

A large, round metal bowl is filled with cracked Boston crab. The crab pieces are bright red and appear to be steamed or boiled. They are garnished with finely chopped green onions. The bowl is set against a blurred background of a wooden table.

# Cracking Boston Crab

Presented by Thomas \$ @ 2018-09-07



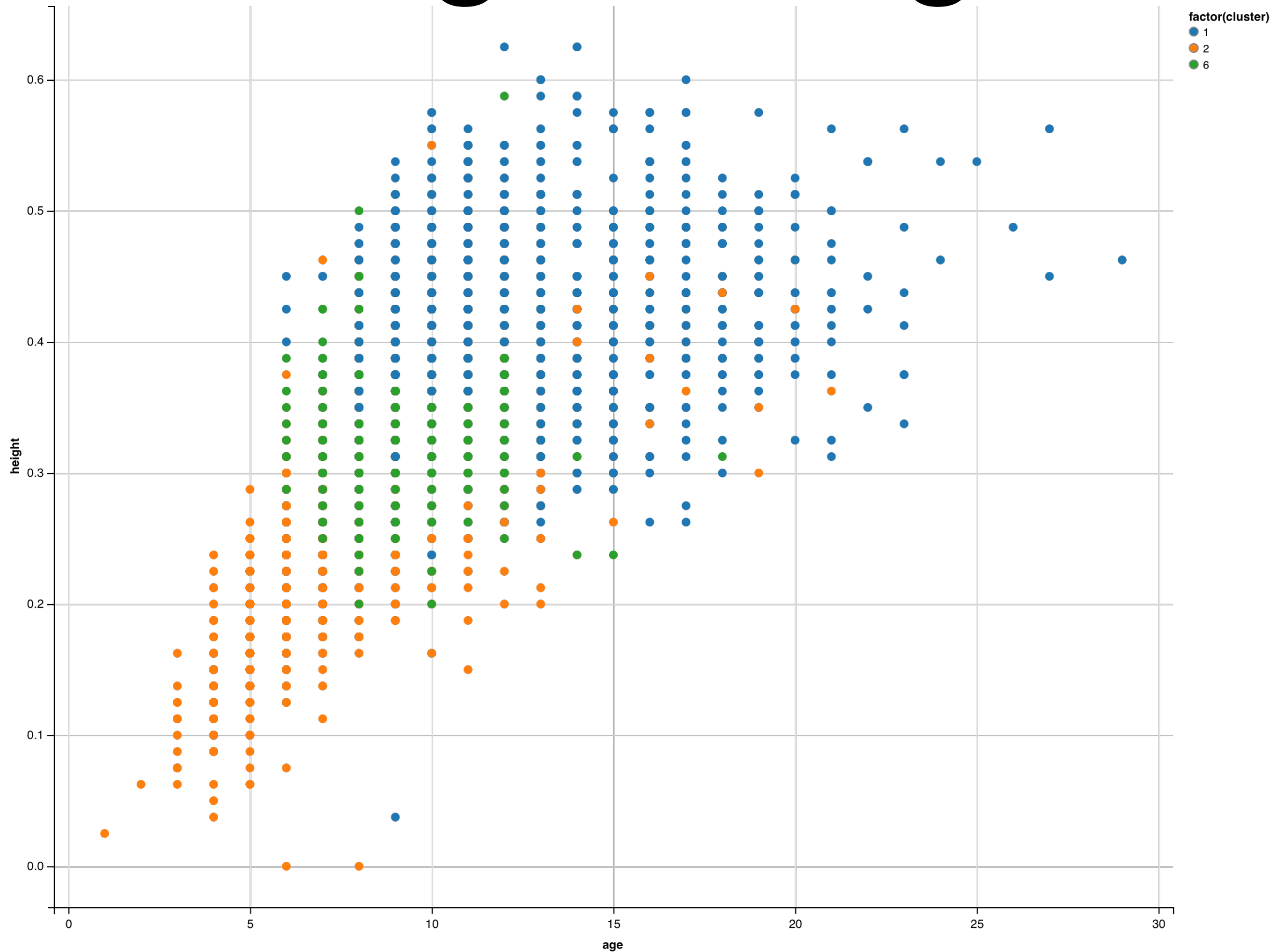
# Objective

- Simply predict the age of the crab in boston based on their size, weight and diameters etc

