

Machine Learning Guidelines for Natural Resource Management Practitioners

Shih-Ni Prim and Natalie Nelson

2024-02-18

Contents

1	Motivation	5
2	Introduction	7
2.1	Supervised Learning	7
2.2	Unsupervised Learning	7
3	Data	9
3.1	Data Exploratoary Analysis	9
3.2	Data Requirement	13
4	Evaluation	15
4.1	Training vs Testing	15
4.2	Metrics	15
4.3	Cross Validation	15
5	Machine Learning Methods	17
5.1	Tree-based Methods	17
6	Presentation	19
6.1	Table	19
6.2	Figure	19
7	Ethical Considerations	21
7.1	Reproducibility	21
7.2	Decision making	21

8	Appendix	23
8.1	Do's and Don'ts	23

Chapter 1

Motivation

As machine learning (ML) has become a powerful tool, it is noted by some that ML has not been widely used in environmental studies. This booklet is meant to provide a concise guide for natural resource management practitioners. This book serves as a starting point rather than a comprehensive resource, so that practitioners can have a basic understanding of how ML works and how to utilize it to analyze data and answer research questions. When appropriate, we provide case studies and R code as well as other online resources to help the readers on the journey of gaining one powerful tool that seems to be omnipresent in the research world.

Chapter 2

Introduction

What is machine learning? Essentially, machine learning teaches computer models to look for patterns or make predictions. This might sound like magic or it might seem complicated, but you can think of machine learning models as finding underlying formulas that the data come from. To solve for such formula, many, many mathematical calculations are involved. As we human beings are prone to mistakes, as long as we can identify a framework, we can give the framework and data to a computer model. It is best at repeating meticulous calculations to find a best guess based on our believes of the system and the data we observed.

2.1 Supervised Learning

2.2 Unsupervised Learning

Chapter 3

Data

When you have data, don't feel so rushed to jump into data analysis yet. Some steps can help you know your data better and, more than often, avoid problems down the road.

3.1 Data Exploratory Analysis

After you read in data, do some checks. Below we use the embedded `mtcars` as an example for illustration.

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.4
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2    3.4.4      v tibble    3.2.1
## v lubridate  1.9.3      v tidyr     1.3.0
## v purrr      1.0.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
data("mtcars")
str(mtcars)
```

```
## 'data.frame':   32 obs. of  11 variables:
## $ mpg : num  21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...
```

```
## $ cyl : num  6 6 4 6 8 6 8 4 4 6 ...
## $ disp: num 160 160 108 258 360 ...
## $ hp  : num 110 110 93 110 175 105 245 62 95 123 ...
## $ drat: num 3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...
## $ wt  : num 2.62 2.88 2.32 3.21 3.44 ...
## $ qsec: num 16.5 17 18.6 19.4 17 ...
## $ vs  : num 0 0 1 1 0 1 0 1 1 1 ...
## $ am  : num 1 1 1 0 0 0 0 0 0 0 ...
## $ gear: num 4 4 4 3 3 3 3 4 4 4 ...
## $ carb: num 4 4 1 1 2 1 4 2 2 4 ...
```

As seen above, the command `str` allows you to see the variables and their types. For this dataset, all variables are numerical. If some are categorical but they should be numerical, make sure you transform them into the right type of data. (Sometimes numbers can be saved as characters, and the analysis would not be correct if the datatype remains as character.)

Next, try to make some plots—typically histograms for numerical variables and barplots for categorical variables.

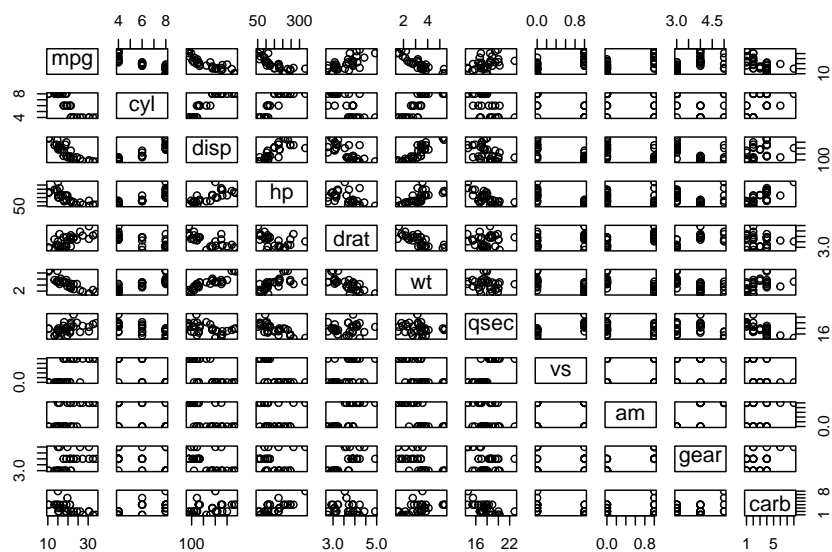
```
par(mfrow = c(3,3), mar = c(2,2,2,2))
for (i in 1:ncol(mtcars)){
  hist(mtcars[,i], main = paste0("Histogram of ", colnames(mtcars)[i]))
}
```





You can also look at the paired plots to see if two variables are too perfectly correlated, which could cause problems in regression models.

