# Practical Machine Learning in R

## Introduction

Lars Kotthoff[1][2]
larsko@uwyo.edu

---

# What is Machine Learning?

▷ "gives computes the ability to learn without being explicitly programmed" (Wikipedia)

# What is Machine Learning?

▷ "gives computes the ability to learn without being explicitly programmed" (Wikipedia)

▷ "A computer program is said to learn from experience $E$ with respect to some class of tasks $T$ and performance measure $P$ if its performance at tasks in $T$, as measured by $P$, improves with experience $E$." (Tom Mitchell)
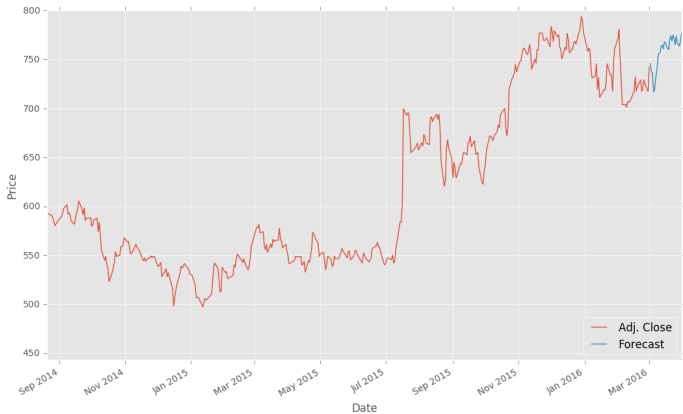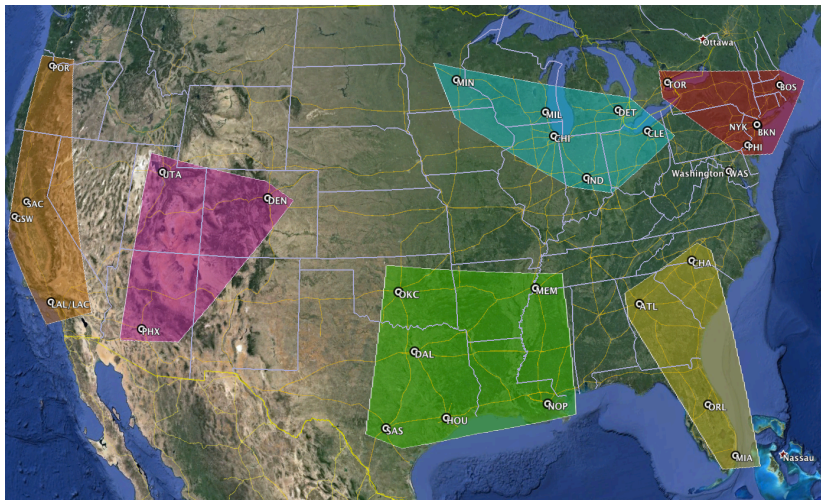
# Examples

# Examples



https://pythonprogramming.net/forecasting-predicting-machine-learning-tutorial/

# Examples



https://squared2020.com/2015/09/09/redefining-nba-divisions-by-clustering/

# Supervised Learning

▷ learn the relationship between input $x$ and output $y$

▷ training data with labels available – $y$ known for given $x$

▷ can see this as function approximation – find an $f$ such that

$$y \approx f(x)$$

# Supervised Learning

- $\triangleright$ $x$ are features or attributes
- $\triangleright$ $y$ is the ground truth
- $\triangleright$ denote predictions $f(x) = \hat{y}$
- $\triangleright$ loss function $L(y, \hat{y})$ measures how good predictions are, e.g.

$$L(y, \hat{y}) = (y - \hat{y})^2$$

- $\triangleright$ want to minimize loss given training data $X_{\text{train}} = \{(x_i, y_i)\}^n$:

$$\arg\min \sum_{i=1}^{n} L(y_i, \hat{y}_i)$$

# Supervised Learning

▷ want to learn a general function that is predictive on new data

▷ second set $X_{test}$ that is not used in training to test generalization performance:
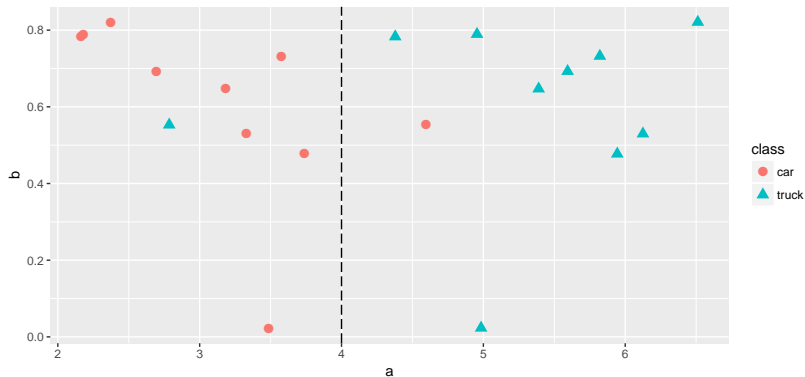
$$\sum_{i=1}^{n} L(y_i, \hat{y}_i)$$

