

Note: *keep learning!*

- Homework #4 due next week!
- **Complete the Exercise #2: Final Session**
 - HW Assignment: Complete FP#2 on MeanShift and make your final submission of your code and results (screen shots/short doc report) via Canvas **by tomorrow midnight.**

CSC872: PAMI – Kazunori Okada (C) 2025

1

1

Note:

- Fast Prototyping Exercise #3 on LDA *Linear Discriminant Analysis* classification starts next week.
 - HW Assignment: **Read the reference paper:**
<https://bidal.sfsu.edu/~kazokada/csc872/PD3.pdf>
 - **Download:**
https://bidal.sfsu.edu/~kazokada/csc872/DATA/FaceClassification_Data.zip

CSC872: PAMI – Kazunori Okada (C) 2025

2

2

PF: Machine Learning Framework

CSC 872
Pattern Analysis and Machine Intelligence

Sources

Kevin Murphy, UBC
<http://www.cs.ubc.ca/~murphyk/>
Andrew Moore, CMU
<http://www.cs.cmu.edu/~awm/>

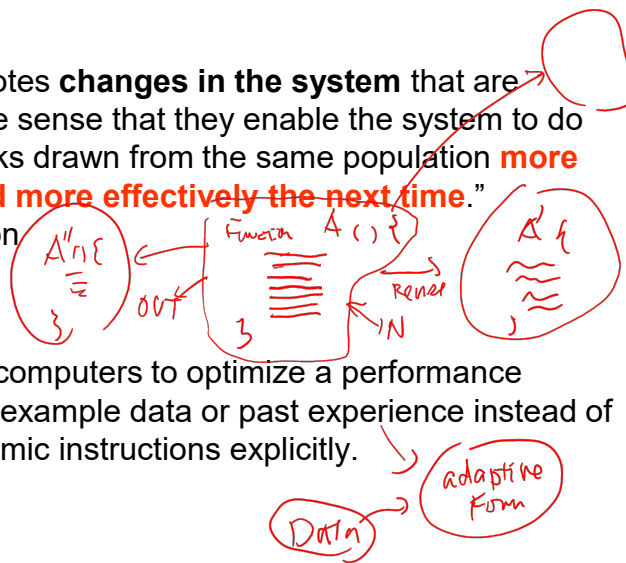
CSC872: PAMI – Kazunori Okada (C) 2025

3

What is Machine Learning?

- **For What?**
- “Learning denotes **changes in the system** that are **adaptive** in the sense that they enable the system to do the task or tasks drawn from the same population **more efficiently and more effectively the next time.**”
--Herbert Simon

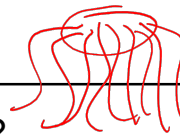
- **How?**
- Programming computers to optimize a performance criterion using example data or past experience instead of coding algorithmic instructions explicitly.



CSC872: PAMI – Kazunori Okada (C) 2025

4

Why Learning?



- What are tasks that calls for learning ???
- There is no need to learn to calculate payroll !!!
sort
- **Human expertise does not exist**
 - navigation on Mars
- **Humans are unable to explain their expertise**
 - speech recognition
- **Solution changes in time**
 - computer network routing
- **Solution needs to be adapted to particular cases**
 - user biometrics

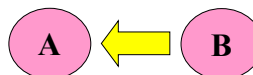
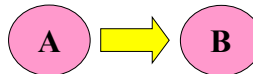
CSC872: PAMI – Kazunori Okada (C) 2025

5

5

Inductive vs Deductive Reasoning

- A entails B
 - Recall our lectures on logic
 - Deductive Reasoning
 - **Deriving B as a consequence of A**
 - All apples are fruit and All fruits grow on trees
 - Therefore all apples grow on trees
 - Inductive Reasoning
 - **Inferring some A from multiple instances of B**
 - Fuji apples grow on trees and so do McIntosh apples
 - Therefore all apples grow on trees
 - How likely that all apples grow on trees?*
- logically invalid!*
- THIS IS LEARNING !!!**