```
pairs(mtcars)
```



Next, take a look at the summary statistics.

```
for (i in 1:ncol(mtcars)){
  print(paste0("***** Summaries of ", colnames(mtcars)[i], " *****"))
  print(summary(mtcars[,i]))
}
```

```
## [1] "***** Summaries of mpg *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  10.40   15.43   19.20   20.09   22.80   33.90
## [1] "***** Summaries of cyl *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   4.000   4.000   6.000   6.188   8.000   8.000
## [1] "***** Summaries of disp *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   71.1   120.8   196.3   230.7   326.0   472.0
## [1] "***** Summaries of hp *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   52.0   96.5   123.0   146.7   180.0   335.0
## [1] "***** Summaries of drat *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   2.760   3.080   3.695   3.597   3.920   4.930
## [1] "***** Summaries of wt *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   1.513   2.581   3.325   3.217   3.610   5.424
## [1] "***** Summaries of qsec *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   14.50   16.89   17.71   17.85   18.90   22.90
## [1] "***** Summaries of vs *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   0.0000   0.0000   0.0000   0.4375   1.0000   1.0000
## [1] "***** Summaries of am *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   0.0000   0.0000   0.0000   0.4062   1.0000   1.0000
## [1] "***** Summaries of gear *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   3.000   3.000   4.000   3.688   4.000   5.000
## [1] "***** Summaries of carb *****"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   1.000   2.000   2.000   2.812   4.000   8.000
```

The summaries can show, for example, ranges and means of the variables. They can give you a better understanding of where the values are and whether there might be data entry errors.

3.2 Data Requirement

Chapter 4

Evaluation

4.1 Training vs Testing

We should first address the concepts of training and testing.

4.2 Metrics

4.2.1 Continuous Responses

4.2.2 Discrete Responses

4.3 Cross Validation

One very standard way of evaluation is k -fold cross validation, commonly with $k = 5$ or $k = 10$. The idea is simple. Divide the data into k groups. Each time, choose $k - 1$ groups for training, fit the model on the last group, which is the test data, and calculate the desired metrics, such as MSE.

In this way, although less data is used for training, the metrics are more accurate, because now we are not using the same data points for training and testing. Using metrics from cross validation for model selection can ensure that your model does not overfit, which means the model does well with training data but does not generalize well on new data.

