




# Machine Learning & Data Mining @ NuIEEE

Tcharammmmmmm - Final session

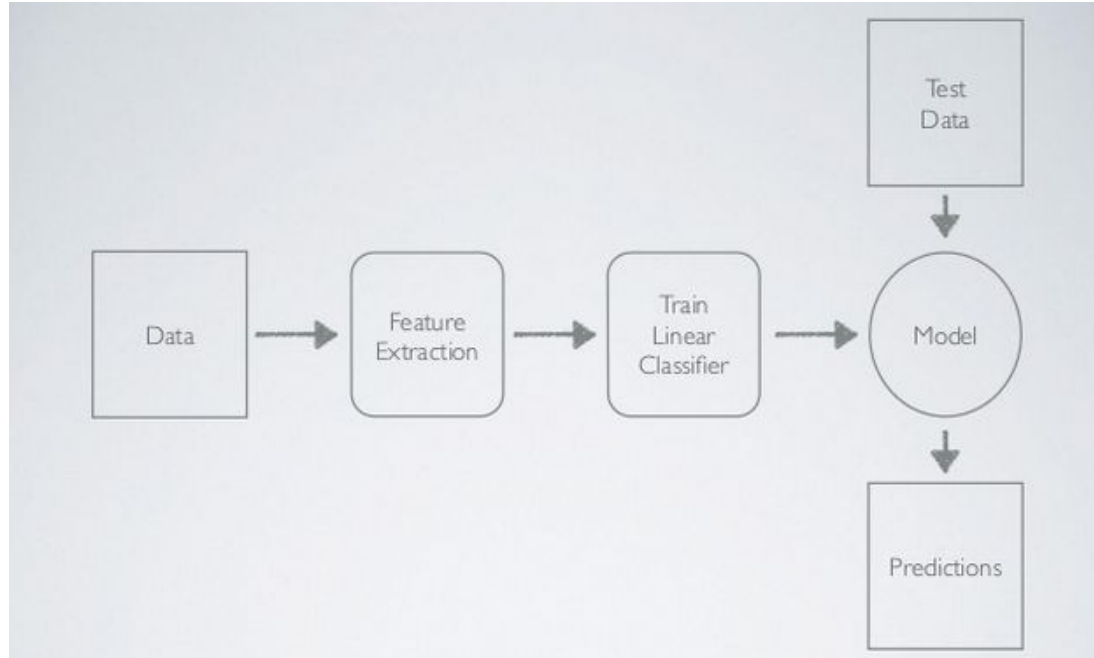
Miguel Sandim & Paula Fortuna



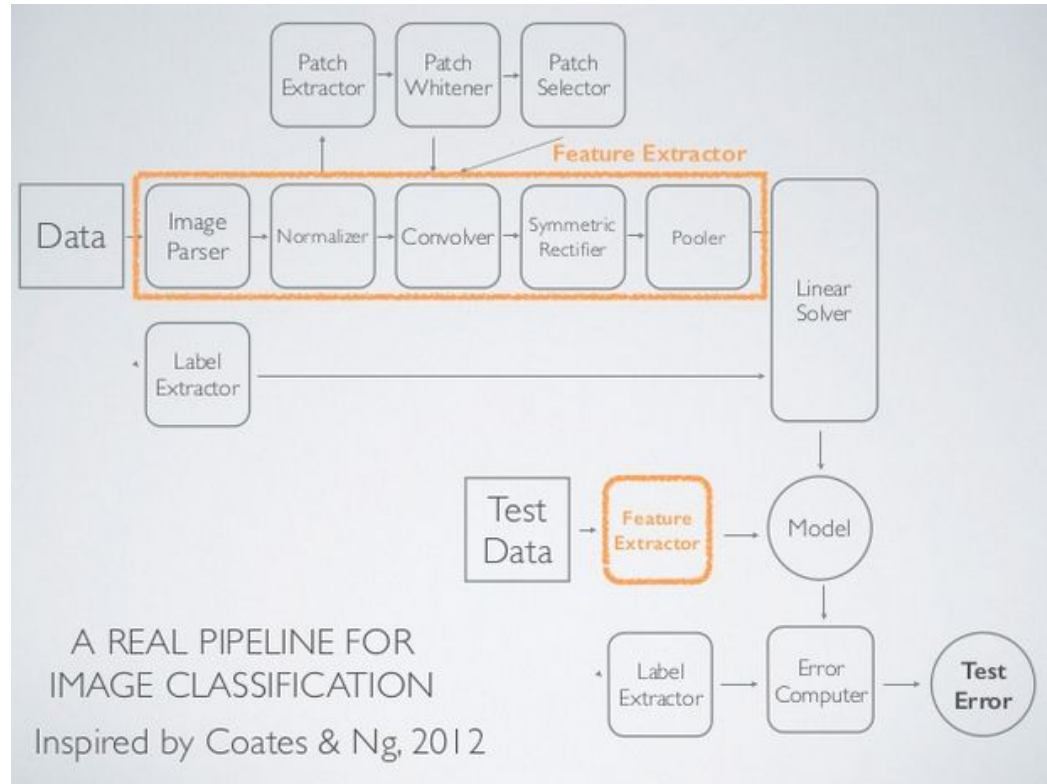
# Final session topics

- ML pipeline
- How to choose a family of algorithms?
- How to evaluate if an algorithm is working
- Overfitting & underfitting
- Validation Methods - K-folds

# Machine Learning Pipeline

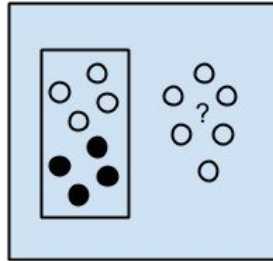


# Machine Learning Pipeline

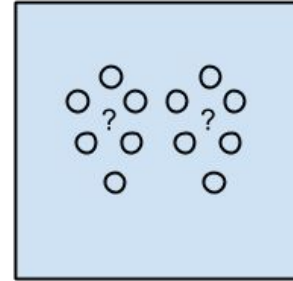


# Machine Learning Model - Algorithms

## Learning Style



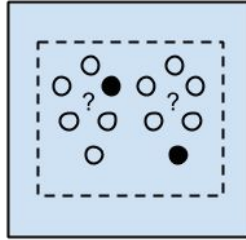
Supervised Learning  
Algorithms



Unsupervised Learning  
Algorithms

# Machine Learning Model - Algorithms

## Learning Style

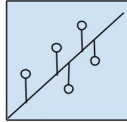


Semi-supervised  
Learning Algorithms

# Algorithms Grouped By Similarity and examples

## Regression Algorithms

- Linear Regression



