

Introduction to Machine- Learning Methods for Patient- Reported Outcomes Data

Yuelin Li

October 19, 2022

ISOQOL Conference, Prague

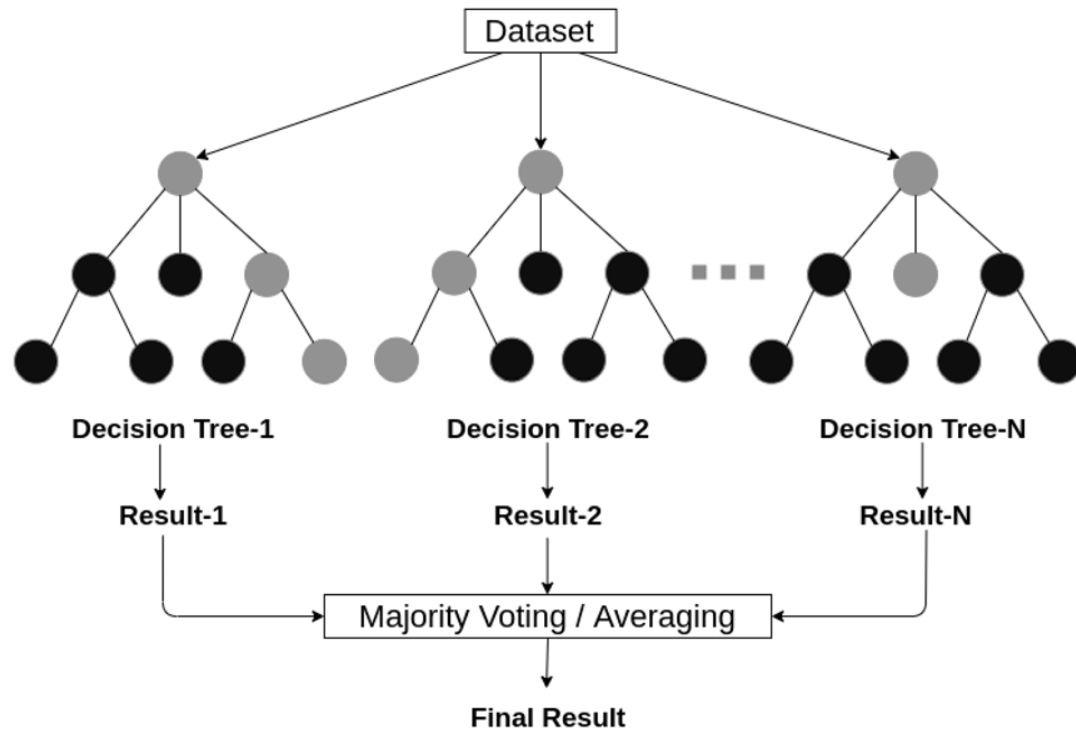
Machine Learning Methods for PRO Data

My goal is to provide a general overview of ML methods

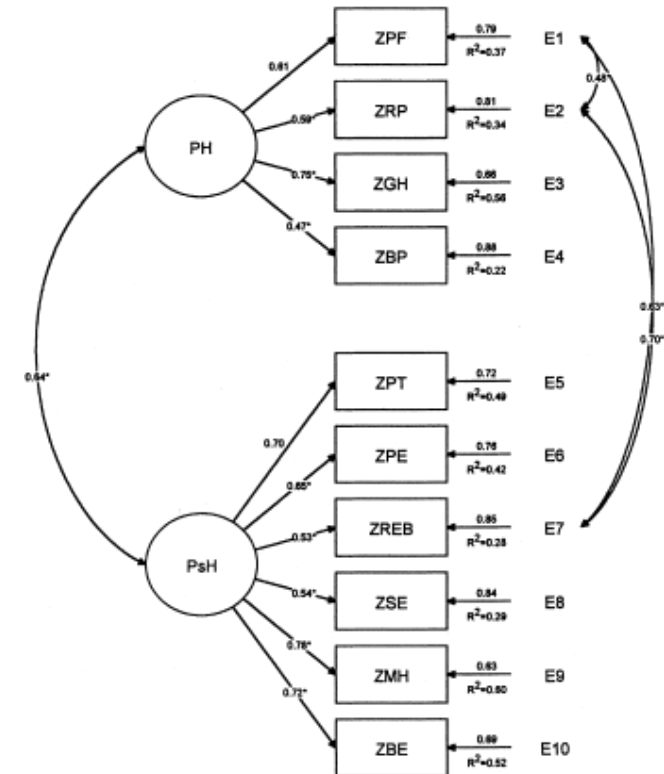
1. Why use ML in PRO data?
2. What is a typical and basic ML workflow?
3. Examples to illustrate the ML workflow
 - Tree-based method
 - Text classification by BERT, a Natural Language Processing (NLP) method
4. How to avoid common pitfalls?

