

Pattern Recognition
ECSE 4410/6410 CAPA
Spring 2021

The Basics of Pattern
Recognition / ML

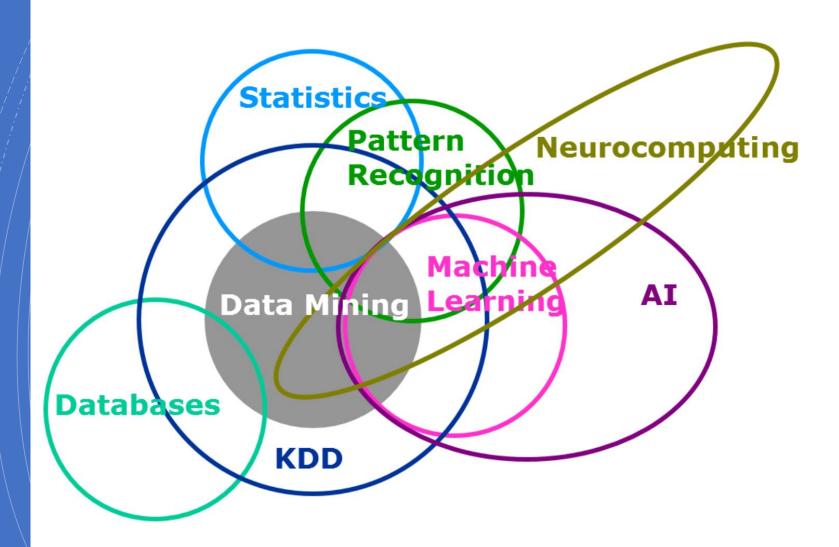
Course Instructor - Thirimachos Bourlai

January to May 2021

OVERVIEW

- How To Talk About Data in Machine Learning
- Data As you Know It
- Statistical Learning Perspective
- Computer Science Perspective
- Models and Algorithms
- Algorithms Learn a Mapping From Input to Output
- Parametric and Nonparametric Machine Learning Algorithms
- Supervised, Unsupervised and Semi-Supervised Learning
- The Bias-Variance Trade-Off
- Overfitting and Underfitting
- A Good Fit in Machine Learning

OVERVIEW



DATA in Pattern Recognition

Introduction

- Data is very important in PR
- It is important to understand and use the right terminology
- In this lecture we will talk about data terminology

Learning outcomes: We will understand the differences below...

- Data terminology used in:
 - General / spreadsheets
 - Statistics
 - Computer/Data science and Engineering

What we know – Tabular Data

PART 1

	А	В	С	D	Е	F
1		Column 1	Column 2	Column 3	Column 4	Column 5
2	Row 1	1	2	3	4	5
3	Row 2	3	4	5	6	7
4	Row 3	5	6	7	8	9
5	Row 4	7	8	9	10	11
6	Row 5	9	10	11	12	13

Column: Data Type; Attribute; Feature

Rows: Observations; Examples

Cells: Single value (numeric; categorical)

Statistics

	А	В	С	D
	x1	x2	х3	Υ
1	1.1	2.23	4.1	1
2	1.2	2.21	4.3	1
3	1.1	2.01	4.1	0
4	1	2	4.7	0

Output = Function (Input)

- Function (f) is the hypothesis that the ML algorithm needs to learn.
- The question is: given some input (features, characteristics, variables, attributes), i.e. X
 (x1,x2,x3...), what is the predicted output (Y), and how close is this output to reality

Statistics

	А	В	С	D
	x1	x2	х3	Υ
1	1.1	2.23	4.1	1
2	1.2	2.21	4.3	1
3	1.1	2.01	4.1	0
4	1	2	4.7	0

Basic Concept

- Output = Function (Input)

Typically:

Output_variable = Function (Input Vector)

As we know in Statistics

Output*DEPENDENT*_Variable =
 = Function (Input*INDEPENDENT*Variables)

Based on the framing of the prediction problem: "the output is dependent", also called - "a function of" the input

Statistics

	А	В	С	D
	x1	x2	х3	Y
1	1.1	2.23	4.1	1
2	1.2	2.21	4.3	1
3	1.1	2.01	4.1	0
4	1	2	4.7	0

To describe our problem, we need to have data

We need to have a set of algorithms and equations to describe and process the data

In statistics:

INPUT VARIABLES: we describe them as capital x (X)

- OUTPUT VARIABLE : we describe them as capital y (Y)

- Equation used : Y = f(X)

- Multiple input variables : formed as an input vector, for example

X1, X2 and X3 above.