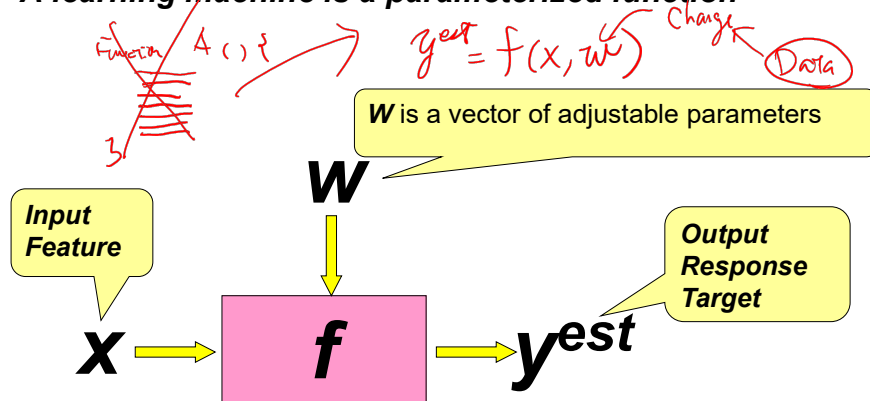
CSC872: PAMI – Kazunori Okada (C) 2025

6

6

Learning Machine *o stem cell!*

- A learning machine f takes an input x and transforms it, somehow using weights w , into a predicted output y^{est}
- **A learning machine is a parameterized function**



CSC872: PAMI – Kazunori Okada (C) 2025

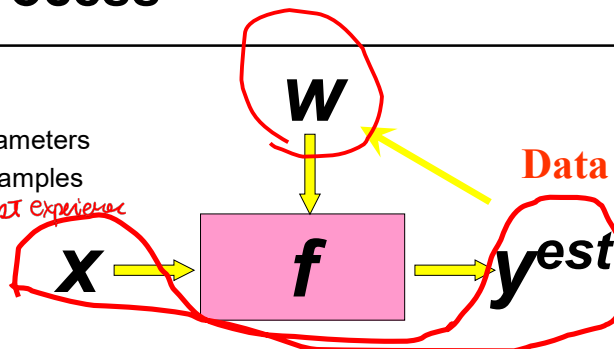
7

7

Learning Process

- **Data** $\rightarrow W$
- Deriving a set of parameters of a function from examples

\rightarrow past experiences



- **What is the nature of your data?**
- **What is the form of your function?**
- **Supervised**
- **Unsupervised**
- **Reinforcement**

CSC872: PAMI – Kazunori Okada (C) 2025

8

8

PF: Unsupervised Learning

- **No output for training = $\{x_i\}$**

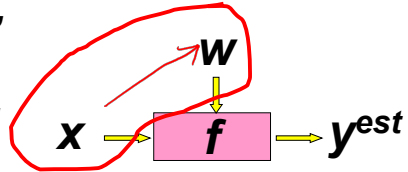
- Learning “data structure”

- “**what normally happens**”

- Probabilistic Density Estimation
 - Parametric: GMM, EM
 - Non-Parametric: KDE, MeanShift

- “**self-similar patterns**”

