

Machine Learning Fundamentals

Practical Machine Learning (with R)

UC Berkeley

Fall 2015

Topics

➤ Administrative

- *Applied Predicative Modeling*, Max Kuhn
- *The Art of R Programming*, Norm Matloff
- *Elements of Statistical Learning*, Hastie, Friedman, Tibshirani

➤ Review Q&A

- Assignments and Grading



REVIEW



EXPECTATIONS

- You have installed **R** and **Rstudio**
- If you are new to **R**, you will have checked out one of the resources and have started becoming familiar with syntax and functions.
- If you are not familiar with source control, you will have investigated **git/github** and/or **sourcetree**

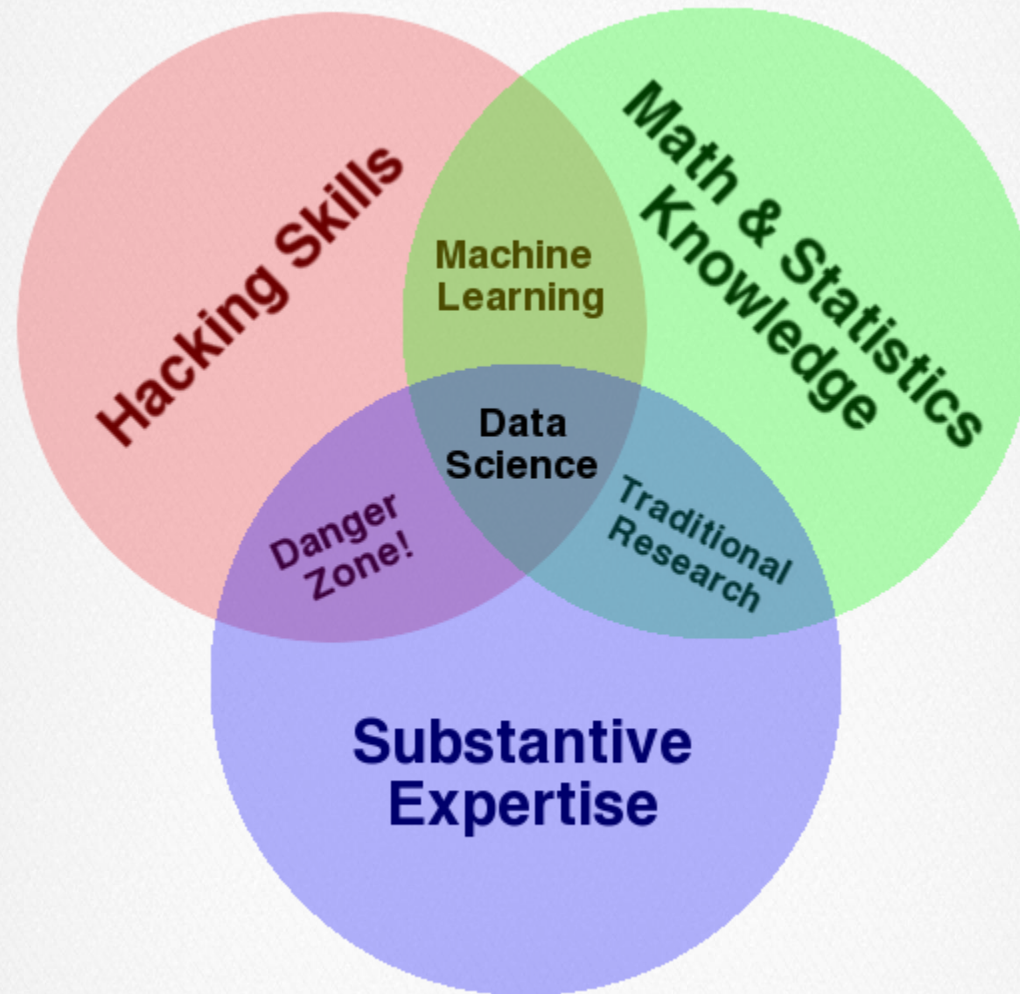


Review

- Class Objective → Practice of ML
- Advantages of R
(popularity, community, extensibility)
- Elite Coding
 - Follow Established Design Patterns
 - Adopt Standards
 - Use Version Control → git(hub) / sourcetree
- Set-up R / Rstudio



