

R Fundamentals

Practical Machine Learning (with R)

UC Berkeley

Agenda

- Administrative
 - Role Call
 - Missing Coordinates
 - Class Google Group
 - Class Location
- Review
- New Topics



REVIEW AND EXPECTATIONS



EXPECTATIONS: R

- 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.



EXPECTATIONS: GIT

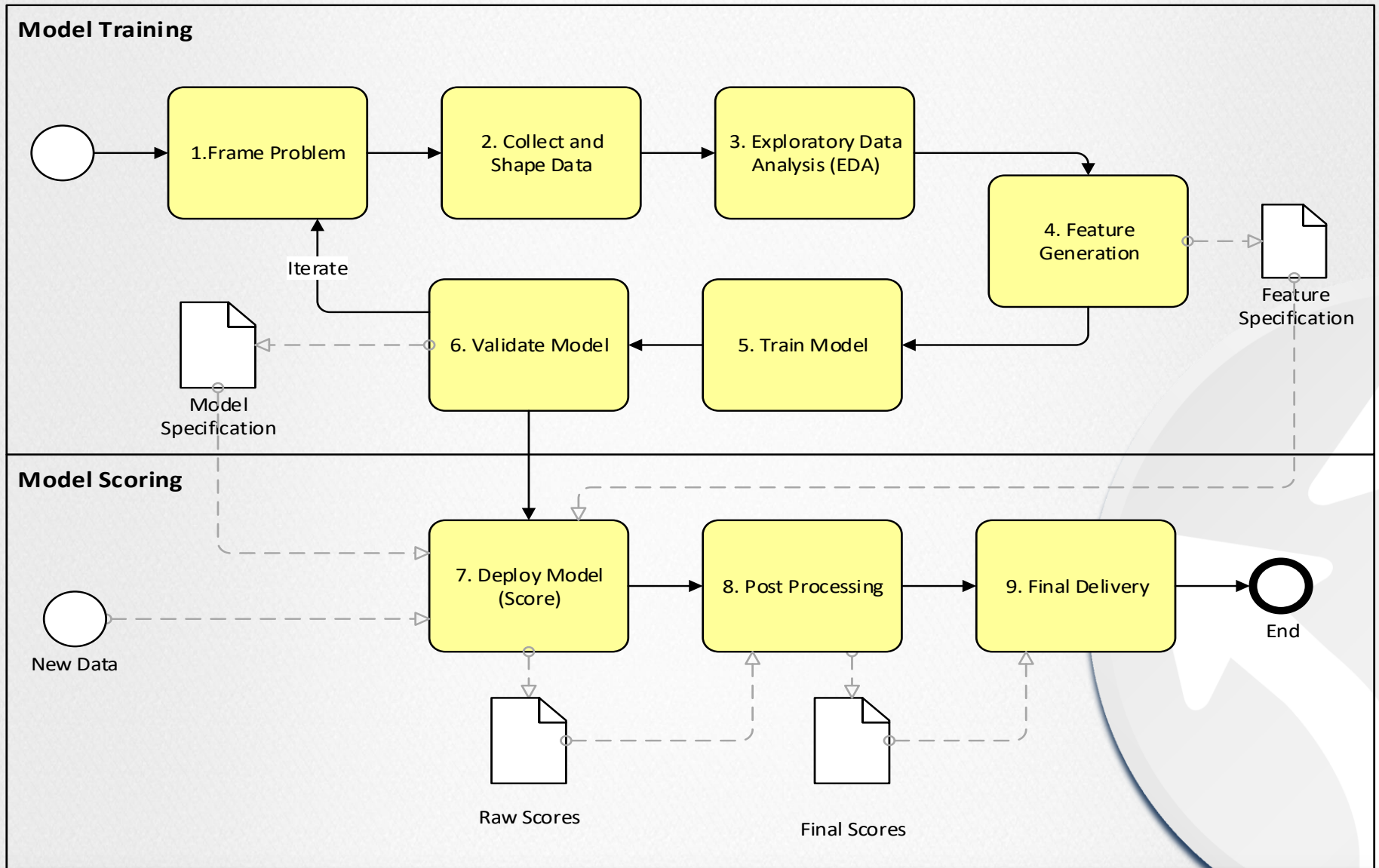
⇒ You understand:

- installed **git** and created a github account
- **fork** the class repository
- **clone** a local copy of the repo
- **pull** new changes
- edit existing files
- **add** and **commit** changes
- **push** the assignment back to your repo

⇒ Now: **pull** upstream changes
`CSX460/CSX460.git`



Expectations: Process



GETTING HELP IN PRIMER

Help in R

`?, help, ??, apropos`

Operators

`?Arithmetic`

Control Flow

`?Control`

Rstudio Cheatsheets ...Google



\hat{y}

Prediction
Forecast
Estimate

...



