - Suggestions on lecture pace and style
 - Research opportunities

Supervised Learning

Basic formulation

- Data : $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$

- $x_i \in \mathbb{R}^d$ inputs
- $y_i \in \mathbb{R}^m$ outputs / labels
- $N \geq 1$ size of data

- Data distribution : $x_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{I}$.

- Target function:

- Deterministic : $y_i = F^*(x_i)$
- Stochastic : $y_i \sim P^*(\cdot | x_i)$ e.g. $y_i = F^*(x_i) + \epsilon_i$
- General : $(x_i, y_i) \sim \mathcal{I}$

$\epsilon_i \sim N(0, I)$

- Output space

- Regression $y_i \in \mathbb{R}$
 - Classification $y_i \in \{1, 2, \dots, C\}$
- $\xrightarrow{\text{embed}} \mathbb{R}^C \quad \mapsto \begin{pmatrix} \vdots \\ \text{One-hot} \end{pmatrix}$

Goal of supervised learning

Given \mathcal{D} , find $\tilde{F} \approx F^*$

Learn with risk minimization

- Define a hypothesis space

$$\mathcal{H} = \{ F \mid F: \mathbb{R}^d \rightarrow \mathbb{R}^m \}$$

goal is to find $F \in \mathcal{H}$ s.t. $\|F^* - F\|$ is small.

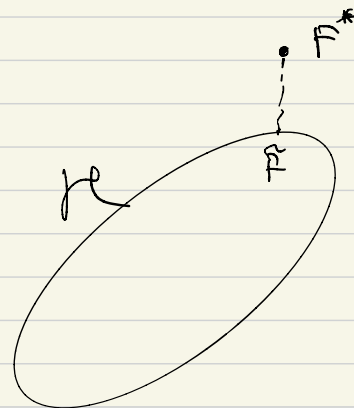
- Define a loss function
 $\Phi: \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$.

- Define the empirical risk

$$R_{\text{emp}}(F) = \frac{1}{N} \sum_{i=1}^N \Phi(F(x_i), y_i)$$

- Empirical risk minimization

$$\hat{F} \leftarrow \arg \min_{F \in \mathcal{H}} R_{\text{emp}}(F)$$



$$y_i = F^*(x_i)$$

- loss function
 - mean-square loss

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

- zero-one loss

$$\Phi(y, y') = \mathbb{1}_{y \neq y'} = \begin{cases} 1 & y \neq y' \\ 0 & y = y' \end{cases}$$

- Population/Expected Risk Minimization (PRM)

$\text{argmin}_{F \in \mathcal{H}} \mathbb{E}_{x \sim \mu} \Phi(F(x), \underbrace{F^*(x)}_y)$ (deterministic)

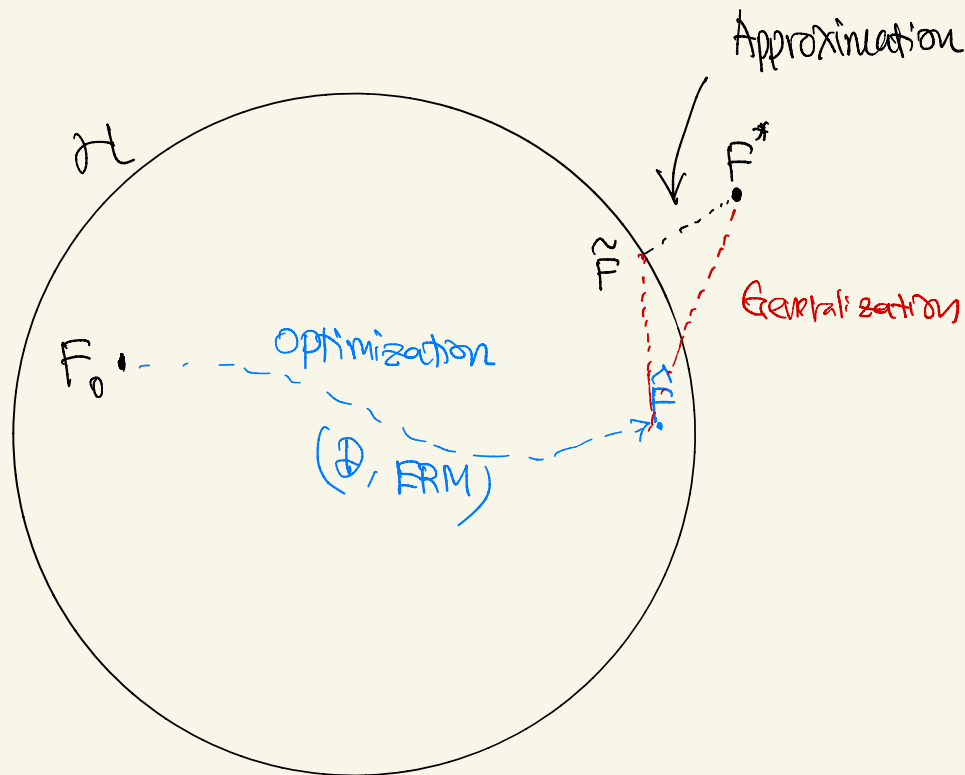
\leftarrow
 $\text{argmin}_{F \in \mathcal{H}} \mathbb{E}_{(x, y) \sim \mu} \Phi(F(x), y)$ (general)