Data Science Venn Diagram



Ref. <http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>

USEFUL R PACKAGES

- ⇒ **ML Framework:** *caret* (Classification and Regression Training)
- ⇒ **Pipe operators:** *magrittr*, *pipeR*, *backpipe* (shiny)
- ⇒ **Tables:** *data.tables*, *dplyr*
- ⇒ **Visualization:** *ggvis*, *ggplot2*
- ⇒ **Reporting:** *knitr*, *rmarkdown*



R Resources

ONLINE

- ➔ (META)CRAN
 - [Packages](#)
 - [Task Views](#)
- ➔ [Stackoverflow.com](#)
- ➔ [r-bloggers.com](#)
- ➔ [Advanced R Programming](#)
- ➔ Github

Offline

The Art of R Programming

Norm Matloff

ISBN-13: 858-2592222227

ISBN-10: 1593273843

R in Action: Data Analysis and Graphics with R

Robert Kabacoff

ISBN-13: 978-1617291388

ISBN-10: 1617291382



COURSE WEBSITE

⇒ <https://github.com/CSX460>



GIT / GITHUB / SOURCE TREE

⇒ Workflow

- clone
- branch
- (work)
- add
- commit (early and often)
 - tag
- push
- Also checkout, status, log



ADMINISTRATIVE



GRADING



GRADES

- ⇒ **Exams and Quizzes (20%)**
- ⇒ **Class Participation and Exercises(30%)**
- ⇒ **Project (50%)**
 - Identify problem you want to tackle
 - Frame the problem
 - Build Features
 - Review linear model and cart
 - Build Model
 - Deploy
- ⇒ **Attendance is Mandatory**



HIGH DIMENSIONAL SPACES



EXAMPLE OF ML ALGORITHM(S)

- Spam Filter
- handwriting recognition (svm)
- Traffic engineering (lights)
- Weather prediction
- Sentiment analysis (social media)
- Netflix Recommender
- Fraud detection (Visa)
- Imaging processing
- Intrusion detection
- Self-driving cars



QUESTION 1

What is machine learning?

A formal **process** for building a **model**



QUESTION 2

What is a model?

a ***function*** that ***estimates*** a ***response***
associated with (a set of) known
predictors

$$\hat{y} = f(\vec{x})$$



QUESTION 1

What is machine learning?

A formal **process** of building a model

How do we find f ?



WHAT ARE THE PROPERTIES OF f

- Should be easy to evaluate
- Takes a one or more values of inputs
- Yields a single output value for each input
- Output, $\hat{\mathbf{y}}$, should be “close” to observed values, \mathbf{y} :