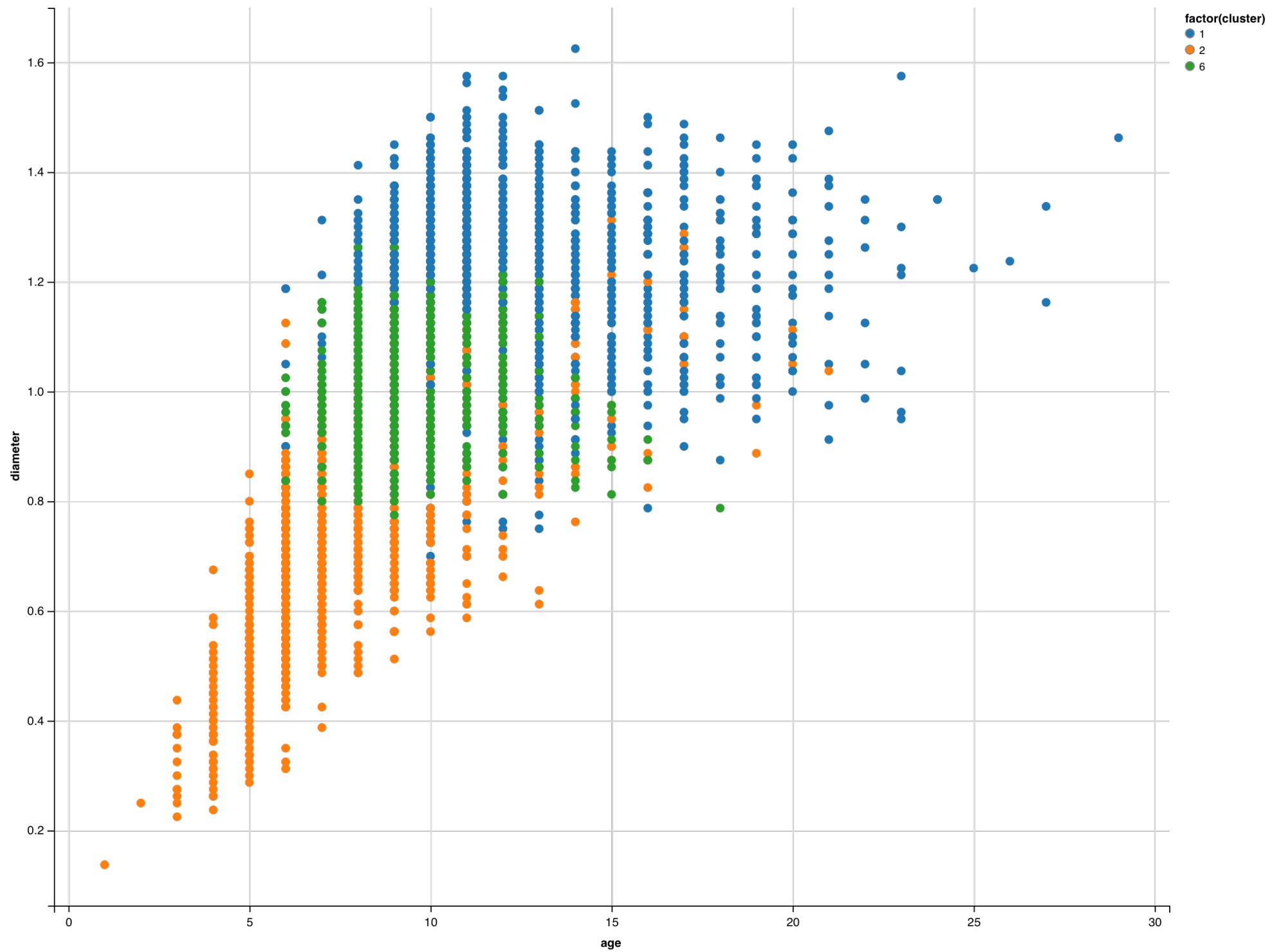
# Challenges

- Read files from pdf with dynamic structure and dirty data
- Lack of strong correlation from the features to crab's age
- Without additional information, assume there are more than one species in the data set, so further unsupervised k-mean is performed
- Data set split 70% (Train) / 30% (Final)  
within the 70% Train > 60% (Test) / 40% (Valid)