▷ usually full data set $X$ is split into non-overlapping train and test sets:
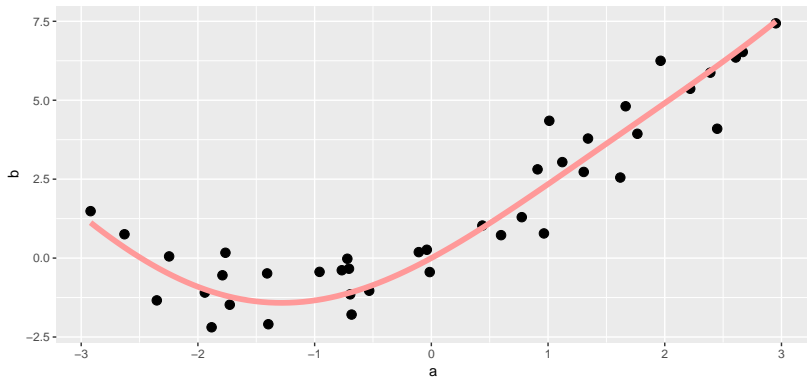
$$X_{train} \cup X_{test} = X$$

$$X_{train} \cap X_{test} = \varnothing$$

# Supervised Classification



Goal: Predict a class (discrete quantity), or membership probabilities

# Supervised Regression

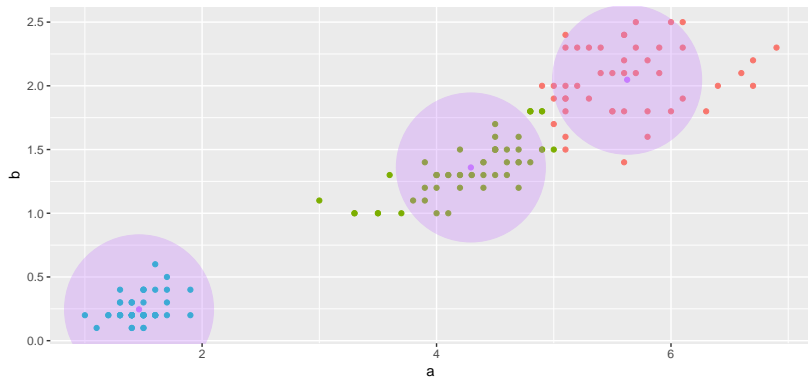

Goal: Predict a continuous quantity

# Unsupervised Learning

▷ no ground truth $y$ available

▷ determine group membership or assign labels

▷ loss function measures properties of groups, e.g. homogeneity wrt. features

▷ still want to minimize loss given training data and generalize

# Unsupervised Clustering



Goal: Group data by similarity, or estimate membership probabilities

# In this Course

- ▷ classification
- ▷ regression
- ▷ clustering
- ▷ data preprocessing (missing values, dimensionality reduction)
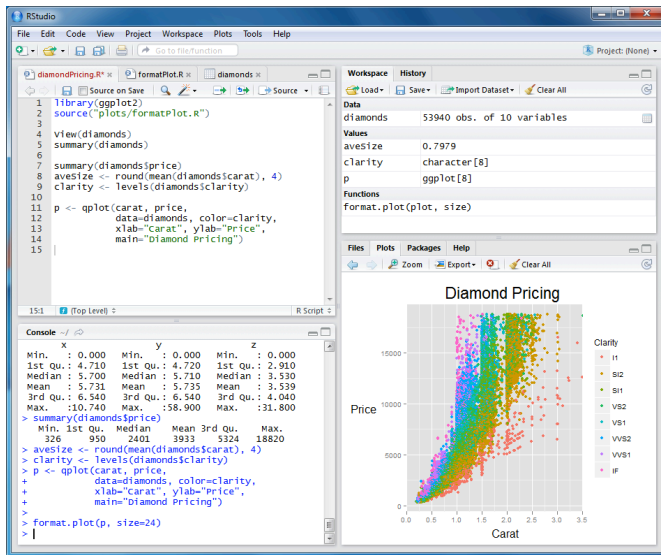- ▷ performance evaluation
- ▷ parameter tuning

# Not in this Course

- ▷ R tutorial
- ▷ details on particular methods
- ▷ deep learning
- ▷ time series
- ▷ Big Data

# What you'll need

# Install RStudio



https://www.rstudio.com/products/rstudio/download/

# Install mlr

▷ on the R console:

  **install.packages**("mlr")

▷ or see http://derekogle.com/IFAR/supplements/
  installations/InstallPackagesRStudio.html

▷ extensive tutorial available: https://mlr-org.github.
  io/mlr-tutorial/devel/html/

# Format

- ▷ meetings roughly every week
- ▷ half lecture, half practical exercises
- ▷ happy to discuss specific problems