- Clustering
- Dimensionality reduction
 - PCA

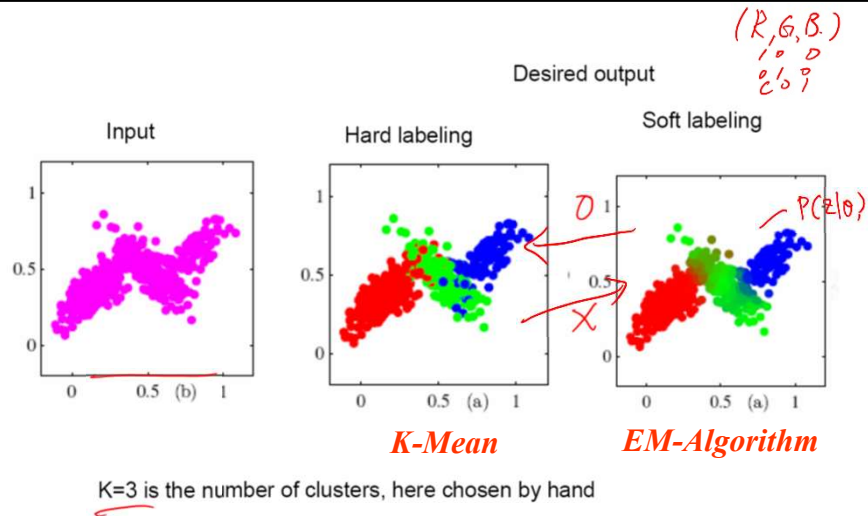


CSC872: PAMI – Kazunori Okada (C) 2025

9

9

Examples: Clustering

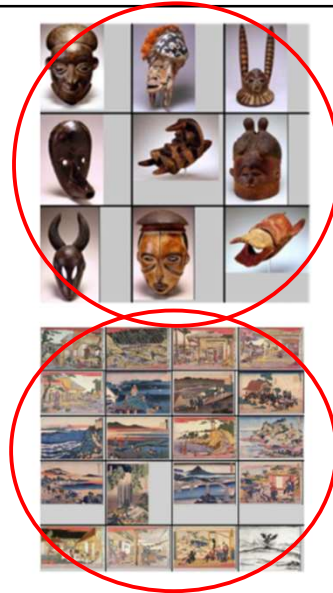


CSC872: PAMI – Kazunori Okada (C) 2025

10

10

Examples: Clustering Arts



Model each cluster by
dimensional reduction
method like PCA

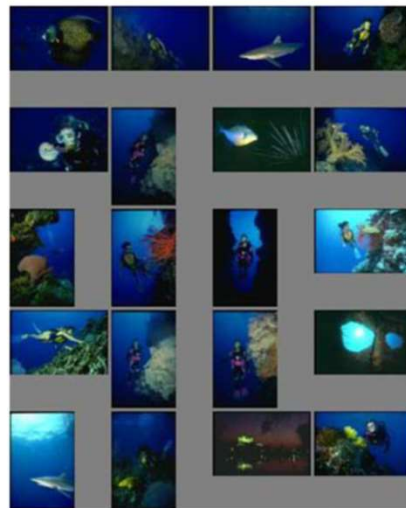
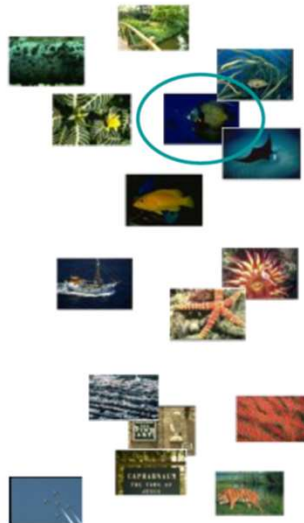
CSC872: PAMI – Kazunori Okada (C) 2025