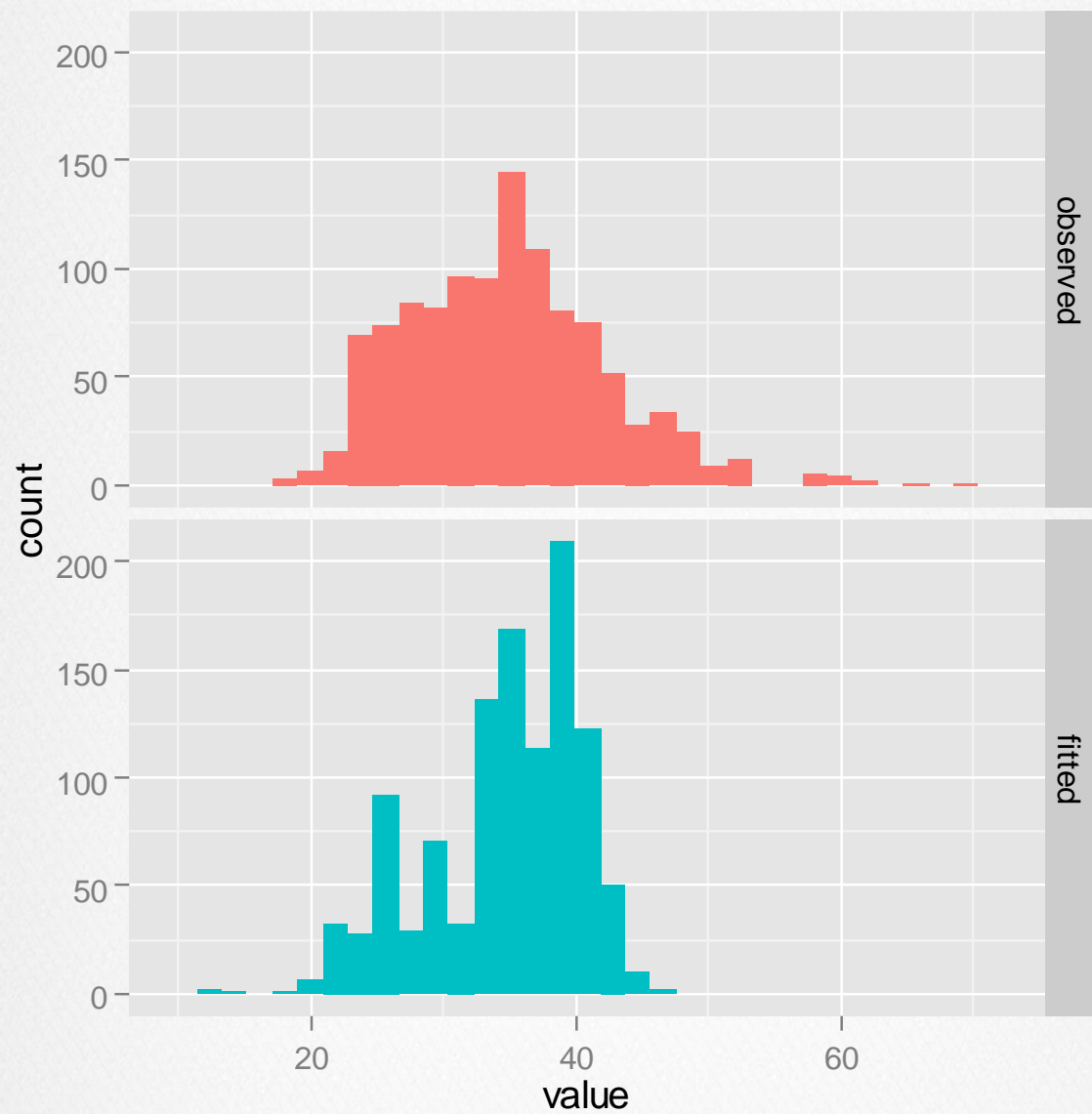
$$\hat{\mathbf{y}} \sim \mathbf{y}$$

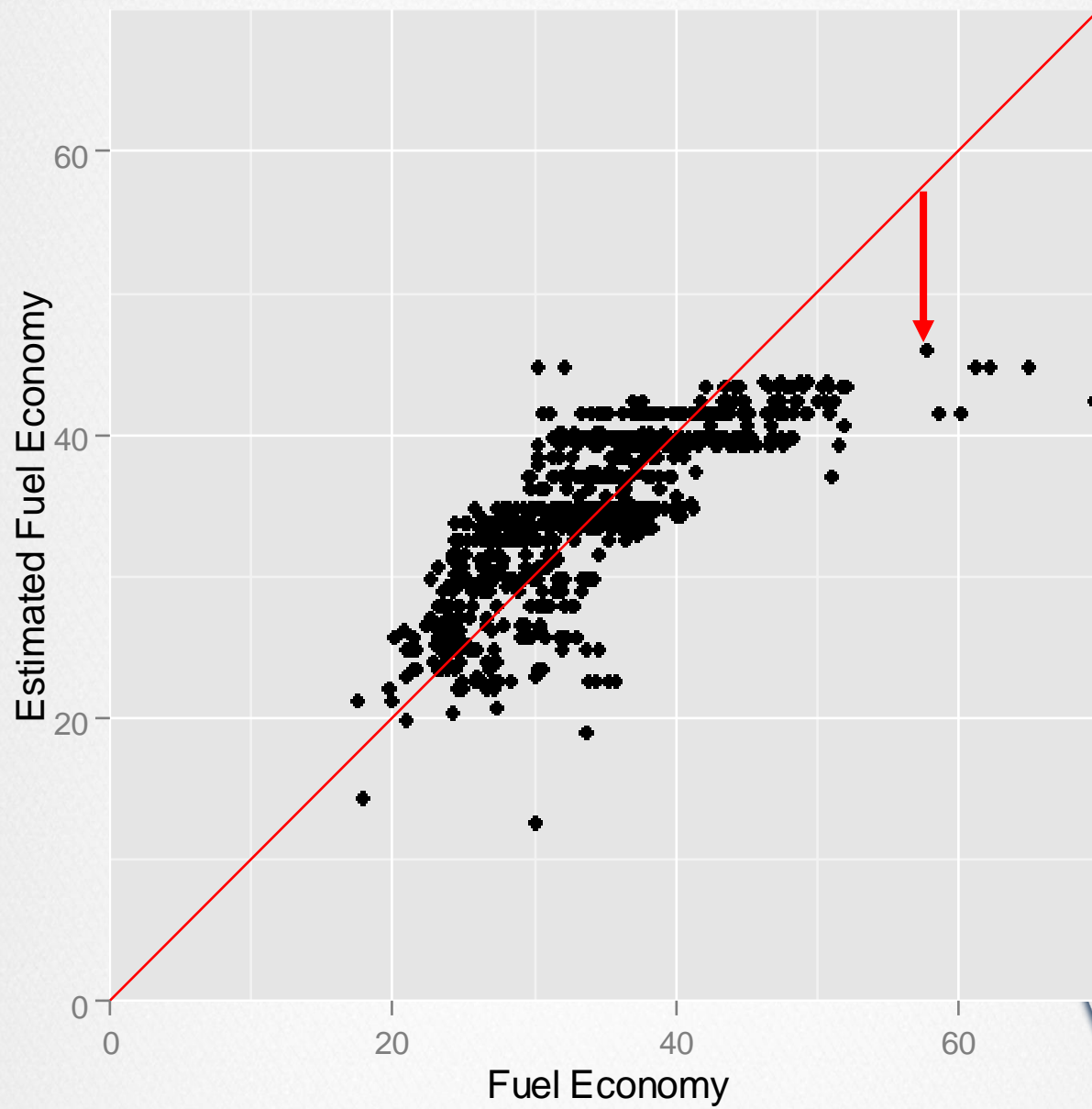


QUESTION 3

How good is the model?







MEASUREMENTS FOR ERROR FOR MODEL



WHAT ARE THE PROPERTIES OF f

- ⇒ Should be easy to evaluate
- ⇒ Takes a one or more values of inputs
- ⇒ Yields a single output value for each input
- ⇒ Can measure the error
- ~~⇒ Output, \hat{y} , should be “close” to observed values, y : $\hat{y} \sim y$~~

What else



→ The number of functions available?



3 REQUIREMENT FOR ALGORITHM

- A method for evaluating how well the algorithm performs (**ERRORS**)
- A restricted class of function (**MODEL**)
- A process for proceeding through the restricted class of functions to identify the functions (**SEARCH/OPTIMIZATION**)

LINEAR REGRESSION

- ⇒ Errors : Minimize Squared Error
- ⇒ Model:

