

Report 3

Advay Vyas

September 18, 2025

Contents

Introduction	2
Data processing	2
Load and preprocess data (borrowed from other file)	2
Support Vector Regression (SVR)	4
Theoretical background	4
Implementation	4
Bayesian Additive Regression Trees	4
Theoretical background	4
Implementation	4
Elastic Net / LASSO / Ridge	6
Theoretical background	6
Implementation	6
Conclusion	6

```
library(dplyr)
library(tidyverse)
library(ggplot2)
library(lubridate)
library(patchwork)
library(corrplot)
library(mosaic)
library(moderndiver)
library(effectsize)
library(tidyr)
library(caret)
library(purrr)
library(fastDummies)
library(GGally)
```

Introduction

This week, I plan to use the methods I mentioned last week to estimate the WIS differences. I'll process and prepare the data and then apply several models to try and see which work best.

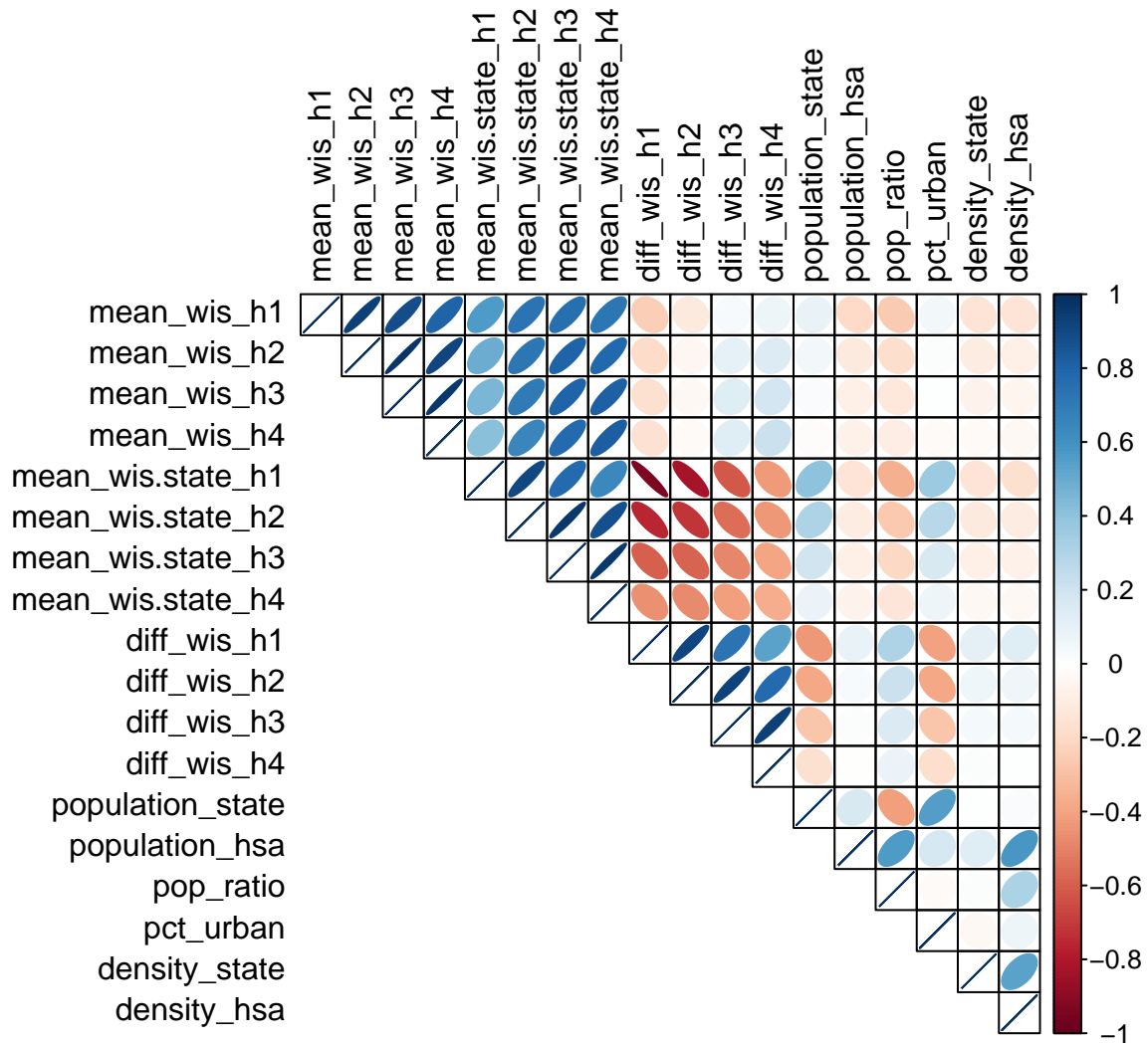
Data processing

Load and preprocess data (borrowed from other file)