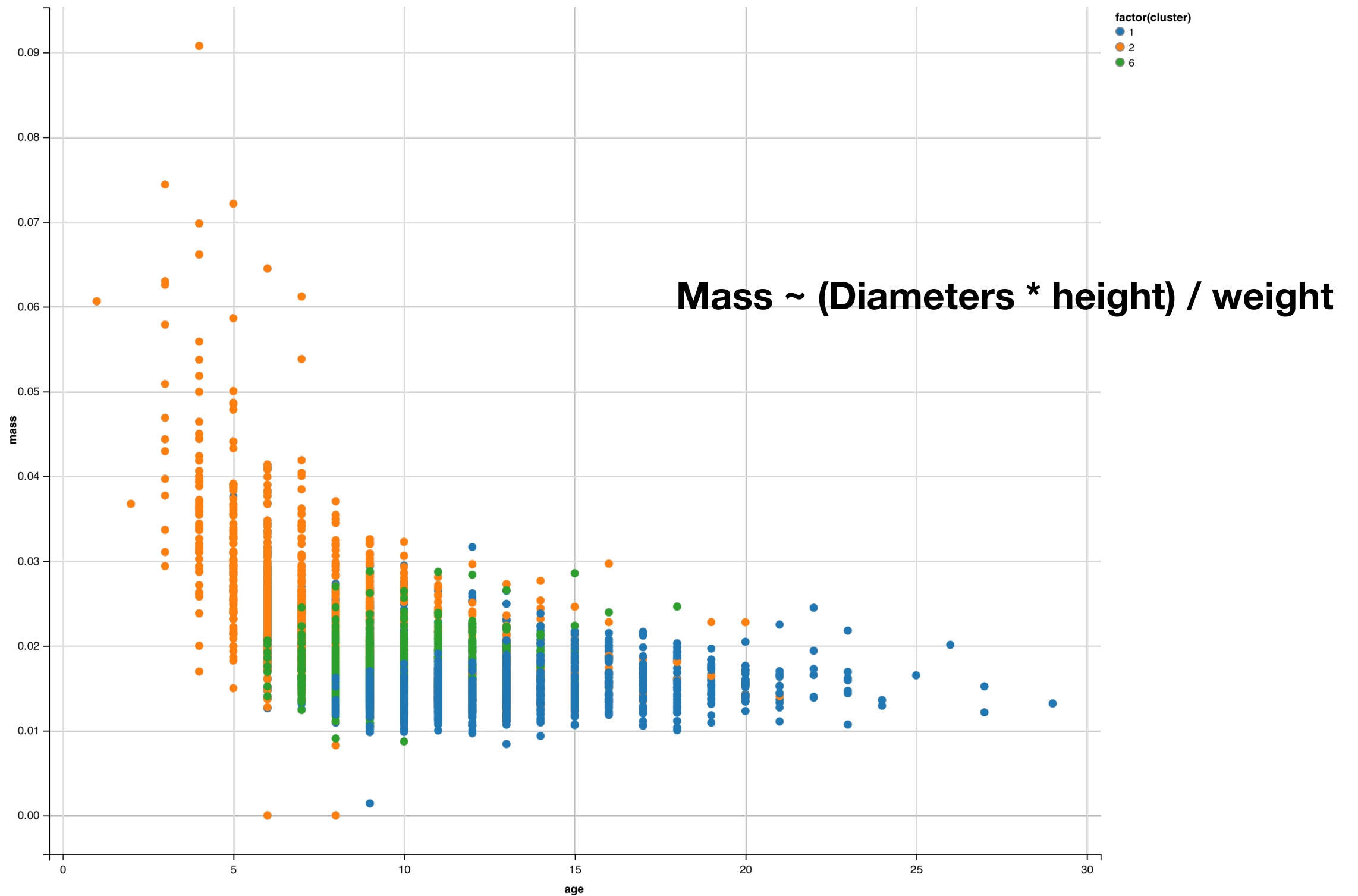
# Height vs Age

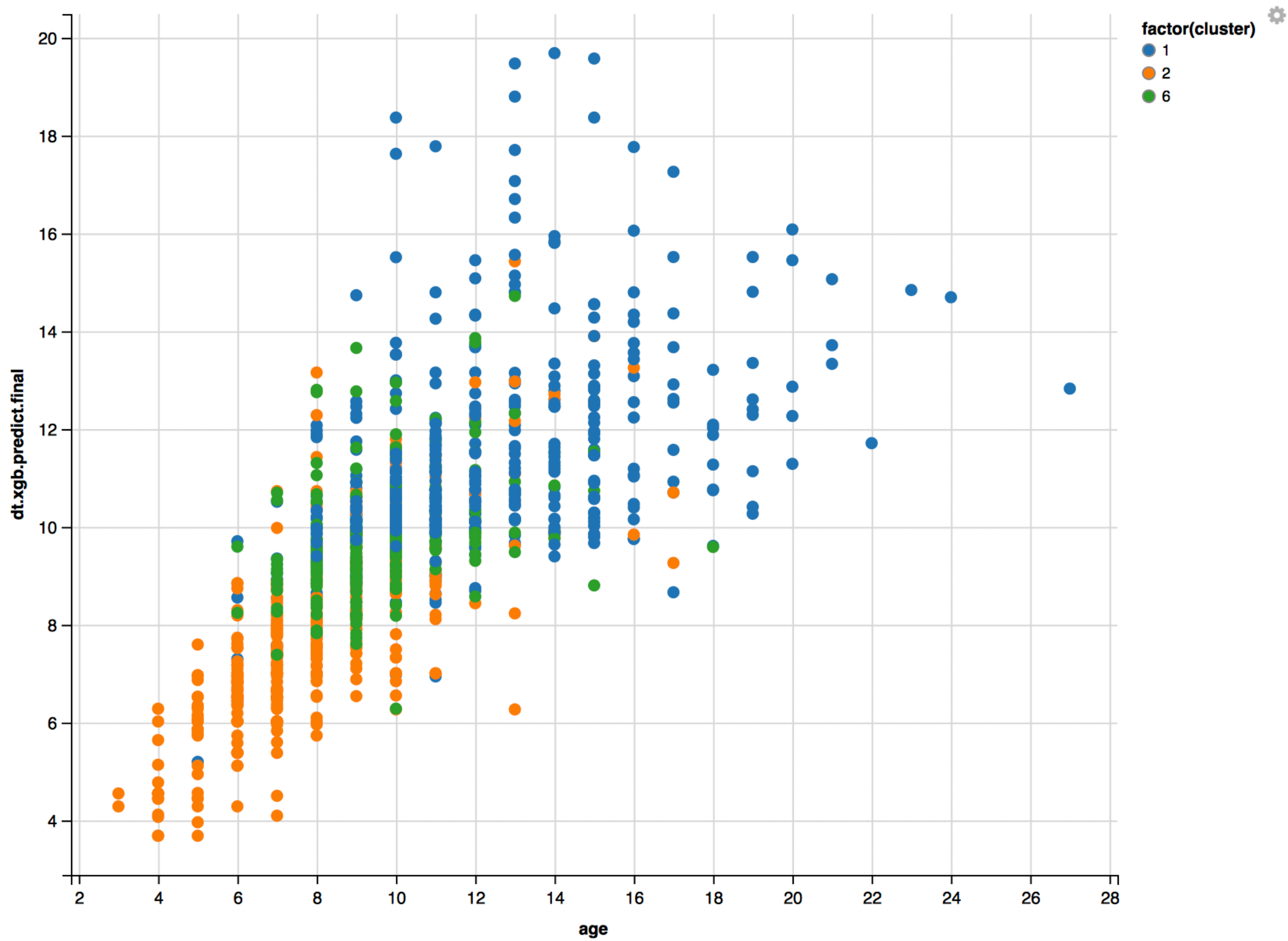


