

# The Learning Problem

Fabio A. González  
Machine Learning 2024-1

# What is a machine learning model?

- A function (Linear function, neural network, etc)
- A probability distribution.
  - Mixture of Gaussians distribution
  - Naive Bayes
- An algorithm
  - Decision tree
  - K-nearest neighbor classifier

# Models as functions

Assume a supervised learning model

Predictor:  $f: X \rightarrow Y$       $f: \mathbb{R}^n \rightarrow \{0, 1\}$

Hypothesis space: Set of functions from where the predictor is chosen

Example: Linear regression

$$H = \{f_{w,b}(x) = wx + b \mid w \in \mathbb{R}^n, b \in \mathbb{R}\}$$

$\Theta = \{w, b\} \rightarrow$  Parameters that determine  $f_\theta$

Learning: finding  $\theta$      how to find it?

# Models as probability distributions

- Model the probability distribution of the data
  - $P(x, y)$
  - $P(y|x)$
- Use the probability to make a decision

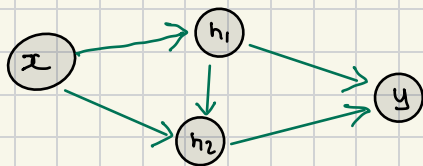
e.g. 
$$f(x) = \begin{cases} 1 & \text{if } P(y=1|x) \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

## Types of probabilistic models

- Parametric estimation  $P(x, y) = \mathcal{D}_{\theta}$  distribution with parameters  $\theta$

- Bayesian estimation  $\theta \sim \mathcal{D}'_{\theta'}$   
Use Bayesian inference to estimate  $\theta, \theta'$