11

11

Examples: Content-based Image Retrieval

Image Search

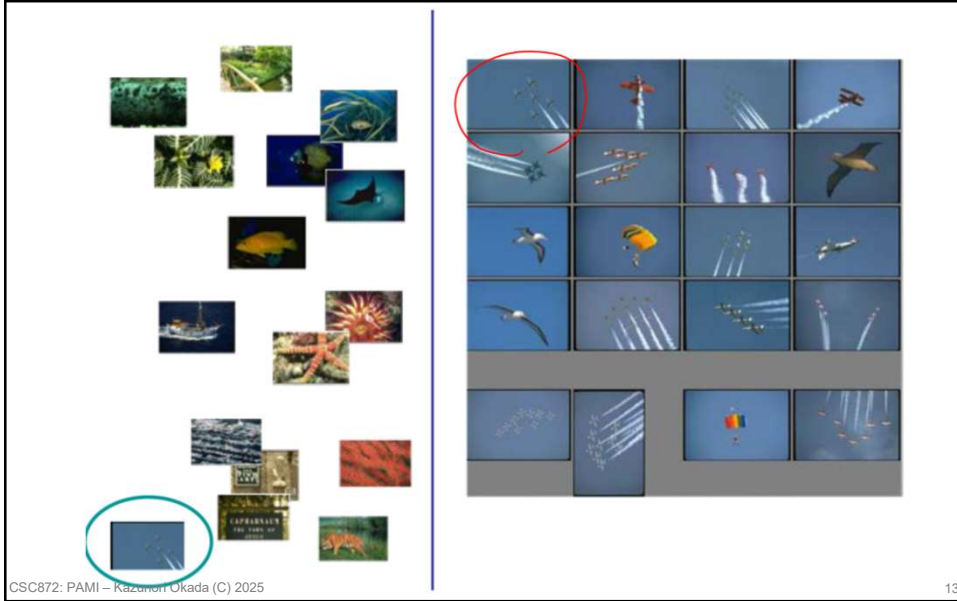


CSC872: PAMI – Kazunori Okada (C) 2025

12

12

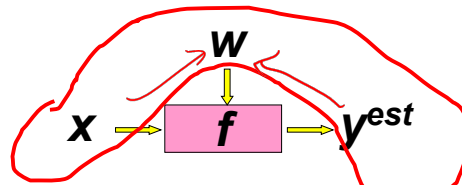
Examples: Content-based Image Retrieval



13

PF: Supervised Learning

- **Training data pairs = $\{(x_i, y_i)\}$**
- **Learning “a mapping from x to y”**
 $\{(x_1, y_1), (x_2, y_2), \dots\} = \{(x_i, y_i) \mid i=1, \dots, n\}$
- **Methods differ in terms of**
 - **The form of function** (hypothesis space)
 - **The estimation method** (used to find the best hypothesis)
 - Bayesian Classifier
 - Bayesian Regression
 - Neural Nets
 - Decision Tree/Forest
 - SVM
 - Logistic Regression
 - Boosting (Meta/Ensemble Learning)
 - Convolutional Neural Nets



CSC872: PAMI – Kazunori Okada (C) 2025

14

14

What Supervised Learning good for?

- **Prediction of future cases:**
 - Use the rule to predict the output of future inputs
 - E.g., *local weather prediction*
- **Knowledge extraction:**
 - Learn the rule that is easy to understand
 - E.g., *financially low risk if high both in income/savings*
- **Compression:**
 - Learn the rule that is simpler than the data it explains
 - E.g., *Fed's interest rate increase mortgage rate*
- **Outlier Detection:**
 - Exceptional cases that are not covered by the rule
 - E.g., *fraud*

CSC872: PAMI – Kazunori Okada (C) 2025

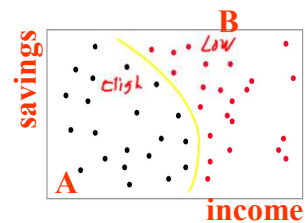
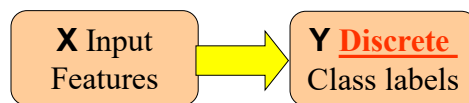
15

15

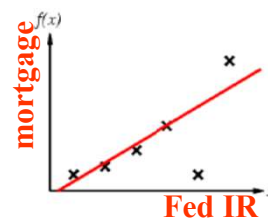
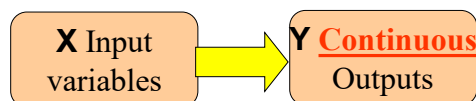
Classification vs Regression

- Nature of output different!!!