Computer Science

	А	В	С	D
	x1	x2	х3	Υ
1	1.1	2.23	4.1	1
2	1.2	2.21	4.3	1
3	1.1	2.01	4.1	0
4	1	2	4.7	0

CS and Statistics terminology → overlap



	А	В	С	D
	Attribute 1	Attribute 2	Attribute 3	OUTPUT
Observation 1	1.1	2.23	4.1	1
Observation 2	1.2	2.21	4.3	1
Observation 3	1.1	2.01	4.1	0
Observation 4	1	2	4.7	0
	А	В	С	D
	Fearure 1	Fearure 2	Fearure 3	OUTPUT
Instance 1	1.1	2.23	4.1	1
Instance 2	1.2	2.21	4.3	1
Instance 3	1.1	2.01	4.1	0
Instance 4	1	2	4.7	0

Computer Science

	A	В	С	D
	Attribute 1	Attribute 2	Attribute 3	OUTPUT
Observation 1	1.1	2.23	4.1	1
Observation 2	1.2	2.21	4.3	1
Observation 3	1.1	2.01	4.1	0
Observation 4	1	2	4.7	0
	А	В	С	D
	Fearure 1	Fearure 2	Fearure 3	OUTPUT
Instance 1	1.1	2.23	4.1	1
Instance 2	1.2	2.21	4.3	1
In atomica 2	4.4	2.01	4.1	0
Instance 3	1.1	2.01	4.1	

PART 1

Computer Science perspective:

 $Output_Attribute = Program(Input_Attributes)$

- Row: Observation about an entity -- of different Attributes
- Column: each column → all Observations about each Attribute
 OR

Output_Feature = Program (Input_ Features)

- Row: Observation about an entity -- of different Features
- Column: each column → all Observations about each Feature

Attributes or Features > must be extracted from RAW DATA

Computer Science

	А	В	С	D
	Attribute 1	Attribute 2	Attribute 3	OUTPUT
Observation 1	1.1	2.23	4.1	1
Observation 2	1.2	2.21	4.3	1
Observation 3	1.1	2.01	4.1	0
Observation 4	1	2	4.7	0
	А	В	С	D
	Fearure 1	Fearure 2	Fearure 3	OUTPUT
Instance 1	1.1	2.23	4.1	1
Instance 2	1.2	2.21	4.3	1
Instance 3	1.1	2.01	4.1	0
Instance 4	1	2	4.7	0

PART 1

Computer Science perspective:

- Input = a set of observations or instances (data rows) that is consider "an instance"
- Each data row is considered: "a single example" or "single instance" of data observed or generated by the data source

Models and Algorithms

There is a confusion about the difference

- "A mathematical model is an abstract model that uses mathematical language to describe the behavior of a system."
- "Eykhoff (1974) defined a mathematical model as 'a representation of the essential aspects of an existing system (or a system to be constructed) which presents knowledge of that system in usable form'"

Models and Algorithms

In Computer Science – Engineering:

Model = the representation learned from data Algorithm = the process for learning this representation

Model = Algorithm(Data)

EXAMPLE:

- Model = A decision tree
 - Algorithm = C5.0
- Models = a set of coefficients
 - Algorithm = the Least Squares Linear Regression

DATA in Pattern Recognition

What have we learned thus far?

- A. Tabular data are seen in a spreadsheet as ______
- B. The statistical terms of input and output variables that maybe denoted as ____ and ___ respectively.
- C. The computer science terms of variable are names as:
 - _____, ____ and _____.
- D. When talking of models, we mean the ______.
- D. When talking of algorithms, we mean the ______.

What are we going to learn

PART 2

- 1. The mapping problem
- 2. What Predictive Modeling is
- 3. Different ML algorithms represent different

strategies for learning the mapping function

Learning a Function

Machine learning algorithms:

Learn a target function (f) that maps X → Y

$$Y = f(X)$$
 (1)

What is the "learning task"?