If μ is empirical measure

$$\mu = \frac{1}{n} \sum_{i=1}^n \delta_{(x_i, y_i)}$$

ERM.

Three paradigms of Supervised Learning



Example (Linear Models)

Regression $y \in \mathbb{R}$

Linear model hypothesis space

$$\mathcal{H} = \left\{ F : F(x) = \sum_{j=0}^M \omega_j \phi_j(x), \omega_j \in \mathbb{R}, j = 0, 1, \dots, M-1 \right\}$$

↑ Basis fns, feature maps.
 $\phi_j : \mathbb{R}^d \rightarrow \mathbb{R}$.

Examples of Basis Function ($d=1$)

- $\phi_0(x) = 1, \phi_1(x) = x \Rightarrow F(x) = w_0 + w_1 x$.
- Polynomial : $\phi_j(x) = x^j$
- RBF / Gaussian : $\phi_j(x) = \exp\left(-\frac{1}{2\sigma_j^2} (x - m_j)^2\right)$

- Empirical Risk Minimization

$$\begin{aligned}
 R_{\text{emp}}(w_0, \dots, w_{M-1}) &= \frac{1}{2N} \sum_{i=1}^N \Phi(Ax_i, y_i) \\
 &= \frac{1}{2N} \sum_{i=1}^N \left(\underbrace{\sum_{j=0}^{M-1} w_j \phi_j(x_i)}_{F(x)} - y_i \right)^2
 \end{aligned}$$

Matrix form

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

$w = \begin{pmatrix} w_0 \\ \vdots \\ w_{M-1} \end{pmatrix} \in \mathbb{R}^M$, $y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N$, $(\Phi)_{ij} = \phi_j(x_i) \in \mathbb{R}^{N \times M}$.

$$\min_w R_{\text{emp}}(w) \rightarrow \hat{w}$$

Ordinary Least Squares Formula:

Suppose $\Phi^T \Phi$ is invertible, then

$$\hat{w} = (\Phi^T \Phi)^{-1} \Phi^T y.$$

Proof

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

$$\nabla_w \text{Remp}(w) = \frac{1}{N} \Phi^T (\Phi w - y)$$

$$\text{Set } \nabla_w \text{Remp}(\hat{w}) = 0$$

$$\Rightarrow (\Phi^T \Phi) \hat{w} = \Phi^T y.$$

$$\hat{w} = (\Phi^T \Phi)^{-1} \Phi^T y. \quad \square$$

If $M > N$, $\Phi^T \Phi$ not invertible. $(\hat{w} = \Phi^T y)$ $\Phi^T \in \mathbb{R}^{M \times N}$.

Consider **regularized** problem

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

($\lambda > 0$)

$$F^*: \mathbb{R}^d \rightarrow \mathbb{R}^m.$$

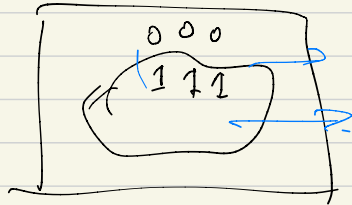
$G \rightarrow$ group of transformations.