Why use ML in PRO Data?

ML: *person-oriented* approach ¹



Traditional, *variable-oriented* approach



1: Bergman & Magnusson (1997) *Development and Psychopathology* <https://doi.org/10.1017/S095457949700206X>

Machine Learning: What is it?

- **Unsupervised Learning**
- **Supervised Learning**

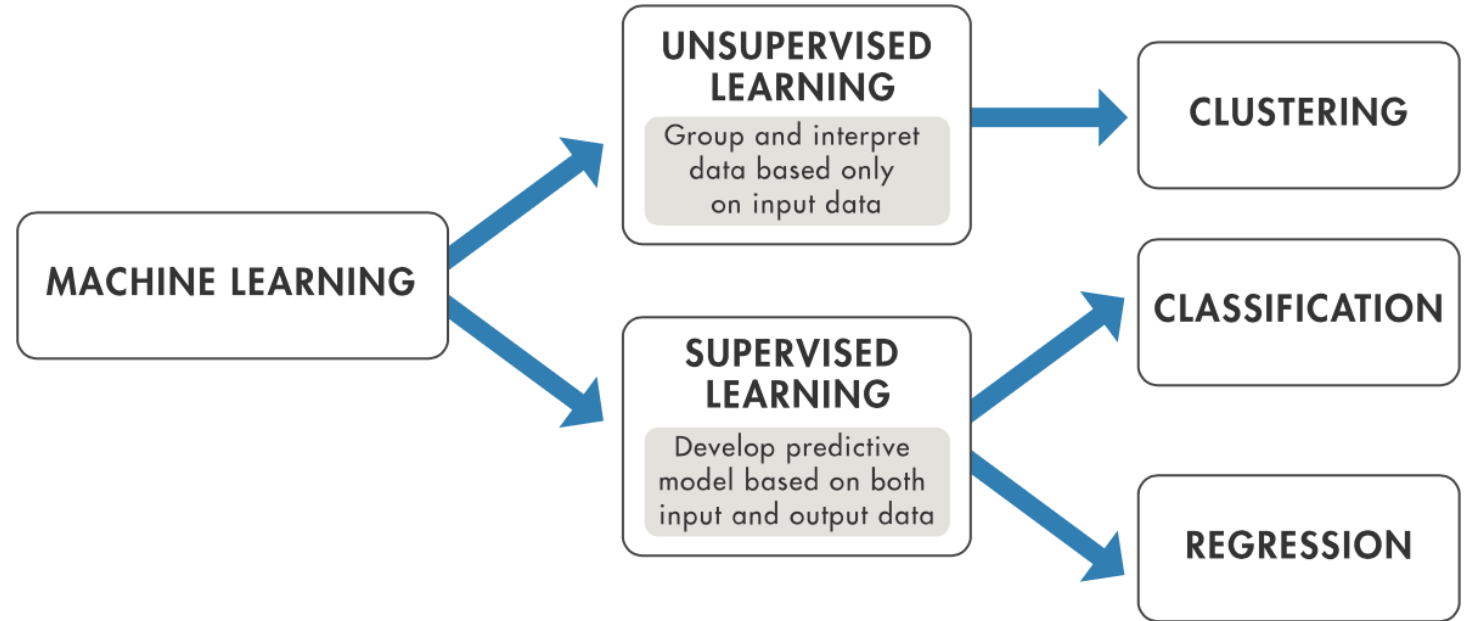


Figure 1. Machine learning techniques include both unsupervised and supervised learning.