- **Classification**



- **Regression**



CSC872: PAMI – Kazunori Okada (C) 2025

16

16

Generative vs. Discriminative

- Classification by learning discriminant function

$$y^{\text{est}} = \hat{f}(x) \quad \text{Discriminative}$$

- Bayesian classification *Generative*

$$y^{\text{est}} = \operatorname{argmax}_y p(y|x)$$

- Discriminative model

$$p(y|x) = \hat{g}(x) \quad \text{Data}$$

- Generative model

$$p(x|y), p(y) \rightarrow \text{Bayes' Rule}$$

CSC872: PAMI – Kazunori Okada (C) 2025

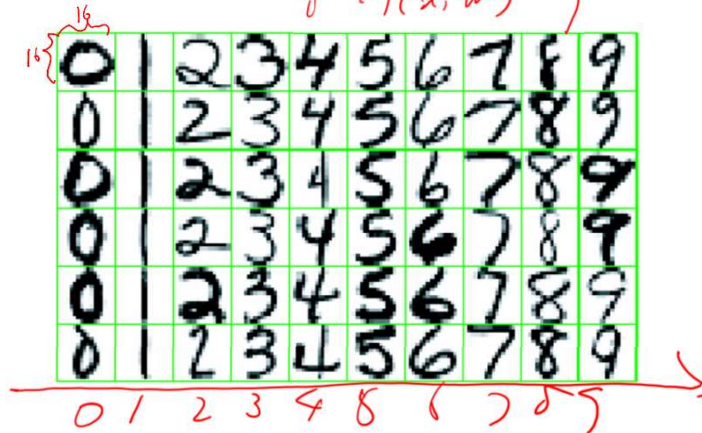
17

17

Examples: Handwritten Digit Recog

$$x_i \in \mathbb{R}^{16 \times 16} \quad y_i \in \{0, \dots, 9\}$$

$$y^{\text{est}} = f(x, w)$$



CSC872: PAMI – Kazunori Okada (C) 2025

18

18

Examples: Face Recognition

Training examples of a person $x_i \in \mathbb{R}^{128 \times 128}$



$y_i \in \{Ken, Sue, \dots\}$

Test images



$$y^{est} = \underset{\substack{\uparrow \\ \text{PCA/MHC}}}{f(\pi, w)}$$

AT&T Laboratories, Cambridge UK
<http://www.uk.research.att.com/facedatabase.html>

CSC872: PAMI – Kazunori Okada (C) 2025

19

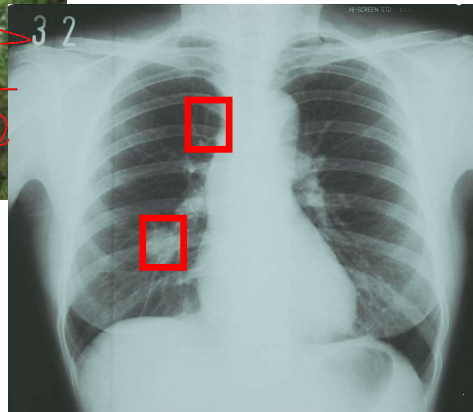
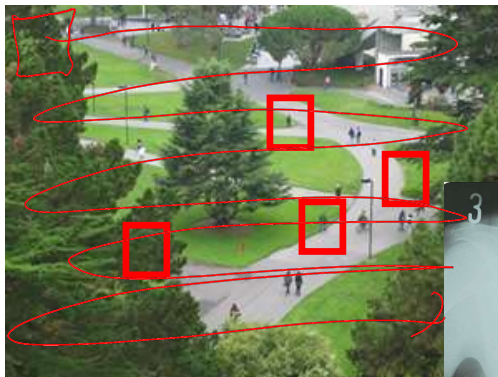
19