$g(x)$ group action.

$$F^*(g(x)) = F^*(x) \quad (\text{invariant})$$

(equivariant) $F^*(g(x)) = \hat{g}(F^*(x))$

\hat{g} action on \mathbb{R}^m .



$$R_{\text{pop}}(\hat{F}) - R_{\text{emp}}(\hat{F}) \leq \frac{C(\hat{F})}{N^\alpha}$$

Neural Network Hypothesis Space

$$\mathcal{H}_M = \{ F : F(x) = \sum_{j=1}^M v_j \underbrace{\sigma(w_j^T x + b_j)}_{\phi_j(x)}, w_j \in \mathbb{R}^d, v_j, b_j \in \mathbb{R}, j=1, \dots, M \}$$

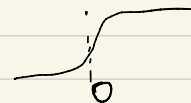
- "adaptive basis model"
- σ is the activation function

$$\sigma: \mathbb{R} \rightarrow \mathbb{R}$$

Examples:

- ReLU (Rectified Linear Unit) $\sigma(z) = \max(0, z)$
- sigmoid $\sigma(z) = \frac{1}{1 + e^{-z}}$
- tanh

leaky-ReLU



- Matrix form

$$F(x) = v^T \sigma(Wx + b)$$

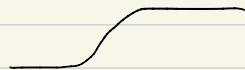
$\mathbb{R}^{M \times d}$

$$[\sigma(z)]_i = \sigma(z_i)$$

Universal Approximation Theorem

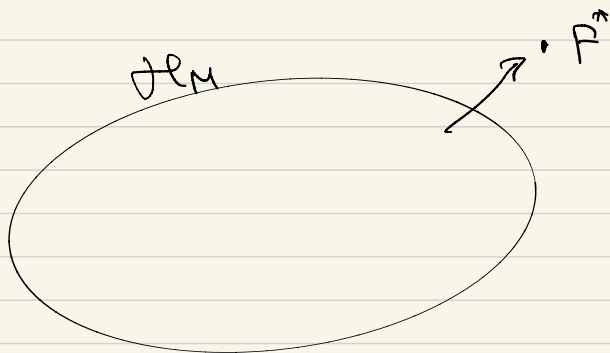
Let $K \subset \mathbb{R}^d$ be compact; $F^* : K \rightarrow \mathbb{R}$ be continuous.

Assume $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is sigmoidal, i.e. it is continuous
and $\lim_{z \rightarrow -\infty} \sigma(z) = 0$ $\lim_{z \rightarrow +\infty} \sigma(z) = 1$



Then, $\forall \varepsilon > 0$, $\exists F \in \cup_M \mathcal{H}_M$ s.t.

$$\|F^* - F\|_{C(K)} = \sup_{x \in K} |F^*(x) - F(x)| \leq \varepsilon.$$



Optimizing / Training NN

- $\mathcal{H} = \{ f_{\theta} : \theta \in \Theta \}$

Euclidean space.

- ERM $\min_{\theta \in \Theta} R(\theta) \leftarrow \min_{f \in \mathcal{H}} R(f)$

- If R is C^1

Necessary condition for optimality's

$$\nabla_{\theta} R(\hat{\theta}) = 0$$

- Gradient descent:

$$\theta_{k+1} = \theta_k - \eta \nabla R(\theta_k) \quad k=0, 1, \dots$$

\uparrow $\eta > 0$ is the learning rate.

Suppose $\theta_k \rightarrow \theta_{\infty}$

$$\cancel{\theta_{\infty}} = \cancel{\theta_{\infty}} - \underbrace{\eta \nabla R(\theta_{\infty})}_{=0}$$

- Local vs Global Minima.