## Clustering Algorithms

- K-Means

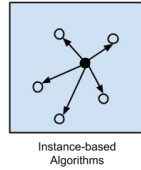


## Ensemble Algorithms

- Random Forests

## Instance-based Algorithms

- k-Nearest Neighbour (kNN)



## Association Rule Algorithms

- Apriori Algorithm

## Feature Selection

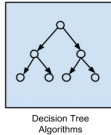
- Wrapper Models

## Regularization Algorithms

- Ridge Regression

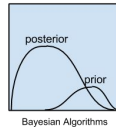
## Decision Trees

- C4.5



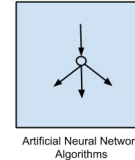
## Bayesian Algorithms

- Naive Bayes



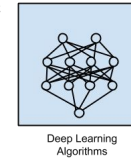
## Neural Networks

- Back-propagation



## Deep Learning

- Deep Boltzmann Machine

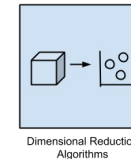


## Other

- SVM

## Dimensionality Reduction

- PCA







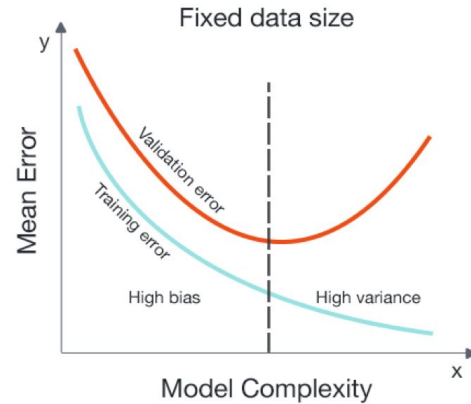
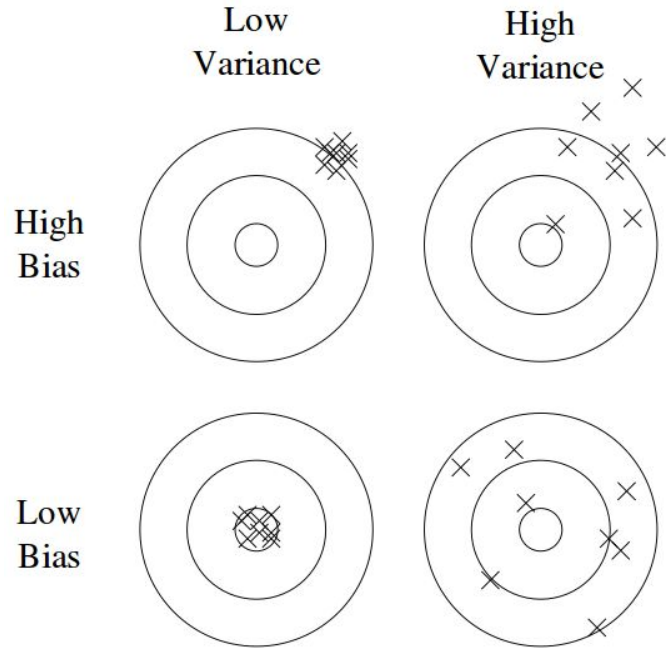
# How to choose the right algorithm?