- To make predictions, i.e. future (Y) when new/unseen examples of input variables (X) are available
- The function (f) is unknown and that is why we need to learn it from data using these ML algorithms

Learning a Function

The mapping is not perfect

thus there is an error (e)

$$Y = f(X) + e \qquad (2)$$

• Error:

- Its <u>independent</u> of the input (X) data
- It exists due to unavailability of enough features that can support the best mapping

Making Predictions

Predictive modeling or predictive analytics

- Learn the mapping Y = f(X) to make predictions of Y for new X
- Goal: make the most accurate predictions possible

We do **not** care about:

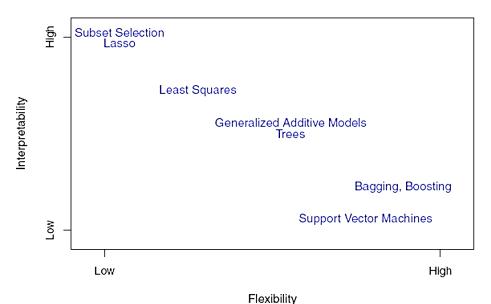
- The shape/form of the function (f) we need to learn We care about:
- Function (f) that needs to make the most accurate predictions

Statistical Inference (SI)

- **SI Goal**: if we care about learning the mapping of Y = f(X) to learn more about data relationships \rightarrow use simple methods
- In SI we care more about "understanding the learned model and its form" rather that to make the MOST accurate predictions

Making Predictions

- "The major difference between machine learning and statistics is their purpose.
- Machine learning models are designed to make the most accurate predictions possible.
- Statistical models are designed for inference about the relationships between variables."



PART 2

Source: https://healthcare.ai/machine-learning-versus-statistics-use/

Learning a Function

- ML = techniques for estimating the target function (f) to predict the output variable Y given X
- Different models → different assumptions about f (linear or nonlinear)
- Different ML algorithms → different assumptions about the function type needed to optimize the model/hypothesis to approximate it

What do we need to do?

- Try different ML algorithms
- Study the problem to determine which algorithms have been used before to efficiently solve the problem under investigation

What have we learned in Part Two?

- A. ML algorithms estimate the _____ of output variables (Y) given input variables (X)
- B. Different ML algorithms make _____about the form of the underlying function.
- C. When we don't know much about the form of the target function, we must try ______ to see what works best.

What are we going to learn

- PML algorithms "simplify the mapping to a known functional form"
- Non-PML algorithms → learn any I/O mapping
- 3. All ML algorithms are organized into PML and Non-PML categories

Parametric ML Algorithms

"A learning model that **summarizes** data with **a set of parameters** of <u>fixed size</u> (*independent of the number of training examples*) is called a parametric model. No matter how much data you throw at a parametric model, it won't change its mind about how many parameters it needs."

Artificial Intelligence: A Modern Approach, page 737

Parametric ML Algorithms:

- Step 1: Select a form for function (f)
- Step 2: Use the training data to learn the function's coefficients

Parametric ML Algorithms

PART 3

Hypothesis: the functional form of a line is used to <u>simplify the</u> <u>learning process</u> of a problem

Basic mapping function (as we will see - used in linear regression) is a line:

$$a_0 + a_1 \cdot X_1 + a_2 \cdot X_2 = 0$$

- a0, a1 and a2: line coefficients (control the intercept and slope)
- X₁ and X₂ are two input variables

Need to do:

- Estimate a0, a1 and a2 to have a **predictive model** for the problem

Parametric ML Algorithms

PART 3

The hypothesis is often assumed to have a functional form that is a linear combination of the input variables

- Thus,...... Parametric machine learning algorithms

→ (called) Linearmachine learning algorithms

CHALLENGE:

- The actual underlying function that best maps the problem may not be a linear function
- It may almost be linear but require some minor transformation of the input data to best represent the problem (data)
- It may NOT be a line and thus, this hypothesis will result in poor system performance

Parametric ML Algorithms - **OVERVIEW**

Example Algorithms	Benefits	Limitations
•Logistic	• Easy to understand	• Constrained to a
Regression	• Fast to learn	specific form
•Linear	• Less Data to Train	 Not suited for
Discriminant	• Work well even	complicated
Analysis	with an approxi-	problems
•Perceptron	mate fit	• Poor Fit to the
		underline function

Non-Parametric ML Algorithms

 "Algorithms that do not make strong assumptions about the form of the mapping function are called Non-Parametric ML Algorithms"

 "By not making assumptions, they are free to learn any functional form from the training data"

"Nonparametric methods are good when you have a lot of data and no prior knowledge, and when you don't want to worry too much about choosing just the right features"

