# Overview of Hands-On Example & XGBoost

Nick Henscheid
Sr. Quantitative Medicine Scientist

# Agenda

- Overview of hands-on portion of the webinar (Thurs)
- Code installation tips
- Introduction to the practicum dataset: predicting breast cancer recurrence
- Conceptual & practical explanation of XGBoost for classification

# Hands-on Practicum: Overview

- This Thursday from 11AM – 1PM EDT there is a hands-on coached tutorial session

- We will have two coaches per breakout room

- Instructions to download and install the code are in your inbox

- The code is written in R

- The GitHub Repository is:

  https://github.com/metrumresearchgroup/ascpt-ml-webinar-2023


*Please allow approximately one hour to install the necessary packages and test the code*

*If you have issues, please email one of the coaches*

*Start after today's webinar to allow enough time to resolve issues before Thursday*

# Hands-on Practicum: Installation Tips

QUANTITATIVE
MEDICINE

- Make sure you have R 4.1.3 or later and appropriate version of Rtools installed

- The .Rproj file provides an isolated environment to install packages

- Packages required:

    *knitr, tidyverse, tidymodels, xgboost, here*

- Use RStudio package installer or install.packages()

- The dataset is contained in the /data folder

- Test package installation with /script/test-script.R

- The workbook and answer key are in /script

📁 data

📁 renv

📁 script

📄 .gitignore

📄 .Rprofile

📄 ascpt-ml-webinar.Rproj

📄 pkgr.yml

📄 README.md

📄 renv.lock

# Example Dataset: UCI Breast Cancer
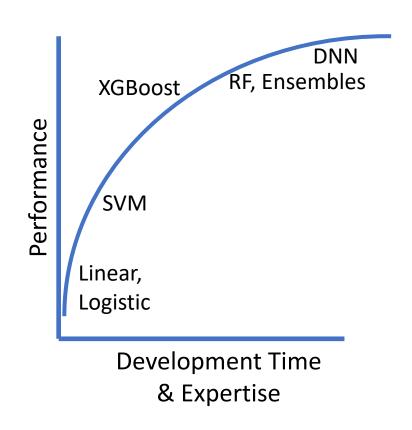
https://archive.ics.uci.edu/dataset/14/breast+cancer

| Variable Name | Role | Type | Values | Notes |
|---|---|---|---|---|
| Class | Target | Binary | no-recurrence-events, recurrence-events | |
| age | Feature | Categorical | 10-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80-89, 90-99 | |
| menopause | Feature | Categorical | lt40, ge40, premeno | |
| tumor-size | Feature | Categorical | 0-4, 5-9, 10-14, 15-19, 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, 50-54, 55-59 | |
| inv-nodes | Feature | Categorical | 0-2, 3-5, 6-8, 9-11, 12-14, 15-17, 18-20, 21-23, 24-26, 27-29, 30-32, 33-35, 36-39 | Number of lymph nodes showing metastases |
| node-caps | Feature | Binary | yes, no | Metastasis contained to lymph cap |
| deg-malig | Feature | Integer | 1, 2, 3 | Histological grade |
| breast | Feature | Binary | left, right | |
| breast-quad | Feature | Categorical | left-up, left-low, right-up, right-low, central | |
| irradiat | Feature | Binary | yes, no | Patient received rad therapy |

| Class | age | menopause | tumor-size | inv-nodes | node-caps | deg-malig | breast | breast-quad | irradiat |
|---|---|---|---|---|---|---|---|---|---|
| no-recurrence-events | 60-69 | ge40 | 25-29 | 0-2 | no | 3 | right | left_low | no |
| recurrence-events | 50-59 | premeno | 15-19 | 0-2 | no | 2 | left | left_low | no |

# XGBoost: Overview

- XGBoost is a powerful *method* and a *software framework*

- It is ideal for large tabular datasets

- XGBoost can perform both classification and regression

- The algorithm is based on regularized gradient boosting but has added features to improve performance, reliability and usability

- XGBoost has many hyperparameters…

- …but tidymodels can help

- In terms of performance versus labor, XGBoost ranks highly: competitive performance across many tasks and datasets with minimal development
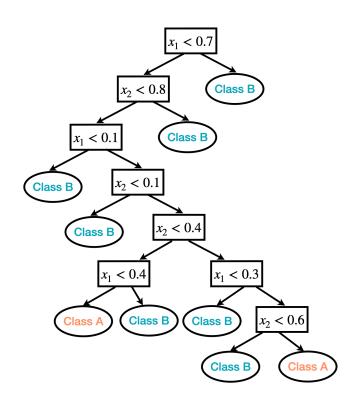
# XGBoost: Algorithm Architecture

- XGBoost is based on the idea of an *ensemble*:

$$f(x) = \sum f_j(x)$$

- Each individual model ($f_j(x)$) is 'weak', but combining them results in a model with good performance

- The ensemble models in XGBoost are trees (by default the first model predicts 0.5 for everything).