<https://www.mathworks.com/discovery/machine-learning.html>

Unsupervised Learning

- Unsupervised learning finds hidden patterns or intrinsic structures in input data
- Clustering
 - K-means
 - Finite mixture modeling
 - Bayesian nonparametric methods ¹
 - NLP (e.g., Latent Dirichlet Allocation ²)
 - Psychometric network models ⁴

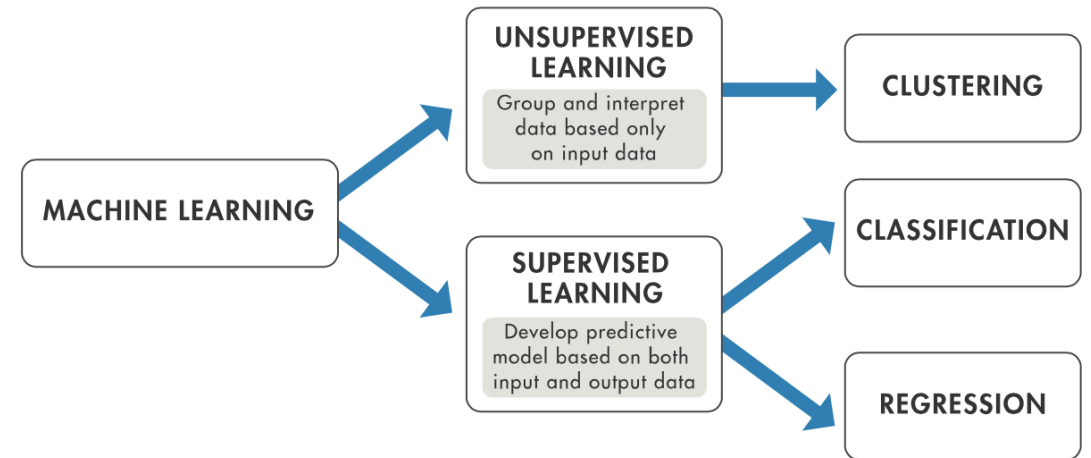


Figure 1. Machine learning techniques include both unsupervised and supervised learning.

Supervised Learning

- Supervised machine learning builds a model that makes predictions based on evidence in the presence of uncertainty
 - Classification methods
 - Classification and Regression Trees (CART)¹
 - SVM (Support Vector Machine)
 - Random Forest
 - XGBoost (eXtreme Gradient Boosting)
 - Artificial neural networks ²
 - NLP (e.g., BERT for Text Classification ³)
 - Regression and its cousins