```
df <- readr::read_csv("GBQR_diff_wis.csv", show_col_types = FALSE)

df_season <- df %>%
  filter(season == "2023/24")

cor_matrix = cor(df[, 4:ncol(df_season)])
# corrplot.mixed(cor_matrix, tl.col = "black", tl.pos = "lt", addgrid.col = TRUE, upper="color", lower=
corrplot(cor_matrix, tl.col = "black", tl.pos = "lt", addgrid.col = TRUE, type = "upper", method = "ell
```



Support Vector Regression (SVR)

Theoretical background

Implementation

Bayesian Additive Regression Trees

Theoretical background

Implementation

```
library(dbarts)
```

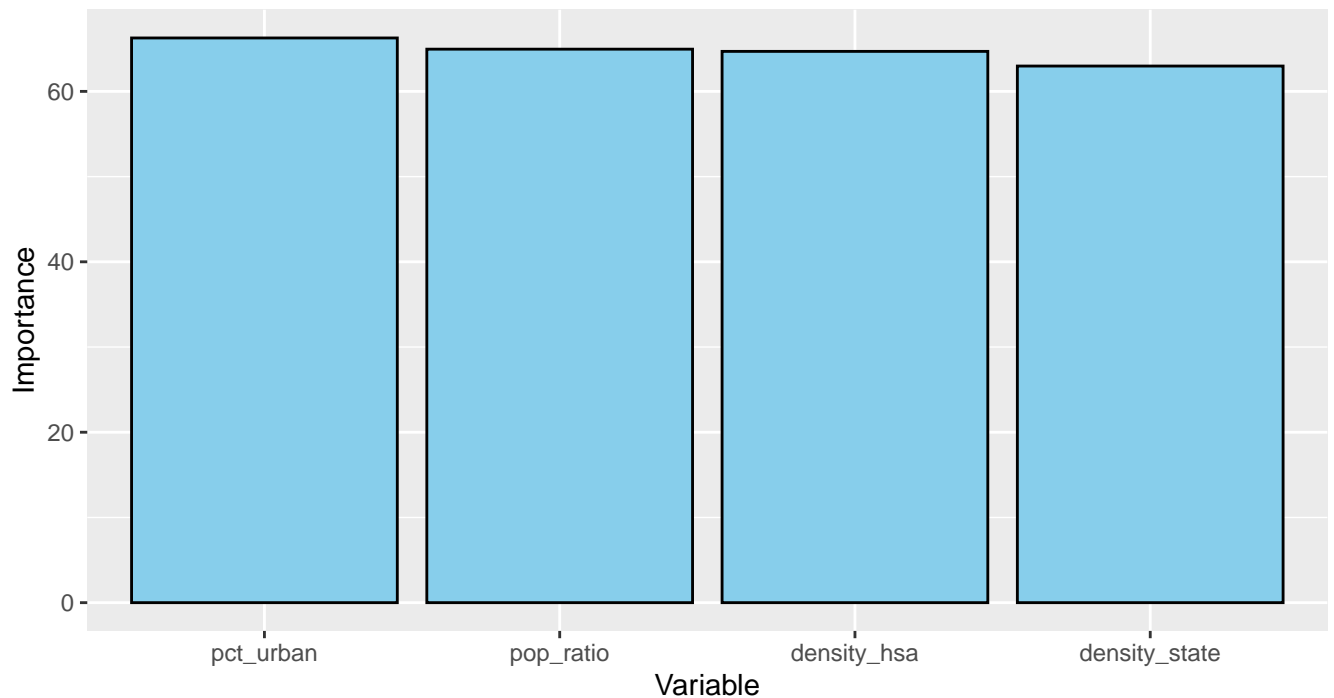
```
##  
## Attaching package: 'dbarts'  
  
## The following object is masked from 'package:tidyr':  
##  
##      extract
```

```
bart_fit <- bart(  
  x.train = df_season[, c("pop_ratio", "pct_urban",  
                          "density_state", "density_hsa")],  
  y.train = df_season$diff_wis_h3,  
  verbose = FALSE, keptrees = TRUE  
)
```