DATA USES

Dependent variable,
Target (variable),
Outcome, **Response,**
Class (classification)

Independent variables, covariates
predictors, attribute,
descriptor, **feature**

...

Y	X ₁	X ₂	X ₃	... X _n

Unit of observation,
Cases,
Instance,
Data Point,
Sample

MAGRITTR: PIPE OPERATOR

```
install.package('magrittr')
```

```
1:10 %>% mean
```

```
1:10 %>% add(2) %>% mean
```

```
x <- 1:10
```

```
x %<>% add(2) %>% mean
```

Notes:

* Use `backpipe` package for `%<%`



DATA.TABLE: FAST DATA FRAMES

```
install.packages('data.table')  
data(iris)  
setDT(iris)
```

```
iris[ i, j, by= , ... ]
```



Row

Note:

- see ?data.table



DATA USES

Dependent variable,
Target (variable),
Outcome, **Response,**
Class (classification)

Independent variables, covariates
predictors, attribute,
descriptor, **feature**

...

Y	X ₁	X ₂	X ₃	... X _n

Unit of observation,
Cases,
Instance,
Data Point,
Sample

DPLYR: DATA PIPELINES

```
install.packages('dplyr')  
data(iris)
```

```
iris %>%  
  filter( Species != "setosa") %>%  
  group_by(Species) %>%  
  summarize(  
    mean(Sepal.Width),  
    mean(Sepal.Length)  
  )
```

Note:

Uses `magrittr`



BACK TO MACHINE LEARNING



MACHINE LEARNING TYPES

⇒ **Type** of Response:

- Continuous → **REGRESSION**
- Categorical* → **CLASSIFICATION**
*Binary is a special case

⇒ **Availability** of “labelled” Responses

- Available → **SUPERVISED**
- Unavailable → **UNSUPERVISED**
- Sometimes available/inferable → **SEMI-SUPERVISED**
- Avail. as training progresses → **ADAPTIVE/REINFORCEMENT**



GOAL FIND A FUNCTION, f

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

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



QUESTIONS?



QUESTIONS:

What do we mean by “close”?

⇒ What is a quantitative way of measuring closeness?

Compare: $\hat{y} \sim y \quad f(\hat{y}, y)$

But how?

USE all observations/estimates?

$$g(f(\hat{y}, y)) = (g \circ f)(\hat{y}, y)$$



QUESTIONS:

What functions are available to be used?



OUR MODEL

Naïve Model

$$\hat{y} = \text{mean}(y)$$

Our Model, a linear model:

$$\hat{y} = \beta_0 + \beta_1 x_1$$



QUESTIONS:

How do we find one? The best one?



SEARCH / OPTIMIZATION

Find the parameters minimize that minimize the loss function ...

SOLVE:

$$\operatorname{argmin}_{\beta} L(\mathbf{y}, \hat{\mathbf{y}})$$

$$\operatorname{argmin}_{\beta} \sum (\mathbf{y} - \hat{\mathbf{y}})^2 \text{ (SSE)}$$

- Direct Solution (special case)
- Recursive Goal Seeking



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**)

LM / MODEL FORMULA



LINEAR REGRESSION MODEL

Abstract to multiple dimensions

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$$

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

Mathy-r !!!



APPENDIX



MLWR CHAPTER 1

- ⇒ Four Parts to “Learning” Process
- ⇒ Five Steps for Modeling
- ⇒ **Types** of Data
- ⇒ Types of Machine Learning Algorithms



MLWR CHAPTER 2

- ⇒ Data structures
- ⇒ Saving/Loading Data With R
- ⇒ Exploring the structure of the Data
 - Numeric variables
 - Categorical variables
 - Relationship Between Variables



QUESTIONS:

- What do we mean by “close”?
- What functions are available to be used?



- How do we find one? The best one?

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**)