- Use Past Experience To Choose An Algorithm
- Start with the simplest approach
- Visualize data
- Use Trial And Error To Choose An Algorithm
  - use diverse algorithms
  - tuning the algorithm parameters
  - using ensemble methods
- Spot Check Algorithms in R - process to automatically apply a set of algorithms

# Variance and Bias

- **Bias** is a learner's tendency to consistently learn the wrong thing (difficulty in finding the "pattern" in the data);
  - **Example:** approximating a real-world problem by a simple linear regression model may lead to high bias.
- **Variance** is a learner's tendency to learn random things irrespective of the "real signal" (learning everything is a "pattern", including the "noise" in the data). Another interpretation is the sensitivity of the model to changes in the train set.
  - **Example:** decision trees learned on different training sets generated by the same phenomenon are often very different, when in fact they should be the same.

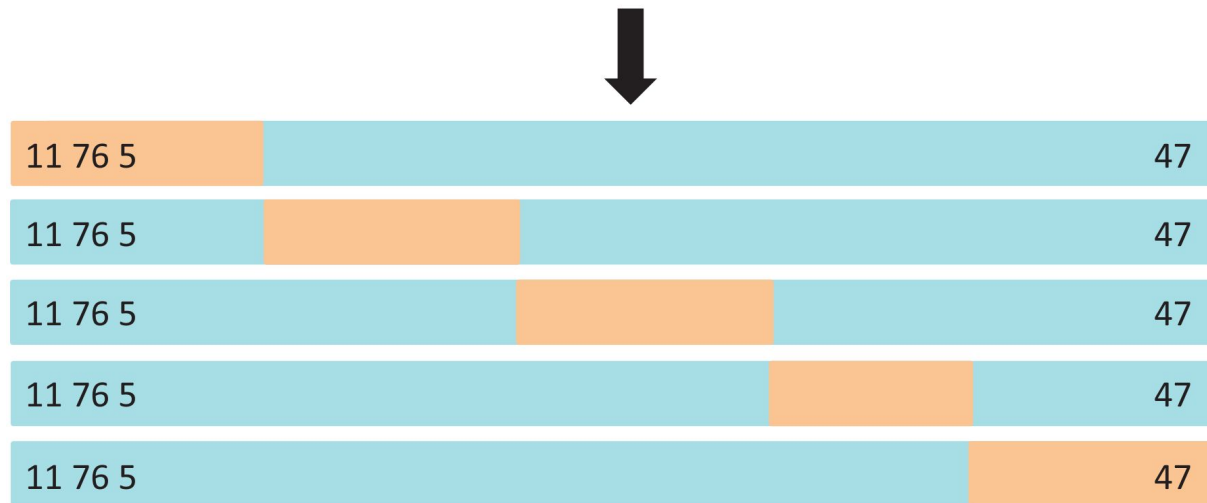
# Overfitting and Underfitting



# Overfitting is the major problem!

- How to solve?
  - **Regularization** - include an additional term that penalizes highly complex models;
  - **Cross-Validation** - choose the best model's parameters based on the error rate on a validation set. This set should be different from the training set (used to train the model) and the test set (used to assess the final model's accuracy).
    - Leave-One-Out Cross-Validation
    - k-Fold Cross-Validation

# k-Fold Cross Validation



$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k MSE_i$$



# References

- [1] <http://www.slideshare.net/jeykottalam/pipelines-ampcamp>
- [2] <http://machinelearningmastery.com/a-tour-of-machine-learning-algorithms/>
- [3] <http://machinelearningmastery.com/evaluate-machine-learning-algorithms-with-r/>
- [4] Domingos, Pedro. "A few useful things to know about machine learning." *Communications of the ACM* 55.10 (2012): 78-87.
- [5] James, Gareth, et al. An introduction to statistical learning. Vol. 112. New York: springer, 2013.