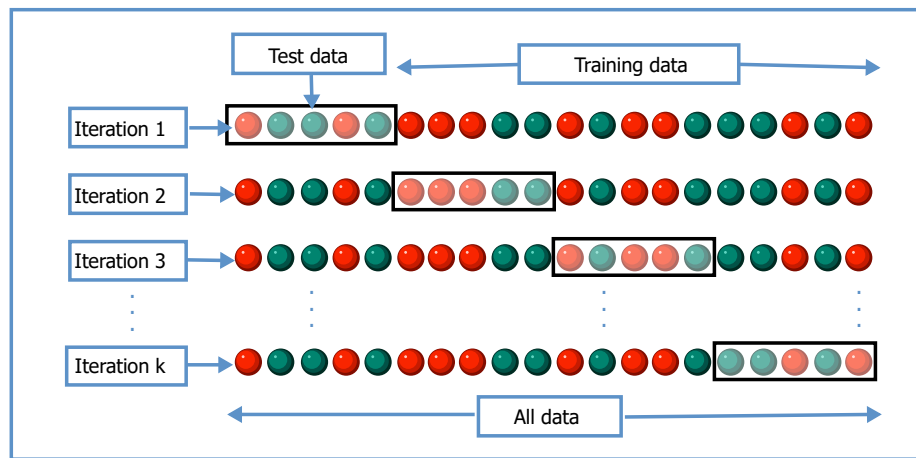


Figure 4.1: Image Source: [https://en.wikipedia.org/wiki/Cross-validation_\(statistics\)](https://en.wikipedia.org/wiki/Cross-validation_(statistics))

Chapter 5

Machine Learning Methods

Here we provide a list of commonly used machine learning methods and some brief discussion.

5.1 Tree-based Methods

Tree-based methods are popular for their ease for interpretation.

5.1.1 Random Forest

Survival of passengers on the Titanic

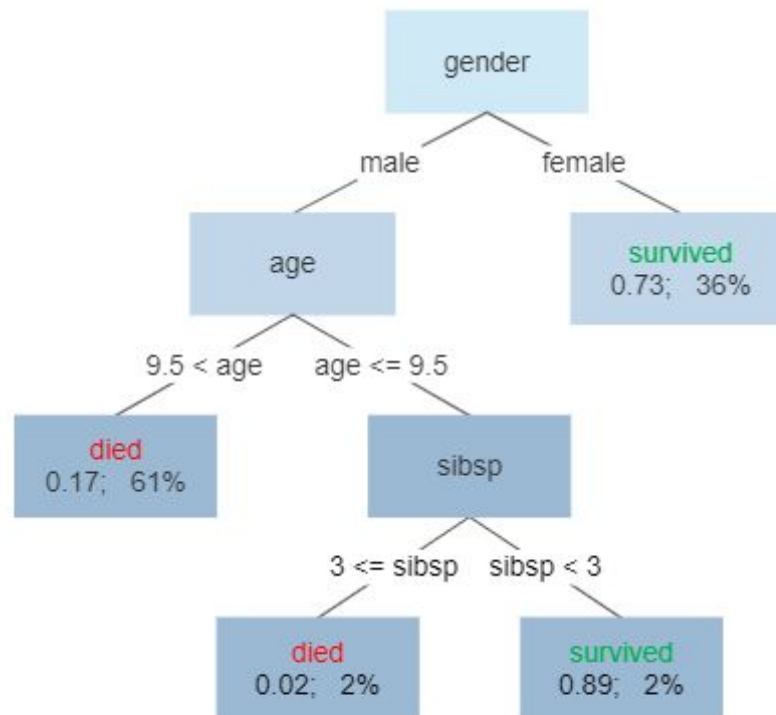


Figure 5.1: Image Source: https://en.wikipedia.org/wiki/Decision_tree_learning

Chapter 6

Presentation

It is also important to present the results in a way that aids rather than impede communication.

6.1 Table

6.2 Figure

Chapter 7

Ethical Considerations

7.1 Reproducibility

To allow for others to reproduce your work, it is important to provide enough details in terms of methods, data processing, code implementation, etc. It is also encouraged to have all the code and data available online in a repository. If parts or all of the data should not be shared publicly, it helps to provide a simulated data set.

7.2 Decision making

Since research related to environmental sciences and natural resources could likely affect decision making, we will now address some topics.

7.2.1 Uncertain qualification

While all ML models can provide point estimates, not all can quantify uncertainty. It is, however, important to show how confident the model is about the estimates. If the conclusion of the study could affect an important policy change, it is crucial to present a full picture of the findings, which include uncertainty quantification. It is wildly different whether the model is 20% or 95% confident about its answer, for instance.

7.2.2 Interpretability

While some models, such as neural networks, are highly efficient, they are more like black boxes and do not lend easily to interpretability. In the case of neural

networks, even if you are able to find all the weights in the hidden layers, there is really no way to interpret them. To ensure that the model arrives at a reasonable conclusion, you might consider using a model that is more interpretable, such as linear regression or tree-based methods. This way, experts with domain knowledge can examine whether the conclusion makes sense. In other words, the findings from ML models can add to researchers' understanding of the field rather than throwing out an answer that is not easily interpreted.

Chapter 8

Appendix

8.1 Do's and Don'ts