# Diameters vs Age



# Mass vs Age





# ML Modeling

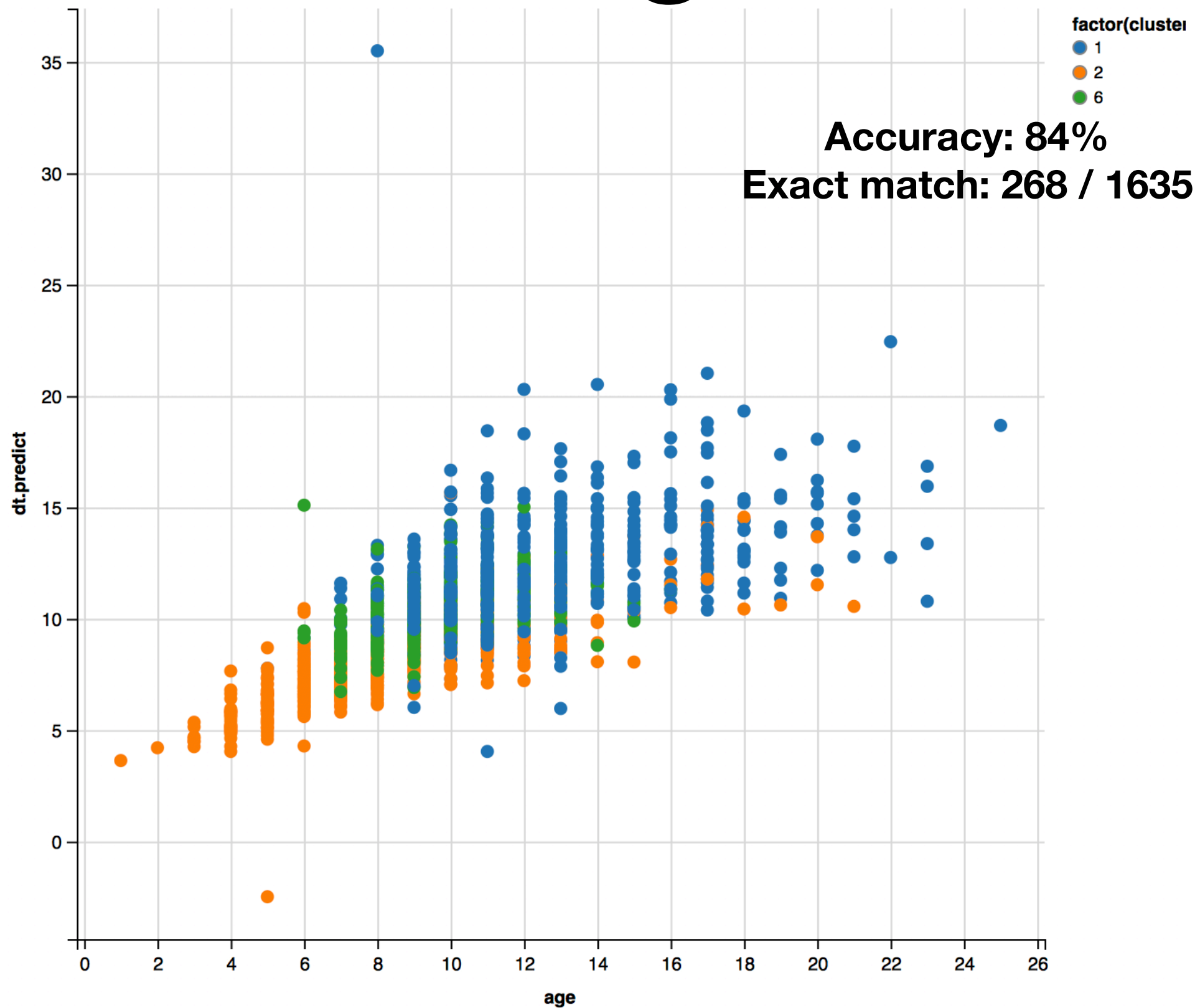
- 2 types of ML models are used in this analysis
  - Linear regression models
  - Xgboost