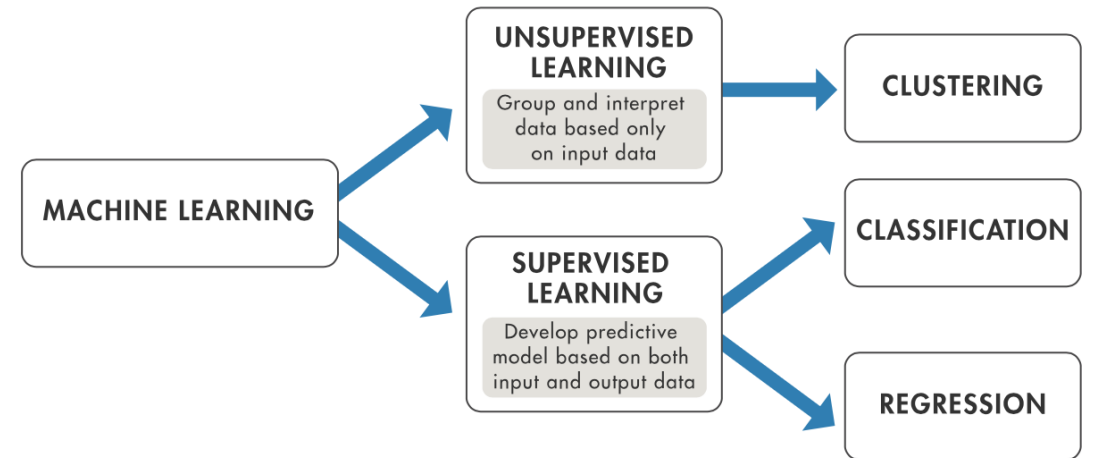
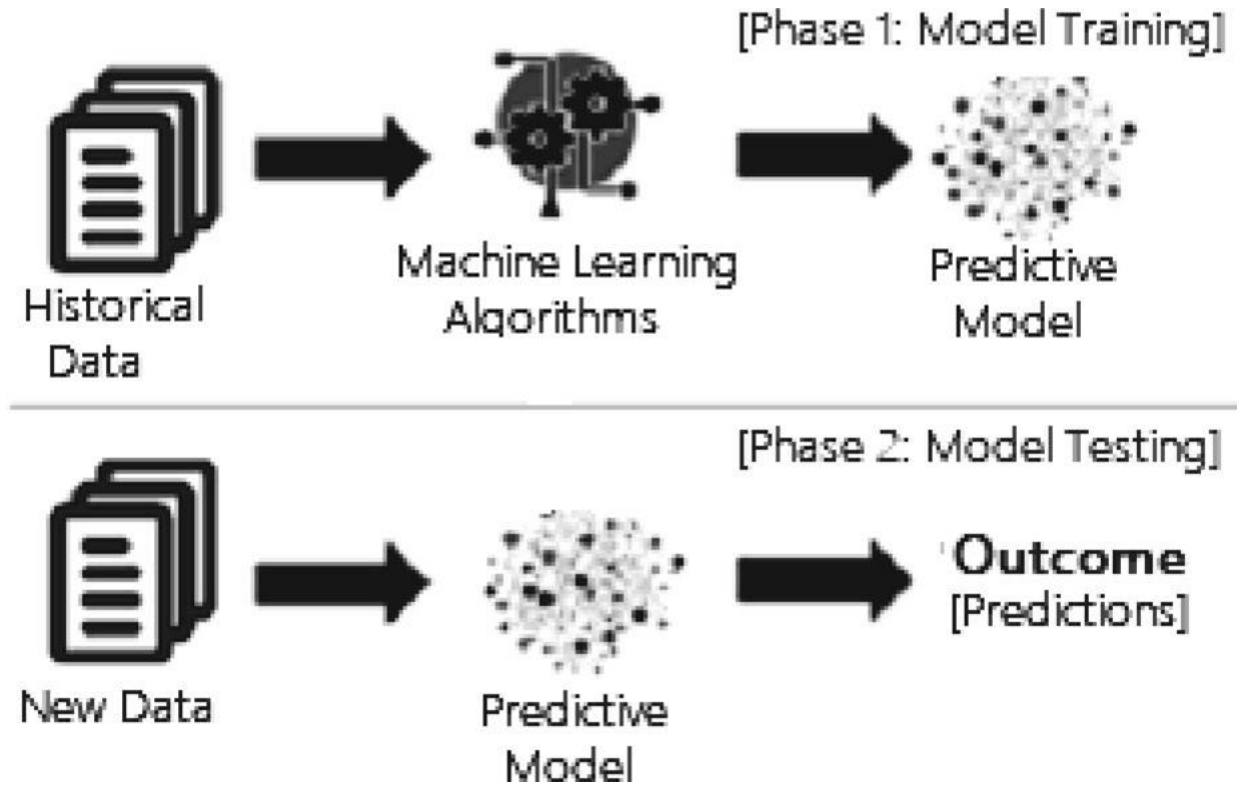


Figure 1. Machine learning techniques include both unsupervised and supervised learning.