Examples: Object Detection

$$y^* = p(x, w)$$

$x_i \in \mathbb{R}^{16 \times 16}$

$y_i \in \{\text{Yes, No}\}$



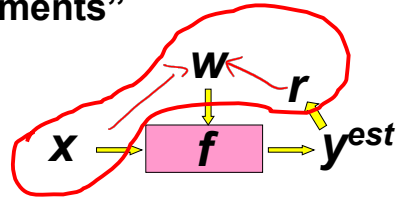
CSC872: PAMI – Kazunori Okada (C) 2025

20

20

STOP: PF: Reinforcement Learning

- **No output but delayed reward = $\{(x_i, r_i)\}$**
- Learning “credit assignments”



- ***E.g., playing chess***
 - “which action led me to winning the game?”
 - *Reward tells you if an action was good or bad*
 - *Try to learn f only with such information...*

CSC872: PAMI – Kazunori Okada (C) 2025

21

21

Major Problem Major Problem !!!

- OK. I have this problem at hand
- I want to do machine learning to solve it.
- I learned in this or that course MANY machine learning algorithms...
- **WHICH ONE AM I GONNA USE???**



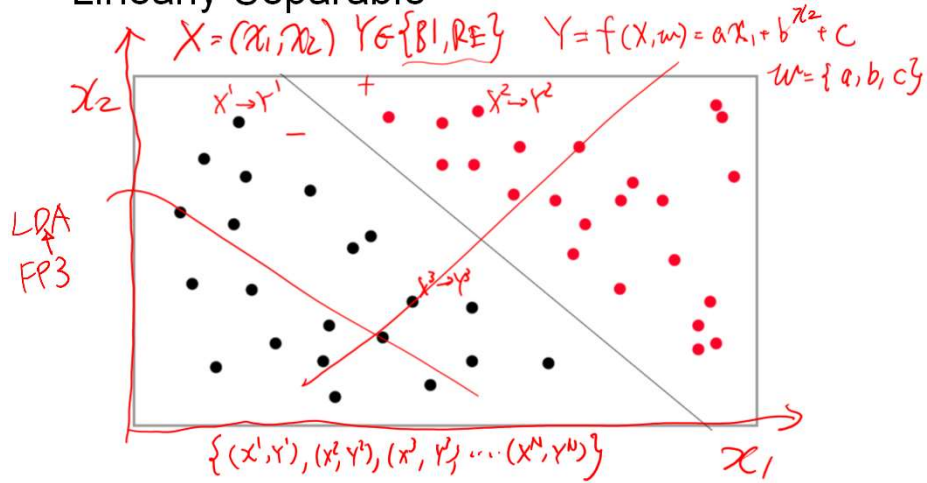
CSC872: PAMI – Kazunori Okada (C) 2025

22

22

Well the answer depends on data... Discriminant function example

- Linearly Separable



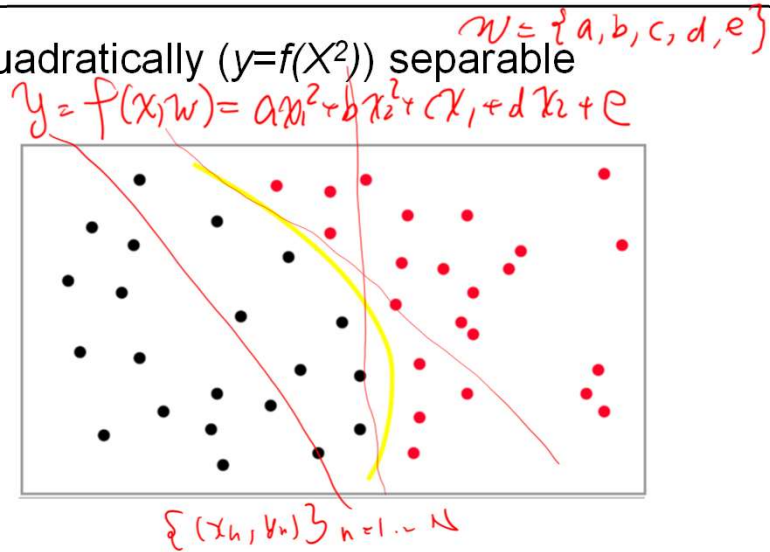
CSC872: PAMI – Kazunori Okada (C) 2025

23

23

Discriminant function

- Quadratically ($y=f(X^2)$) separable



CSC872: PAMI – Kazunori Okada (C) 2025

24

24

Training (Empirical) Errors

Given N training samples $\{(x_n, y_n)\}_{n=1}^N = D$

You learned a function \hat{f}

$$err(D) = \frac{1}{N} \sum_{n=1}^N I(\hat{f}(x_n) \neq y_n)$$

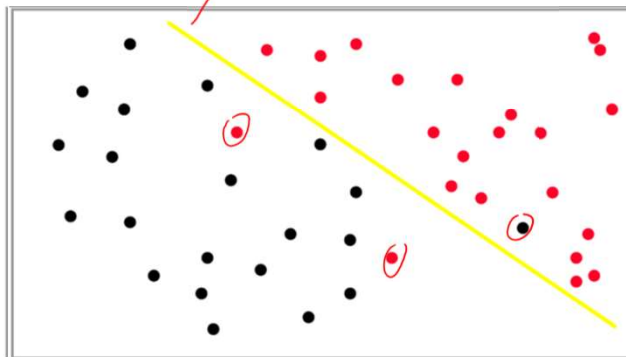