- Graphical



- Non-parametric

Model the probability distribution using the data

## Models based on algorithms

- Decision trees

Hierarchical model

Information gain

- Random Forest

Ensemble model

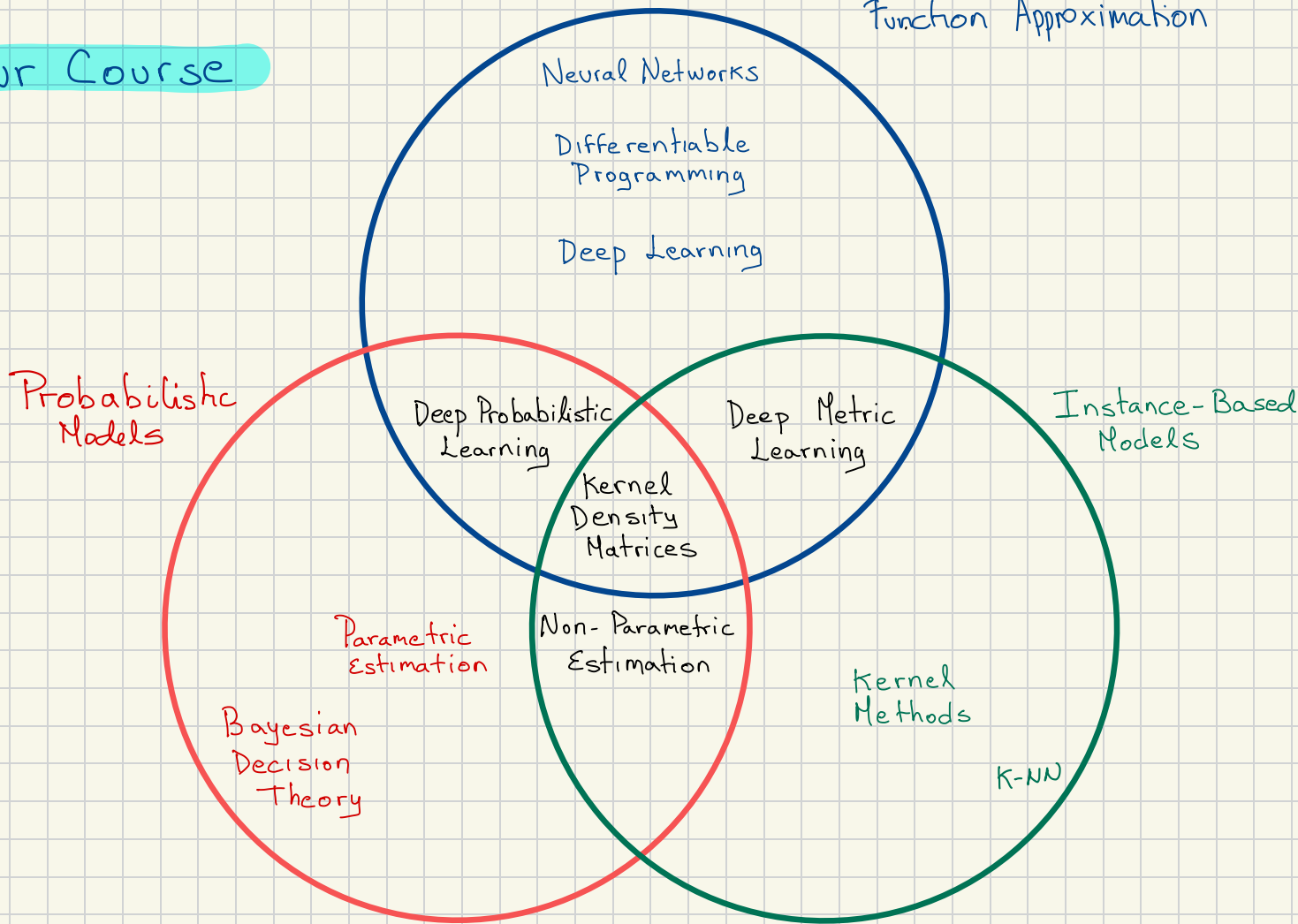
{ Bagging  
Boosting

- K-nearest neighbors

Instance based learning

# Our Course

Function Approximation



## Empirical Risk and True risk

Training data set:  $D = \{(x_i, y_i)\}_{i=1..N}$   $x_i \in X, y_i \in Y$

Empirical Risk:  $R_{\text{emp}}(f_{\theta}, D) = \frac{1}{N} \sum_{i=1}^N l(y_i, f_{\theta}(x_i))$

Loss function:  $l(y_i, \hat{y}_i) = \begin{cases} \sim 0 & \text{if } y_i \sim \hat{y}_i \\ \text{high value} & \text{otherwise} \end{cases}$

True Risk:  $R(f_{\theta}) = \mathbb{E}_{x,y} [l(y, f_{\theta}(x))]$



## Minimizing the empirical risk.

$$\theta^* = \arg \min_{\theta} R_{\text{emp}}(f_{\theta}, D)$$

$f_{\theta^*}$  best predictor on training data.

$$\theta_{\text{True}}^* = \arg \min_{\theta} R(f_{\theta})$$

$f_{\theta_{\text{True}}^*}$  best predictor for all possible data  $\left. \begin{array}{l} \text{seen} \\ \text{unseen} \end{array} \right\}$

Is it really the BEST predictor?

yes, if the Hypothesis space include all possible function

## Bayes optimal predictor (BOP)

Assumption: We know  $P(x, y) \rightarrow P(y|x)$

$$f_{\text{Bayes}}(x) = \begin{cases} 1 & \text{if } P(y=1|x) \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

Optimality:

$$\forall g: x \rightarrow y$$

$$\forall x, y \quad P(g(x) = y | x) \leq P(f_{\text{Bayes}}(x) = y | x)$$

$$R(g(x)) \geq R(f_{\text{Bayes}})$$

$$\text{with } l(y, \hat{y}) = \mathbb{1}(y \neq \hat{y})$$

# Overfitting

The model adjusts too much the training data

$$R_{\text{emp}} \downarrow \quad \text{vs} \quad R \uparrow$$

How to prevent it?

- Control  $R(f)$

Problem: I cannot calculate it

Solution: Estimate it  $\rightarrow$  Cross validation

$$R(f) \approx \hat{R}(f) = R_{\text{emp}}(f, \text{Test}) \quad \text{Test: Test dataset}$$

- Regularization

# Regularization

Add a penalization term to  $R_{emp}$ .

$$\theta^* = \arg \min_{\theta} R_{emp}(f_{\theta}, D) + \underbrace{\alpha \text{Reg}(\theta)}_{\text{Regularization term.}}$$

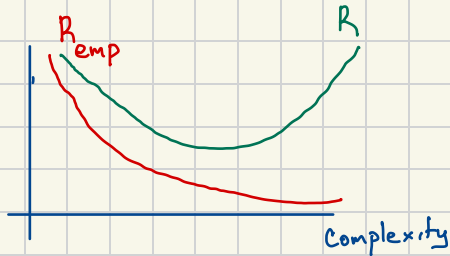
Measure of complexity



Regularization term.

Example: regularized least squares

$$w^* = \arg \min_w \frac{1}{N} \sum_{i=1}^N (f_w(x_i) - y_i)^2 + \alpha \|w\|^2$$



Lasso

$$\alpha \|w\|_1$$

$$\|w\|_1 = \sum_{i=1}^d |w_i|$$

# Vapnik-Chervonenkis Bound

$$R(\alpha) \leq R_{emp}(\alpha) + \sqrt{\left( \frac{h(\log(2l/h) + 1) - \log(\eta/4)}{l} \right)}$$

# References