1: Li & Rapkin (2009, PMID: 19595576); 2. Pfof et al. (2021, PMID: 33914464); 3. Schwartz et al. (2022, JMSS)

What is a typical ML workflow?

- Model training: feed historical data (often large) to ML algorithms to build a predictive model (e.g., logistic regression)
- Model testing: new data (never seen before by predictive model) fed into the model to generate predictions
- Predictive accuracy in new data (e.g., sensitivity/specificity, ROC curves, mean squared error)
- Often, as part of phase 1 model training, you do *cross-validation* to select the desired model



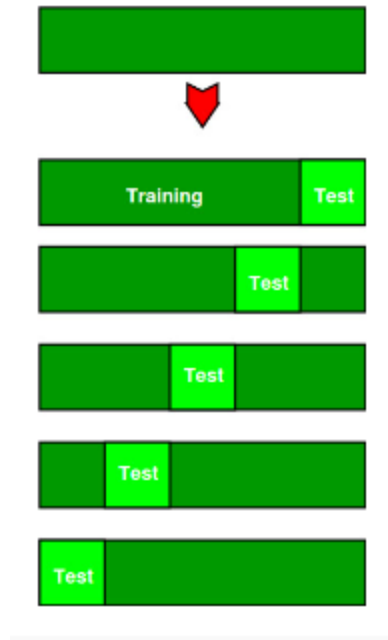
Sarker (2021)

Typical ML workflow

- Training sample, validation sample, and test sample
- Training sample to set the parameters of the predictive models ("train the models" in ML parlance)
- Validation to choose the best model among different models trained on the training sample
- How choose? By selecting (say) model that has the lowest empirical risk on the validation sample
- E.g., choice between a simple vs. a complex model using same training sample. We choose the model that has the lowest mean squared prediction error on the validation sample
- Test sample yields unbiased estimate of risk of model selected in the validation step
- Why test step? Because model selection can also make mistakes
- Sometimes tested by someone else on blinded data (ML competitions)
- Workflow is iterative (e.g., Train-Validation-Test1-Test2)