— <u>Artificial Intelligence: A Modern Approach</u>, page 757

Non-Parametric ML Algorithms

 "Algorithms that do not make strong assumptions about the form of the mapping function are called Non-Parametric ML Algorithms"

"By not making assumptions, they are free to learn any functional form from the training data"

"Nonparametric methods are good when you have a lot of data and no prior knowledge, and when you don't want to worry too much about choosing just the right features"

— <u>Artificial Intelligence: A Modern Approach</u>, page 757

Non-Parametric ML Algorithms

Goal: to best fit the training data in constructing the mapping function, while, at the same time they try to generalize to unseen data

Non-Parametric ML Algorithm can fit many functional forms

Example: k-NN algorithm

Non-Parametric ML Algorithms - **OVERVIEW**

Example Algorithms	Benefits	Limitations	
Decision TreesNaive BayesSVMsNNs	 Flexibility No assumptions about the underlying function Higher Performance prediction models 	 Require more training data More params → Slower to train Overfitting training data is a 	
		high risk	

What have we learned in Part Three?

A. What is a parametric and a non-parametric function

B. Examples, benefits and limitations of **parametric** functions

C. Examples, benefits and limitations of non-

parametric functions

Supervised, Unsupervised and Semi-Supervised Learning

What we are going to learn

- A. About the classification and regression **supervised** learning problems
- B. About the clustering and association **unsupervised** learning problems
- **C. Example algorithms** used for supervised and unsupervised problems

Supervised Learning

$$Y = f(X)$$

SL GOAL:

Approximate the mapping function well \rightarrow predict well the output data (Y) from the input data (X)

SL algorithms are learning from the training data → supervise the training

SL algorithms are grouped into:

- Classification: When the output variable is a category
- Regression: when the output variable is a real value

SL examples:

SVM

• Linear regression → Regression

• Random forest → Classification/regression

→ Classification problems

PART 4

Unsupervised Learning

USL:

- Only have input data (X) → no corresponding output
- There is no teacher to supervise anything
- The algorithms discover natural grouping of data on their own

USL algorithms are grouped into:

- Clustering: Discover the inherent data groupings
- Association: Discover rules that describe large portions of your data

USL examples:

- K-means → clustering problems
- Apriori algorithm → association problems

PART 4

Semi-Supervised Learning

Example: photo archive with some labeled images but all the rest are unlabeled

- Labeling data: Expensive or time consuming
- Unlabeled data: low cost; easy to collect; easy to store.

SSL: When we have a lot of input data (X) and some labeled data (Y)

STEP 1

 USL can be used → Discover the structure in the input variables and make good predictions for the unlabeled data

• STEP 2

- SL can be the grouped data as an input (training) to a SL algorithm
- Use the trained model by the SL to make predictions on new unseen data

SL, USL and Semi-SL

What have we learned in Part Four?

A.	: All data is (labeled/unlabeled)?; Algorithms
	learn the input's data inner structure
В.	: All data is (labeled/unlabeled)? and the
	algorithms learn to predict the output from the input data
C.	: Some data is (labeled/unlabeled)? but most of
	it is (labeled/unlabeled)? and a mixture of SL and USL can be
	used.

PART 5

Bias - Variance

What have we going to learn?

- A. Learning error → is broken down into bias or variance error
- Bias → <u>simplifying algorithm assumptions</u> so that the problem can be solved easier
- C. Variance → model sensitivity to <u>training data</u>
 changes
- D. <u>Model prediction</u> (applied ML) is <u>better</u>understood via the bias vs. variance trade-off



Prediction Error

For any ML algorithm → Prediction

Error depends on:

- Bias Error
- Variance Error
- Irreducible Error

Bias Error

Bias → <u>simplifying algorithm assumptions</u> so that the hypothesis/target function can be solved easier

- Low Bias
 - Less assumptions about the form of the hypothesis function
- High-Bias
 - More assumptions about the form of the hypothesis function

Bias Error

High-bias ML algorithms:

- Linear Regression
- Linear Discriminant Analysis
- Logistic Regression

Linear algorithms

- High bias
- Fast to learn
- Easier to understand but generally less flexible.

BUT:

- They have lower predictive performance (complex problems)
- Do not meet the simplifying assumptions of the algorithmic bias.