Feature Importance

```
bart_vi = as.data.frame(bart_fit$varcount) %>%  
  summarise(across(everything(), ~ mean(.x, na.rm = TRUE))) %>%  
  pivot_longer(cols = everything(),  
               names_to = "variable",  
               values_to = "avg")  
  
ggplot(bart_vi) + geom_col(aes(x = reorder(variable, -avg), y = avg), fill='skyblue', col='black') + labs
```

BART Variables by Importance

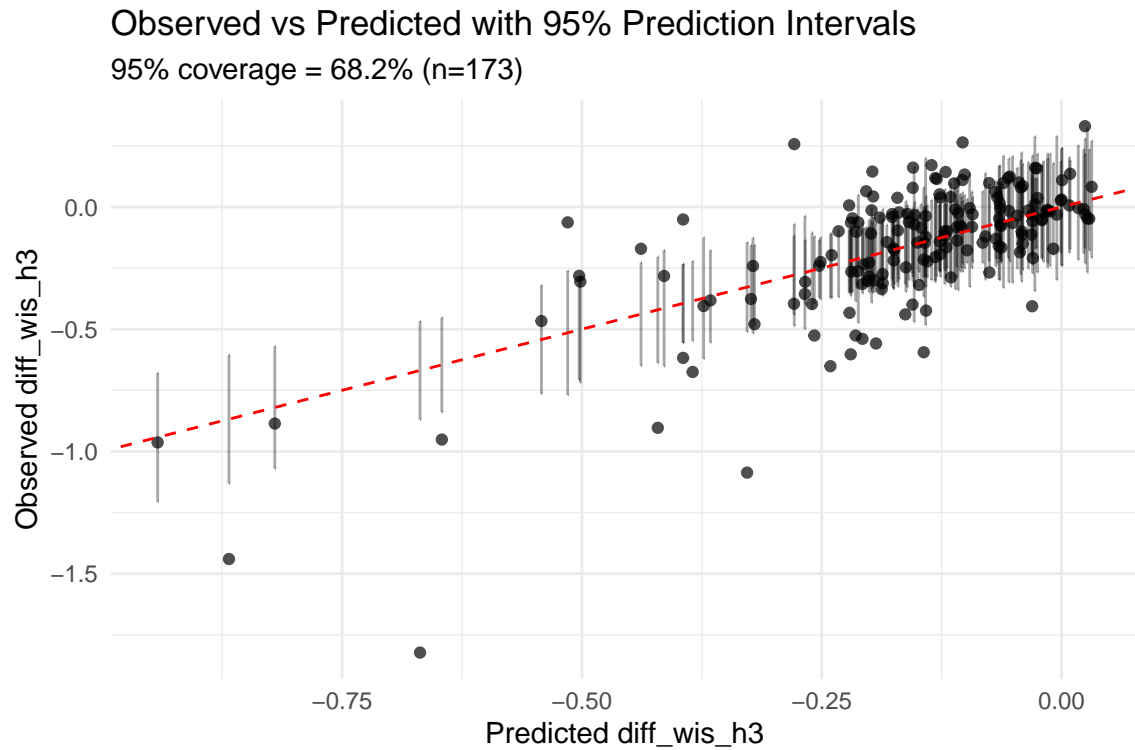


```
pred = predict(bart_fit, newdata = df_season)
bart_predictions = as.data.frame(t(apply(pred, 2, quantile, probs = c(0.025, 0.975))))
colnames(bart_predictions) = c("ci_low", "ci_high")

bart_predictions$estimate = (bart_predictions$ci_low + bart_predictions$ci_high) / 2
bart_predictions$observed = df_season$diff_wis_h3
```

```
coverage_95 = mean(
  bart_predictions$observed >= bart_predictions$ci_low &
  bart_predictions$observed <= bart_predictions$ci_high,
  na.rm = TRUE
)
```

```
ggplot(bart_predictions, aes(x = estimate, y = observed)) +
  geom_errorbar(aes(ymin = ci_low, ymax = ci_high), width = 0, alpha = 0.35) +
  geom_point(alpha = 0.7) +
  geom_abline(slope = 1, intercept = 0, linetype = "dashed", color = "red") +
  labs(x = "Predicted diff_wis_h3", y = "Observed diff_wis_h3",
       title = "Observed vs Predicted with 95% Prediction Intervals",
       subtitle = sprintf("95%% coverage = %.1f%% (n=%d)", 100*coverage_95, nrow(bart_predictions))) +
  theme_minimal()
```



Elastic Net / LASSO / Ridge

Theoretical background

Implementation

Conclusion