Cross validation types

- K-fold cross-validation



- Leave-One-Out (LOO) cross validation

- N-fold cross validation

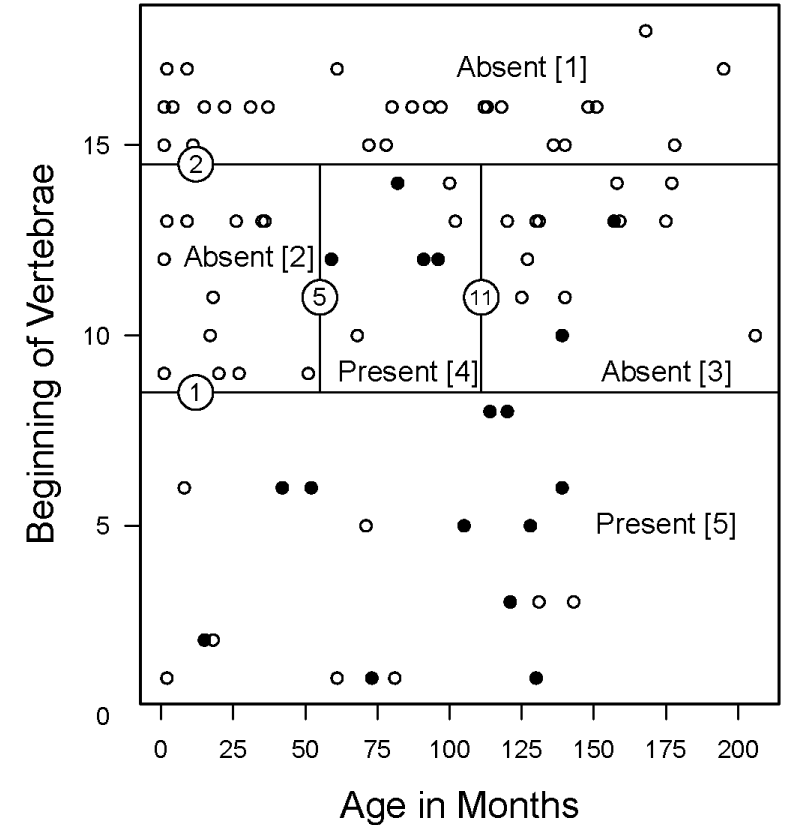
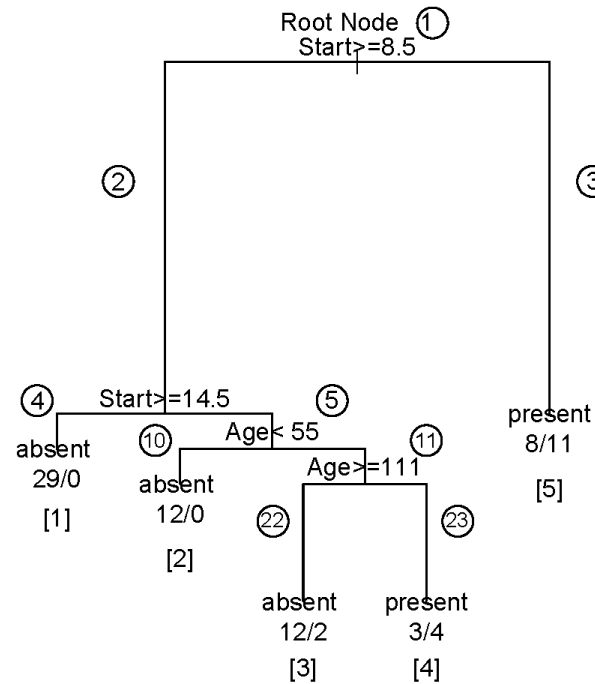
$$\text{MSE} = \frac{1}{N} \sum (y_i - f(x_i))^2$$

- N : Total number of observations
- y_i : Observed outcome value of the i^{th} observation
- $f(x_i)$: Predicted response value of the i^{th} observation