Low-bias ML algorithms:

- Decision Trees
- k-Nearest Neighbors
- Support Vector Machines

Variance Error

Variance → model sensitivity to <u>training data</u> changes

- Training data is used to estimate the hypothesis/target function
- To do that we need a ML algorithm
- ML algorithms are expected to have some variance

- IDEAL SCENARIO When changing datasets, the ML algorithm used will be able to pick out the hidden underlying I/O variable mapping
- **REAL WORLD** When changing datasets, the ML algorithm picks out the hidden underlying I/O variable mapping with an error that depends on many factors

Variance Error

Low Variance: When changing the

training dataset, we can also have small

changes made, when estimating the

hypothesis/target function

High Variance: Suggests large changes...

Variance Error

High Variance (Low Bias) ML algorithms:

- Decision Trees
- k-Nearest Neighbors
- Support Vector Machines

Nonlinear ML algorithms

- High Variance
- A lot of flexibility

For example:

- Decision trees have a high variance → large changes
- Pruning → lowers down variance before they decision trees are used

Low Variance (High Bias) ML algorithms:

- Linear Regression
- Linear Discriminant Analysis
- Logistic Regression

PART 5

Bias vs. Variance TRADE-OFF

What are the GOAL of supervised ML algorithms?

- Achieve "LOW Bias" and "LOW Variance"
- Achieve Good Prediction Performance
- Perform Algorithm Parameterization to balance out bias and variance (increasing one → results in decreasing the other and vice versa NO WAY-OUT RELATIONSHIP)

REMEMBER:

- Linear ML algorithms
 - High bias
 - How variance
- Nonlinear ML algorithms
 - Low bias
 - High variance

Bias vs. Variance TRADE-OFF

Examples with specific algorithms:

k-NN:

- Low bias (less assumptions); High variance (large changes to HTH when changing training data)
- How to change the TRADE-OFF?
 - ➤ Increase k → increases the # of neighbors that contribute to the prediction → increases model bias

SVMs:

- **Low** bias and **High** variance
- How to change the TRADE-OFF?
 - ➤ Increasing the C parameter (impacts the # of margin violations allowed in the training data) → increases the bias and decreases the variance

Bias vs. Variance Reality

The algorithms we choose and their configuration (e.g. kNN with large k) \rightarrow impact the trade-off

In practice:

- Calculating the real bias and real variance error is not possible → the actual underlying target function is unknown
- However, the bias / variance framework can help us understand the behavior of ML algorithms when aiming to determine the highest predictive performance.

BIAS - VARIANCE

What have we learned in Part Five?

- A. "The simplifying assumptions made by the model to make the target function easier to approximate" is called________
- B. "The amount that the estimate of the target function will change when changing training data" is called_______
- C. The B and V Trade-off relationship can be avoided → YES / NO?

Overfitting and Underfitting

What have we going to learn?

A. Overfitting: learning the training data too well and thus, not being able to perform well using new data

- B. Underfitting: insufficiently learning how to solve the problem from the training data
- C. Overfitting is in practice the most common problem in practice and the ways to address it.

Learn about Generalization

- Inductive learning: using the training data to learn the target function.
- Induction: learning general concepts from specific examples (supervised learning)
- [Deduction: learning specific concepts from general rules]
- Generalization: how well the concepts learned by a model apply to specific examples not seen by the model when it was learning.

Overfitting and Underfitting

Machine learning model GOAL: generalize well, i.e. the model built from the training data should also work with any data from the problem domain"

 This allows us to make predictions in the future on data the model has never seen.

- Overfitting and underfitting: terminology used in ML when we talk about how well a model learns and generalizes to new data.
- Overfitting and underfitting: main causes for poor performance of ML algorithms.

PART 6

EXTRA READING: https://towardsdatascience.com/what-are-overfitting-and-underfitting-in-machine-learning-a96b30864690

Statistical Fit

PART 6

 "Goodness of fit test is a statistical hypothesis test to see how well sample data fit a distribution from a population with a <u>normal</u> distribution"

 "The test shows if your sample data represents the data you would expect to find in the actual population or if it is somehow skewed"

Statistical Fit

PART 6

 "Goodness-of-fit establishes the discrepancy between the observed values and those that would be expected of the model in a normal distribution case"

Methods: "There are <u>multiple methods for</u> <u>determining goodness-of-fit</u>. Some of the **most popular methods** used in statistics include the chi-square, the Kolmogorov-Smirnov test, the Anderson-Darling test and the Shipiro-Wilk test"

