

Training/Testing and Regularization

Il-Chul Moon
Dept. of Industrial and Systems Engineering
KAIST

icmoon@kaist.ac.kr

Weekly Objectives

- Understand the concept of bias and variance
 - Know the concept of over-fitting and under-fitting
 - Able to segment two sources, bias and variance, of error
- Understand the bias and variance trade-off
 - Understand the concept of Occam's razor
 - Able to perform cross-validation
 - Know various performance metrics for supervised machine learning
- Understand the concept of regularization
 - Know how to apply regularization to
 - Linear regression
 - Logistic regression
 - Support vector machine

CONCEPT OF BIAS AND VARIANCE

Up To This Point...

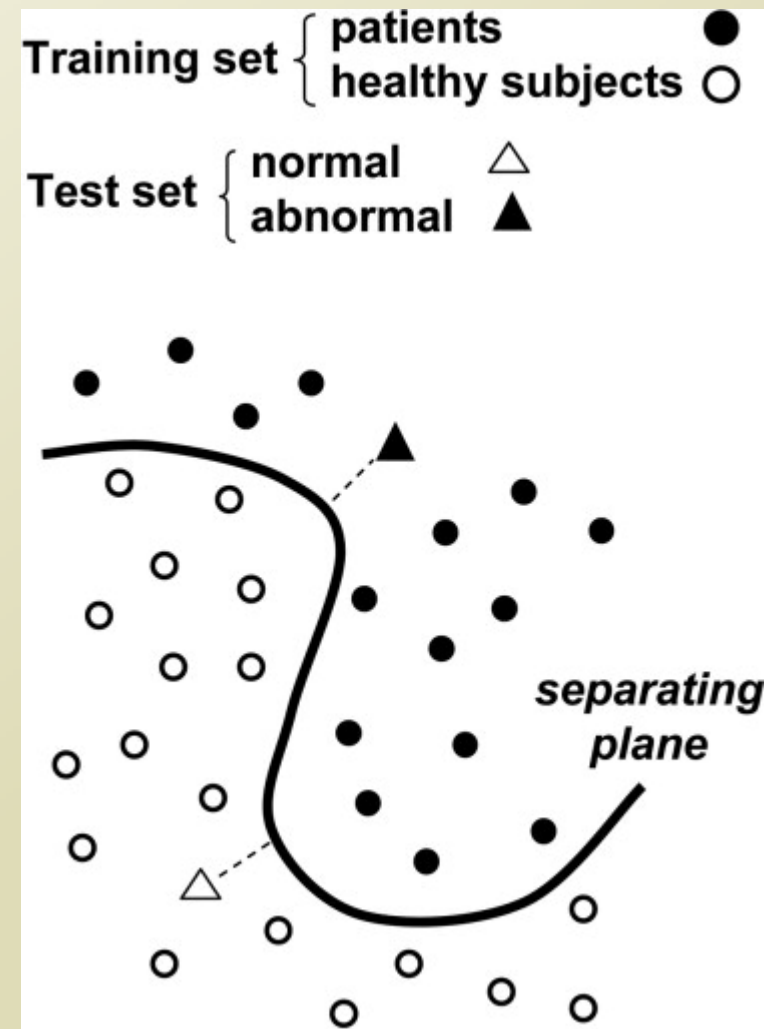
- Now, you are supposed to have some knowledge in classifications
 - Naïve Bayes
 - Logistic Regression
 - Support Vector Machine
- SVM is still a commonly used machine learning algorithm for classifications
- Functioning is *kind of* done
- Efficiency and accuracy now becomes a problem

Better Machine Learning Approach?

- Accurate prediction result
 - Ex) with this NB classifier, I can filter spams with 95% accuracy!
- Is this a right claim?
 - The validity of accuracy
 - No clear definition
 - Why not use other performance metrics? Such as Precision/Recall, F-Measure
 - The validity of dataset
 - Spams??
 - How many spams?
 - Where did you gathered?
 - Big variance in the spams?
 - Is the spam mail evolving?
 - From Nigerian prince scheme to something else?

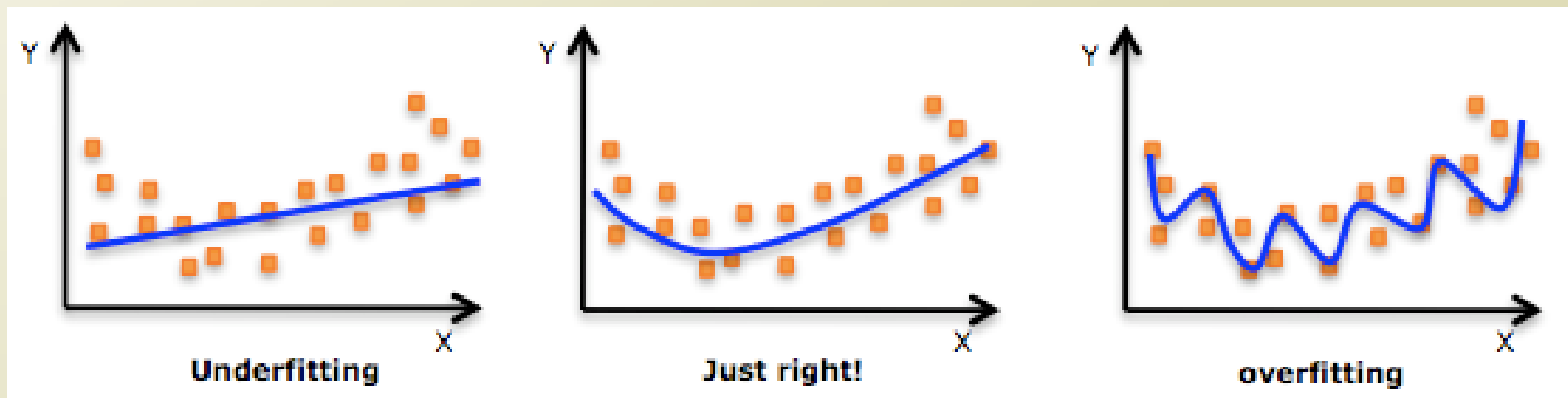
Training and Testing

- Training
 - Parameter inference procedure
 - Prior knowledge, past experience
 - There is no guarantee that this will work in the future
 - ML's Achilles gun is the stable/static distribution of learning targets.
 - Why ML does not work in the future?
 - The domain changes, or the current domain does not show enough variance
 - The ML algorithms inherently have problems
- Testing
 - Testing the learned ML algorithms/the inferred parameters
 - New dataset that is unrelated to the training process
 - Imitating the future instances
 - By setting aside a subset of observations



Over-Fitting and Under-Fitting

- Imaging this scenario
 - You are given N points to train a ML algorithm
 - You are going to learn a simple polynomial regression function
 - $Y=F(x)$
 - The degree of F is undetermined. Can be linear or non-linear
- Considering the three F s in the below, which looks better?



Tuning Model Complexity

- One degree, two degree, and N degree trained functions
 - As the degree increases, the model becomes complex
 - Is complex model better?
- Then, where do we stop in developing a complex model?
 - Is there any measure to calculate the complexity and the generality?
- There is a trade-off between the complexity of a model and the generality of a dataset.