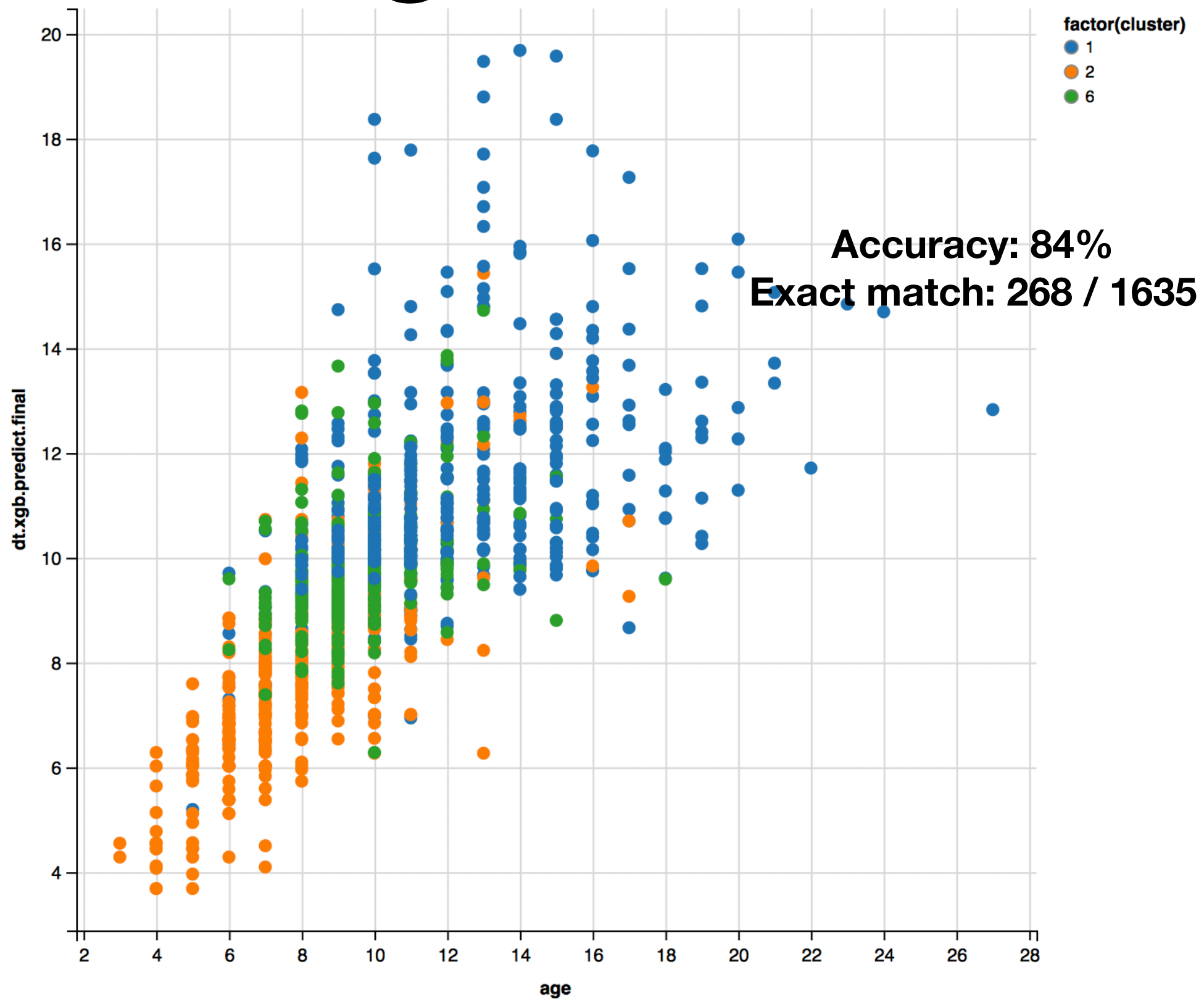
# ML Modeling

- 2 types of ML models are used in this analysis
  - Linear regression models
  - Xgboost
- Accuracy defined as:  
$$1 - ( \text{abs}(\text{Actual Age} - \text{Predicted Age}) / \text{Actual Age} )$$
- Exact Match  
$$\text{Sum}( \text{Round}(\text{Predicted Age} ) = \text{Age} ) / \text{Total \# of Results}$$

# Linear Regression



# Xgboost



# Conclusion

- In net regression model was able to achieve great accuracy in both valid and final data sets and scored higher in the accuracy and # of exact matched evaluation criteria.
- As such the artifacts of linear regression has been selected as the core model for the program for prediction in rScript which can be used for MI reporting, Data ETL.