- XGBoost trees are built iteratively by solving a regularized optimization problem designed to reduce errors made by previous model

- Trees are also pruned to reduce overfitting

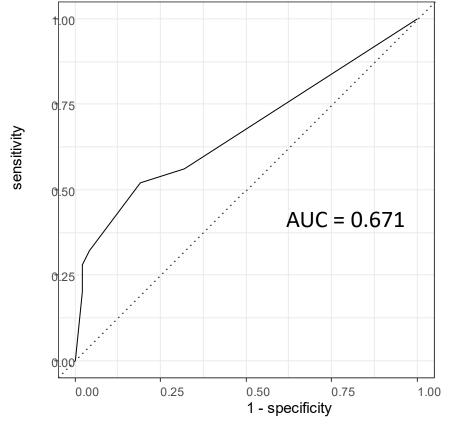- XGBoost has native parallel & distributed capabilities
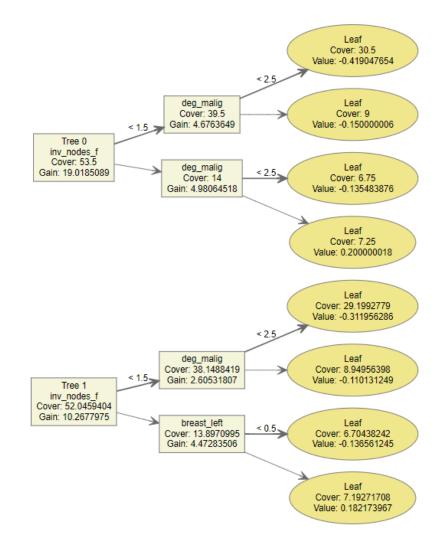
# XGBoost: Hyperparameters

- XGBoost has an array of hyperparameters that influence how many trees are built, how deep the trees are, and how aggressively to regularize

- The full list of hyperparameters is found in the documentation; here are a few:
  - nrounds: controls the number of trees to build
  - max.depth: controls the maximum depth of each tree
  - eta: controls the weight that each tree is given in the final prediction (aka the learning rate)
  - gamma: limits the growth of each tree by eliminating branches that might overfit

```
fit2 <- xgboost::xgboost(data = x_train,   # Predictors
                         label = y_train,  # Responses
                         objective="binary:logistic", # classification
                         nrounds = 5,      # Number of trees to build
                         max.depth=10,     # Maximum depth of each tree
                         eta = 0.3,        # Learning rate aka Shrinkage
                         gamma = 2         # Tree growth limiter
                         )
```

# XGBoost: Hyperparameter examples



```
fit2 <- xgboost::xgboost(data = x_train,   # Predictors
                         label = y_train,  # Responses
                         objective="binary:logistic", # classification
                         nrounds = 2,      # Number of trees to build
                         max.depth=2,      # Maximum depth of each tree
                         eta = 0.3,        # Learning rate aka Shrinkage
                         gamma = 0         # Tree growth limiter
                         )
```

AUC = 0.671

# XGBoost: Hyperparameter examples

```
fit2 <- xgboost::xgboost(data = x_train,    # Predictors
                         label = y_train,   # Responses
                         objective="binary:logistic",  # classification
                         nrounds = 2,       # Number of trees to build
                         max.depth=5,       # Maximum depth of each tree
                         eta = 0.3,         # Learning rate aka Shrinkage
                         gamma = 0,         # Tree growth limiter
                         )
```



AUC = 0.643

```
fit2 <- xgboost::xgboost(data = x_train,    # Predictors
                         label = y_train,   # Responses
                         objective="binary:logistic", # classification
                         nrounds = 5,       # Number of trees to build
                         max.depth=10,      # Maximum depth of each tree
                         eta = 0.3,         # Learning rate aka Shrinkage
                         gamma = 3          # Tree growth limiter
                         )
```

AUC = 0.701

# XGBoost: Parameter Tuning with tidymodels

QUANTITATIVE MEDICINE

XGBoost can be challenging to tune, but tidymodels makes this easy.
See Julia Silge's video https://www.youtube.com/watch?v=hpudxAmxHSM
and blog post https://juliasilge.com/blog/xgboost-tune-volleyball/

- Number of trees
- Maximum tree depth
- Minimum # of data points to split
- Minimum loss reduction to split
- % of samples for training
- % of predictors sampled
- Learning rate

```
281   xgb_spec_all <- boost_tree(
282       trees = 25,
283       tree_depth = tune(),
284       min_n = tune(),
285       loss_reduction = tune(),
286       sample_size = tune(),
287       mtry = tune(),
288       learn_rate = tune()
289   ) %>%
290       set_engine("xgboost") %>%
291       set_mode("classification")
```

# Additional Resources

- https://xgboost.readthedocs.io/en/stable/index.html

- XGBoost on StatQuest

- hackerearth tutorial

- XGBoost Presentation by Tianqi Chen (one of the XGBoost developers)

- List of examples and challenge wins with XGBoost

- In-depth explanation of XGBoost (with math)