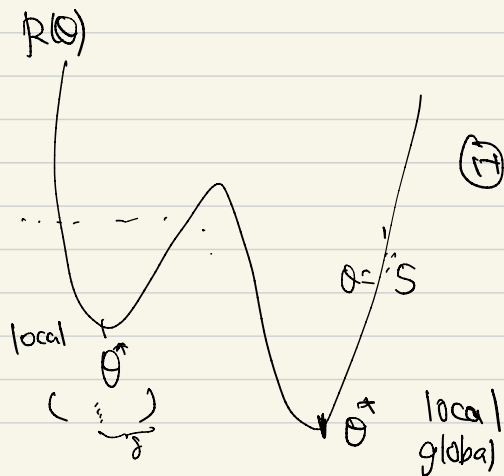
$\theta^* \in \mathbb{T}$ is a local minimum of R if

$\exists \delta > 0$ s.t.

$$R(\theta^*) \leq R(\theta) \quad \forall \|\theta - \theta^*\| \leq \delta.$$

... .. global

$$R(\theta^*) \leq R(\theta) \quad \forall \theta \in \mathbb{T}$$



$$\mathbb{T} = \mathbb{R}.$$

$$\mathbb{T} = [5, \infty)$$

The notion of local/global minima depends on

- \mathbb{T}
- δ

Deep Neural Networks

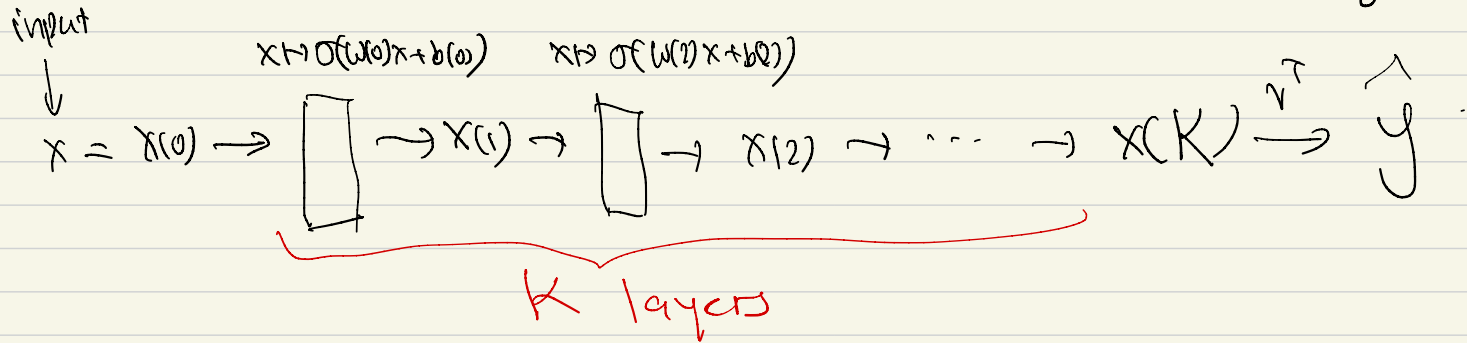
$$\mathcal{H} = \{ F : F(x) = v^T x(K), v \in \mathbb{R}^{d_K} \}$$

$$\text{where } x^{(k+1)} = \sigma(w^{(k)}x^{(k)} + b^{(k)})$$

$$k = 0, 1, \dots, K-1$$

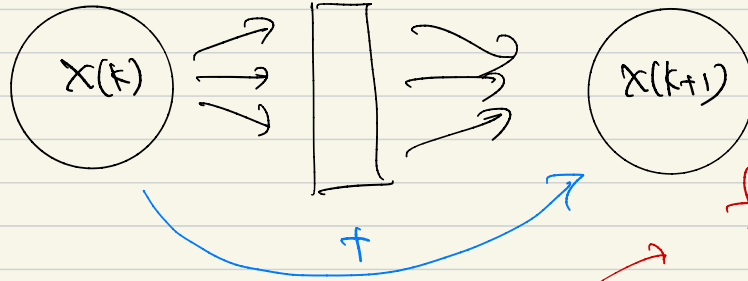
$$x^{(0)} = x \quad (\text{input})$$

$$w^{(k)} \in \mathbb{R}^{d_{k+1} \times d_k}, b^{(k)} \in \mathbb{R}^{d_{k+1}}$$



Residual NN

$$x(k+1) = x(k) + \sigma(W(k)x(k) + b(k))$$



$$f(k, x, \theta) = \sigma(Wx + b)$$

Generally, we can write.

$$x(k+1) = x(k) + f(k, x(k), \theta(k))$$

hidden state weights / trainable params

layer id.

$$k = 0, 1, 2, \dots, K-1$$

$$f: \mathbb{Z}^+ \times \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}^d$$

of layers or depth