Statistical Fit – ML Language

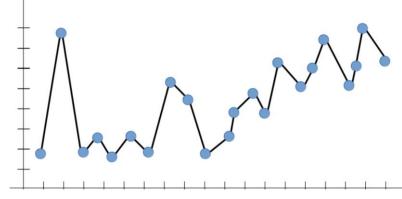
 A statistical fit (SFit) refers to the measures used to estimate how well a target function is approximated

Machine Learning:

- SL algorithms "seek to approximate the unknown underlying mapping function for the output variables given the input variables"
- While SFit methods are useful in ML (e.g. calculating the residual errors), some of them assume we know the form of the target function we are approximating, which is not the case in ML.

15

Overfitting



PART 6

 Overfitting: "a modeling error that occurs when a function is too closely fit to a limited set of data points. Overfitting the model generally takes the form of making an overly complex model to explain idiosyncrasies in the data under study."

In Practice: the data is often impacted with some degree of error or random noise within it → when we use a model that conforms too closely to slightly inaccurate (training) data can impact the model with substantial errors and reduce its predictive power (performance) on new data and negatively impact the model's ability to generalize.

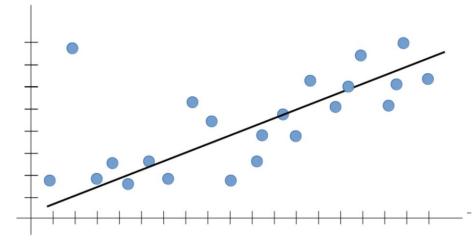
Overfitting – Where?

- More likely with nonparametric and nonlinear models that have more flexibility when learning a target function
- Many such nonlinear algorithms include <u>parameters or</u>
 <u>techniques</u> to **limit and constrain** how much detail the model learns.

Example – Decision Trees (nonparametric algorithm)

- DTs are flexible and subject to overfitting
- By pruning a tree (after it has learned) in order to remove some of the detail it has picked up, overfitting is controlled.

Underfitting



- An Ufit model can neither model the training data nor generalize to new data
- An underfit model is not a suitable model → has poor performance
 on the training data
- It is often not discussed as it is easy to detect given a good performance metric
- The remedy is to move on and try alternate machine learning algorithms
- It provides a good contrast to the problem of concept of overfitting

What is a Good Fit in ML?

In Practice

- We need a model between underfitting and overfitting → results in the highest performance on NEW DATA
- Is it possible? Is it easy?

How do we achieve this goal:

- Monitor the over time performance of the ML algorithm we selected (as it is learning the training data).
- Plot the performance <u>over time</u> (training data; validation; test data)

As the algorithm learns:

- 1. The model error on the training data goes down
- 2. The model error on the **test data**...... goes down

Challenge - Generalization is at stake - if we train too long

- (1) Error continues to decrease ... thus, the model overfits
- (2) At a certain point, the test error starts increasing

What is a Good Fit in ML?

Not good practice

- For every ML algorithm: we do not stop learning, right before error on the test set starts to increase, and thus we expect to work well on unseen data.
- The test set is <u>no longer unseen</u> or a standalone objective measure → knowledge leaks into the training process.

Good practice – techniques to use for the model to work well on unseen data:

- Use resampling methods
- Use a validation dataset

Dealing with Overfitting

Overfitting & underfitting \rightarrow poor model performance

Overfitting is by far the most common problem in practical ML – does not work as expected on unseen data.

When evaluating ML algorithms, you can limit overfitting by:

- Using resampling techniques (to estimate model accuracy)
- 2. Use a validation set

We will revisit all this topic again....

PART 6

Dealing with Overfitting

k-Fold Cross-Validation: The most popular resampling technique

 "It allows you to train and test your model ktimes on different subsets of training data and build up an estimate of the performance of a machine learning model on unseen data"

Validation dataset: a training data subset that you hold back to further tune your model.

- ❖ Step 1: Training tuning your ML algorithms
- Step 2: Evaluate the learned models on the validation dataset
- Step 3: Assess model on test set

Overfitting and Underfitting

What have we learned in Part SIX?

A. ______: Good performance on the training data,
poor generalization to other data.
B. ______: Poor performance on the training data and
poor generalization to other data.

Questions?