Handwritten notes:
 - $I(\hat{f}(x_n) \neq y_n)$ indicates function true or false
 - $\hat{f}(x_n)$ predicted output of x_n
 - y_n true output of x_n

A learning algorithm is designed to reduce this error

Over-fitting & Under-fitting data

- Noisy Data

Over-fitting: memorize irrelevant details of data

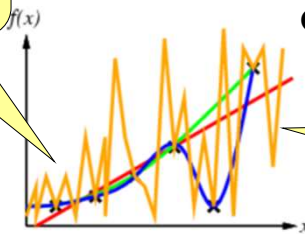


Under-fitting: ignore essential details of data

Regression Example

Each learned function is called a **hypothesis h** . A set of hypothesis form a **hypothesis space H** !

$$H = \{R, G, B, Y\}$$
$$E(R) > E(G) > E(B) = E(Y) = 0$$



Or even this!!!

$$V \subseteq H$$

Version space is formed by hypotheses with **zero** training errors.

Higher Order Polynomial

$$V = \{B, Y\}$$

These two hypotheses in the version space has the same training error !!!

CSC872: PAMI – Kazunori Okada (C) 2025

27

27

No Free Lunch Theorem

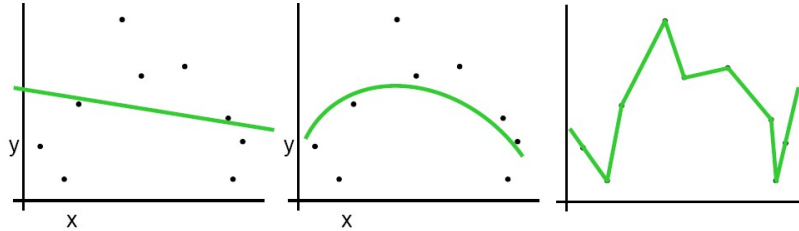
- ***Any hypothesis that agrees with all your data is as good as any other on average***
 - Wolpert and Macready (1995)
 - Wolpert (2001)
- You have two machine learning algorithms
- Ask which one is better in general?
- NFL: on average over all problems, they are the same!
- **No general take-it-all solution!**

CSC872: PAMI – Kazunori Okada (C) 2025

28

28

Choose the best given data?



Why not choose the method with the best fit to the data?