$$\hat{y} = \beta_0 + \sum_{i=1}^p \beta_i x_i$$

$$\hat{y} = \beta_0 + \mathbf{x}^T \boldsymbol{\beta}$$



OPTIMIZATION TECHNIQUES

- Direct Solution
- Recursive Goal Seeking



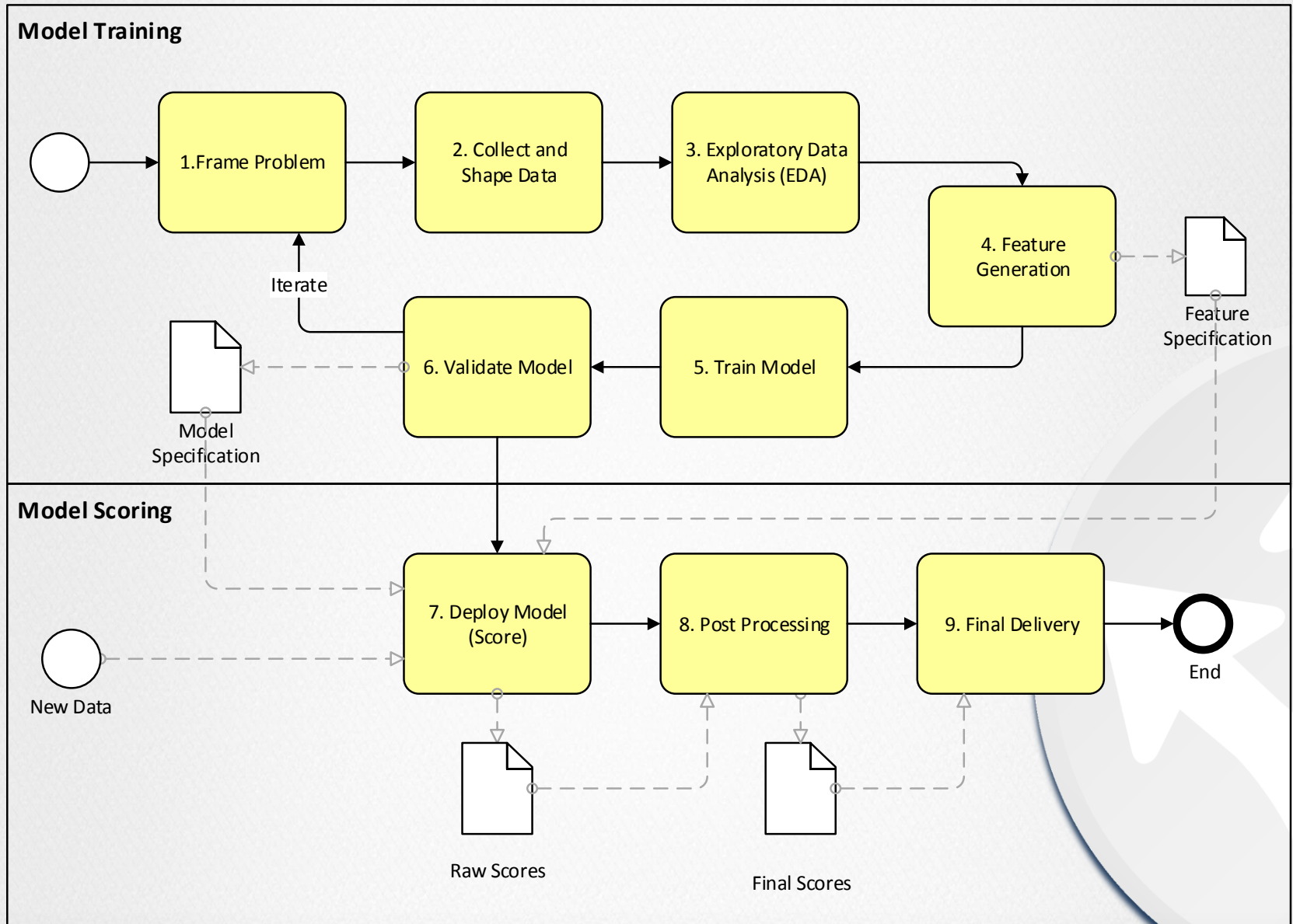
MORE THAN LEARNING

- Making a practical model entails more than learning...

What are the other requirements?



Comprehensive ML Process





Problem

Write a function to
calculate the RMSE,

Write a function to
calculate the MAE



MACHINE LEARNING EXAMPLE

⇒ AppliedPredictiveModelling::FuelEconomy



APPENDIX



Given a vector of numbers (x), write a function (f) that returns a vector of numbers containing the *product* of every other number excluding the current index.

Example:

```
> x <- c( 1, 5, 2, 8 )
```

```
> f(x)
```

```
[1] 80 16 40 10
```

```
# 5*2*8, 1*2*8, 1*5*8, 1*2*5
```

Given a vector of numbers (x)
write a function (f) that returns a
vector of numbers containing the
product of every other integer
excluding the current index.

Example:

```
> x <- c( 1, 5, 2, 8 )  
> f(x)  
[1] 80 16 40 10  
  
# 5*2*8, 1*2*8, 1*5*8, 1*2*5
```

Solution:

```
f <- function(x) prod(x) / x
```


CLASS OVERVIEW : 1

- Introduction to R, setting up the ML developers environment
 - Installing R
 - Installing R Studio
 - Installing packages from CRAN, Bioconductor and Github
 - Exercises



CLASS OVERVIEW : 2

- ➔ Fundamentals of Machine Learning
 - Machine learning overview
 - Regression and classification
 - Supervised, unsupervised, and semi-supervised
 - Algorithm types and requirements
 - Exercises



CLASS OVERVIEW : 3

➤ Linear Regression

- OLS Regression
- Data partitioning
- Model evaluation and tuning
- Exercises



CLASS OVERVIEW : 4

- ⇒ Logistic Regression
 - Logistic Regression
 - Exercises



CLASS OVERVIEW : 5

- Advanced Techniques: Partitioning Methods
 - CART/Regression Trees
 - Clustering
 - K Nearest Neighbors
 - Exercises



CLASS OVERVIEW : 6

- Advanced Techniques: Partitioning Methods
 - CART/Regression Trees
 - Clustering
 - K Nearest Neighbors
 - Exercises



CLASS OVERVIEW : 7

➔ Advanced Techniques

- Bagging
- Bagged Trees / Random Forests
- Exercises



CLASS OVERVIEW : 8

- ➔ Advanced Techniques: Boosting
 - Boosting
 - Neural Networks
 - Support Vector Machines
 - Exercises



CLASS OVERVIEW : 8

⇒ Deployment

- Diving into the data lake
- Optimization
- Delivery and Production



CLASS OVERVIEW : 9

⇒ Final Lecture

- Exercises
- Exam