Cross-Validation Example

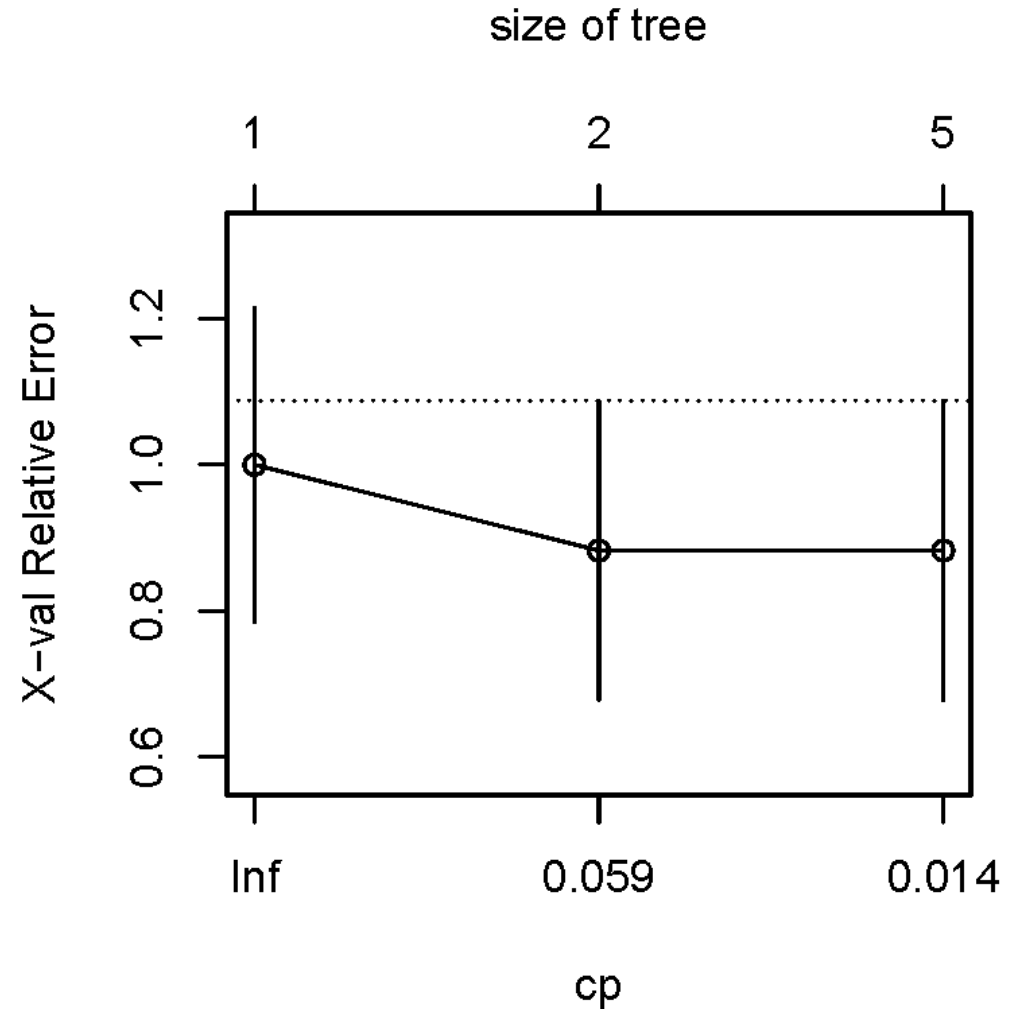
Why x-validation?

- For example, stock Kyphosis dataset
- Recursive partitioning divides 2D space between Age and Start into subgroups
- However, individuals falling near the boundaries can be misclassified by chance
- Cross-validation error provides a measure of this problem by running the same model K times
- Now, we will do this in RStudio



Prune rpart tree

- On the previous slide, `rpart()` grows a tree with 5 terminal nodes
- `rpart()` does a 10-fold cross-validation under the hood
- The `plotcp()` command plots the size of tree and its corresponding relative error to help prune the tree
- Like a traditional scree plot
- Splitting sample to 2 nodes reduces relative error (good)
- Further into 5 nodes not much worse than 2
- Rationale for a 5-node model



Natural Language Processing

Natural Language Processing by BERT

- BERT stands for ‘Bidirectional Encoder Representations from Transformers’
- Unlike earlier NLP models, BERT’s artificial neural nets take the order of words into account
- E.g., ‘work to live’ and ‘live to work’ are semantic opposites
- In late 2018, Google released BERT ¹
- Other technology companies such as Facebook and OpenAI have also joined forces to pre-train BERT using 3.3 billion words total with 2.5 billion words from Wikipedia and 0.8 billion from BooksCorpus.
- As a result, this “pre-trained” BERT has learned how the English language is structured and how the words are typically arranged in the context of sentences
- Learn more in their blog post ²
- Carolyn Schwartz and I tested BERT in coding interview transcripts (close to the performance of a human coder) ³
- BERT is now an integral part of Google search
- You might have noticed that you get better Google search results by complete sentences

Text Classification by BERT

- Next, we turn to Python to show you how to do text classification using a “pre-trained” BERT you can download

To Summarize This Part of Workshop

- Countless other points on ML not covered here
- But I hope I have given you a fundamental intuition on ML
- Importance of cross-validation and NLP parameter tuning as an integral part of data analysis in writing a paper
- You must interrogate your models (e.g., cross-validation), not just to accept ML defaults
- Make you feel more confident and efficient in learning more on your own
- Potentially rediscover finer details of ML yourself