Empirical Risk Minimization: find a hypothesis that minimizes the training error

NFL tells us that **it is not a good idea !**

CSC872: PAMI – Kazunori Okada (C) 2025

29

29

Model Selection

- Problem:
Choose the best model/algorithm given a task/data
 - Linear, Quadratic, Polynomial functions?
 - Generative vs Discriminative?
 - CNN vs SVM vs AdaBoost vs Random Forest?
- **Learning is an ill-posed problem; data is not sufficient to find a unique solution**
- We need an *inductive bias / assumptions / heuristics* about the hypothesis space ...
- **AN OPEN PROBLEM!**

CSC872: PAMI – Kazunori Okada (C) 2025

30

30

What-to-do? #1 Ockham's Razor

- Bias to prefer the **simplest hypothesis** consistent with data!
- Why?
- A simpler hypothesis is less likely to be correct "by chance" and is therefore more likely to be correct on unseen new data (generalization)
- Other reasons?
 - Simpler to use (lower computational complexity)
 - Easier to train (lower sample complexity)
 - Easier to explain (more interpretable) *LDA*

CSC872: PAMI – Kazunori Okada (C) 2025

31

31

Generalization and Test Errors

Given N training samples $\{(x_n, y_n) | n = 1, \dots, N\}$

You learned a function \hat{f}

- Generalization Errors: average error for ALL possible data in domain D (**cannot calculate**)

$$E[err(D)] = \sum_{x,y \in D} p(x,y) I(\hat{f}(x) \neq y)$$

- Test Errors: approximate generalization error...

$$\widehat{err} = \frac{1}{\#\text{TestSet}} \sum_{k \in \text{TestSet}} I(\hat{f}(x) \neq y)$$

not used in Training

Prepare & use different data set

CSC872: PAMI – Kazunori Okada (C) 2025

32

32

Model Selection using Test Data

- **Goal:** *Choose the model that generalize best*
- **Generalization:** *How well a model performs on new data?*
- Steps:
 - Prepare “test data” that is different from your “training data”
 - Compute “test error” for each model (approximating “generalization error”)
 - Pick the model with lowest “test error”
- **Question: how to prepare the test data?**



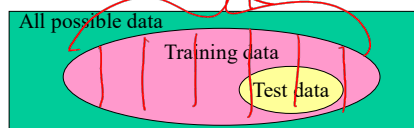
CSC872: PAMI – Kazunori Okada (C) 2025

33

33

Cross Validation (CV)

- We hold back a fraction of the training data, called a validation set, and measure performances on that using test error.



- **Random Test Set CV**
 - Randomly sample X% of data as a validation set & repeat to get an average error.
- **Leave-One-Out CV** \rightarrow # of data small
 - Validation set contains only one data point & repeat N times
- **K-fold CV**
 - Data into K partitions then repeat this K times

CSC872: PAMI – Kazunori Okada (C) 2025

34

34

More Model Selection Criteria

- Can we do model selection without doing test data sampling and CV???
- YES! These only uses training errors

- **AIC**: (Akaike information criterion)
- **BIC**: (Bayesian information criterion)
- **SRMVC**: (structural risk minimization with VC-dimension)

CSC872: PAMI – Kazunori Okada (C) 2025

35

35

Which criterion should I use?

- CV: high variance due to sampling...
- SRMVC, AIC, BIC only uses training errors
- SRMVC wildly conservative
- AIC: asymptotically same as LOO-CV
- BIC: asymptotically similar to K-fold CV
- BIC: best structure instead of best predictor

- **Open problem!!!**

CSC872: PAMI – Kazunori Okada (C) 2025

36

36

Summary

- Machine Learning Framework
 - Unsupervised/Supervised/Reinforcement Learning
 - Classification vs Regression
 - Generative vs Discriminative
 - Model Selection: Occam's Razor, Cross-Validation
- Next
 - Learning for Classification
 - Bayesian Classification
 - Discriminant Analysis
 - Performance Measure
 - FP#3 on Linear Discriminant Analysis: **READ THE PAPER!!!**

CSC872: PAMI – Kazunori Okada (C) 2025

37