# Predicting Food Store Inspection Grades

By André Ruckdaeschel, Kaspar Lichtsteiner & Matthias Steiner

Rating
- A
- B
- C

# Predicting Food Store Inspection Grades

By André Ruckdaeschel,
Kaspar Lichtsteiner &
Matthias Steiner

Card Number _____

Establishment Name _____

Date Issued _____

**NYC**
Health

For additional information
or a copy of an inspection
report, call 311 or visit
nyc.gov/health

Card Number _____

Establishment Name _____

Date Issued _____

**NYC**
Health

For additional information
or a copy of an inspection
report, call 311 or visit
nyc.gov/health

Card Number _____

Establishment Name _____

Date Issued _____

**NYC**
Health

For additional information
or a copy of an inspection
report, call 311 or visit
nyc.gov/health

# The Inspection Data

# The Inspection Data

| Inspection Grade | Store Name | County | Location | Date | ..... |
|---|---|---|---|---|---|
| A | ZUMY 833 INC | Queens | 1007 BRIGHTON BEACH AVE BROOKLYN, NY 11235 (40.578152, -73.959054) | 01/16/2019 | |
| A | MACKALLIE LLC | Suffolk | 209 UNION AVE NEW ROCHELLE, NY 10801 (40.909533, -73.793911) | 04/27/2018 | |
| C | BALS BAGELS | Kings | NA | 12/31/2018 | |
| A | QUICKWAY 68 | Bronx | 324 JACKSON AVE SYOSSET, NY 11791 (40.810935, -73.501746) | 05/08/2018 | |
| B | TARGET 2211 | Kings | NA | 01/16/2019 | |
| ⋮ | ⋮ | | ⋮ | ⋮ | |

# The Inspection Data

| Inspection Grade | Store Name | County | Address | Latitude | Longitude | Date | ..... |
|---|---|---|---|---|---|---|---|
| A | ZUMY 833 INC | Queens | 1007 BRIGHTON BEACH AVE BROOKLYN, NY 11235 | 40.578152 | -73.959054 | 01/16/2019 | |
| A | MACKALLIE LLC | Suffolk | 209 UNION AVE NEW ROCHELLE, NY 10801 | 40.909533 | -73.793911 | 04/27/2018 | |
| C | BALS BAGELS | Kings | 222 HOYT ST BROOKLYN, NY 11217 | NA | NA | 12/31/2018 | |
| A | QUICKWAY 68 | Bronx | 324 JACKSON AVE SYOSSET, NY 11791 | 40.810935 | -73.501746 | 05/08/2018 | |
| B | TARGET 2211 | Kings | 204 LIBERTY ST PENN YAN, NY 14527 | NA | NA | 01/16/2019 | |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | |

# The Inspection Data

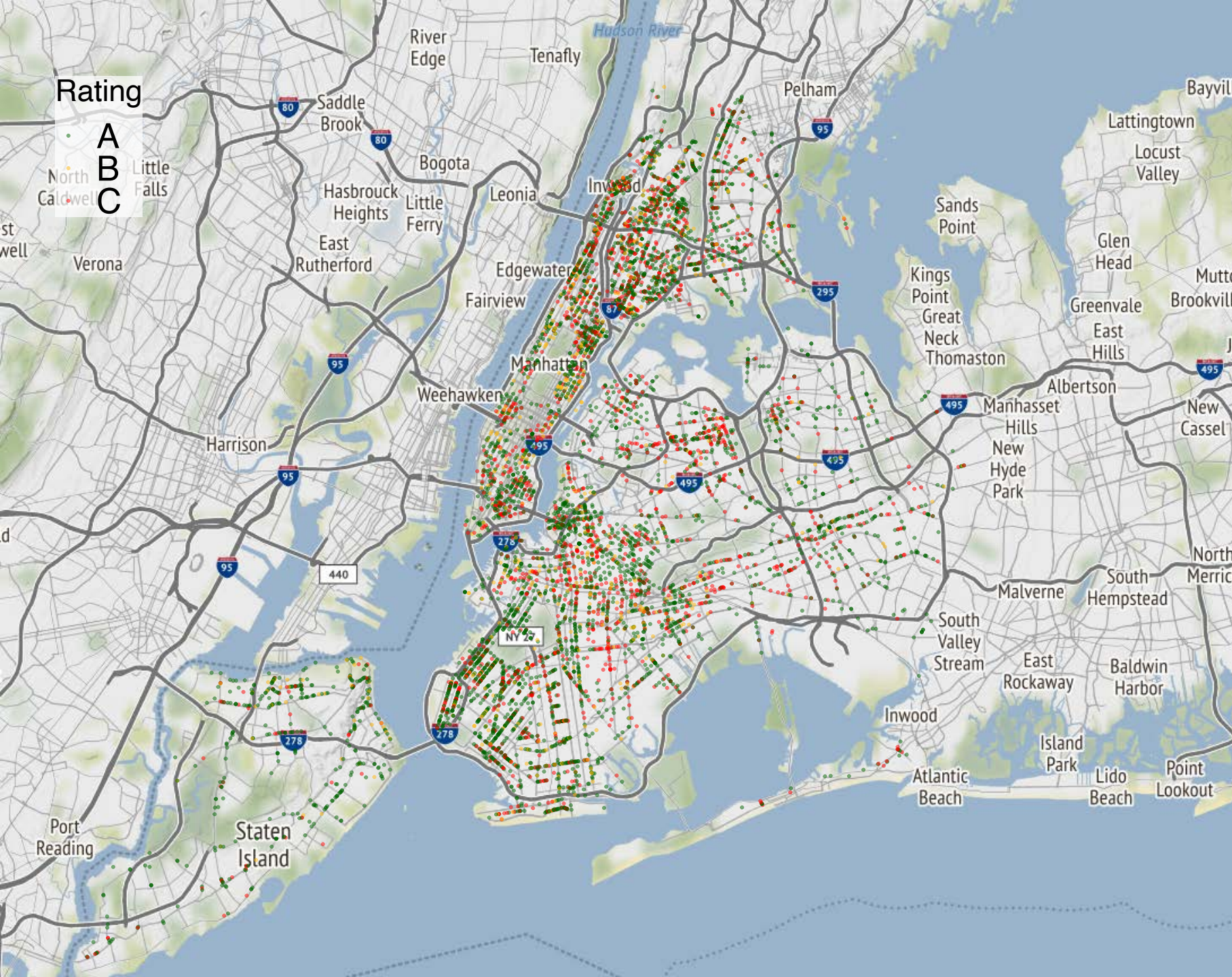| Inspection Grade | Store Name | County | Address | Latitude | Longitude | Date | ⋯⋯ |
|---|---|---|---|---|---|---|---|
| A | ZUMY 833 INC | Queens | 1007 BRIGHTON BEACH AVE BROOKLYN, NY 11235 | 40.578152 | -73.959054 | 01/16/2019 | |
| A | MACKALLIE LLC | Suffolk | 209 UNION AVE NEW ROCHELLE, NY 10801 | 40.909533 | -73.793911 | 04/27/2018 | |
| C | BALS BAGELS | Kings | 222 HOYT ST BROOKLYN, NY 11217 | NA | NA | 12/31/2018 | |
| A | QUICKWAY 68 | Bronx | 324 JACKSON AVE SYOSSET, NY 11791 | 40.810935 | -73.501746 | 05/08/2018 | |
| B | TARGET 2211 | Kings | 204 LIBERTY ST PENN YAN, NY 14527 | NA | NA | 01/16/2019 | |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | |

# Google Maps API for Missing Coordinates

- Free Service offered by Google

- Requires a one-time registration with a valid Email address

- Then the personal API key must be included in the R script

```r
118
119   # use Google maps to get missing coordiantes (takes few minutes an requires API in the head)
120   inspect_data_na <- inspect_data %>%
121     filter(is.na(Latitude)) %>% # all missing coordinates
122     mutate_geocode(Address) %>% # applies Google Maps API
123     mutate(Latitude = lat, Longitude = lon) %>%
124     dplyr::select(-c(lat, lon)) %>%
125     filter(!is.na(Latitude)) # 248 still missing and dropped
126
127   # add new coordinates
128   inspect_data <- inspect_data %>%
129     filter(!is.na(Latitude)) %>%
130     bind_rows(inspect_data_na)
131
132   table(is.na(inspect_data$Longitude)) # no more coordinates with NA
133
134   rm(inspect_data_na)
135
136   save(inspect_data, file = "./data/inspect_data.RData")
137
138   ##################################################################################
```

125:59   🔳 Add Coordinates of shops ⇕                                                          R Script ⇕

**Console**   **Terminal** ✕

~/DSF/ 🔊

```r
> rm(coord)
>
> # 748 coordinates are missing
> table(is.na(inspect_data$Longitude))

FALSE   TRUE
16508   748
>
> # create address column
> inspect_data <- inspect_data %>%
+   mutate(Address = str_c(Street, Zip.Code, sep = ", ")) %>%
+   mutate(Address = str_c(Address, City, sep = " ")) %>%
+   mutate(Address = str_c(Address, State.Code, sep = ", "))
> register_google(key = "AIzaSyCnb_afuEHvqD4CR-xBY_u9Z4El21KpQus")
>
```
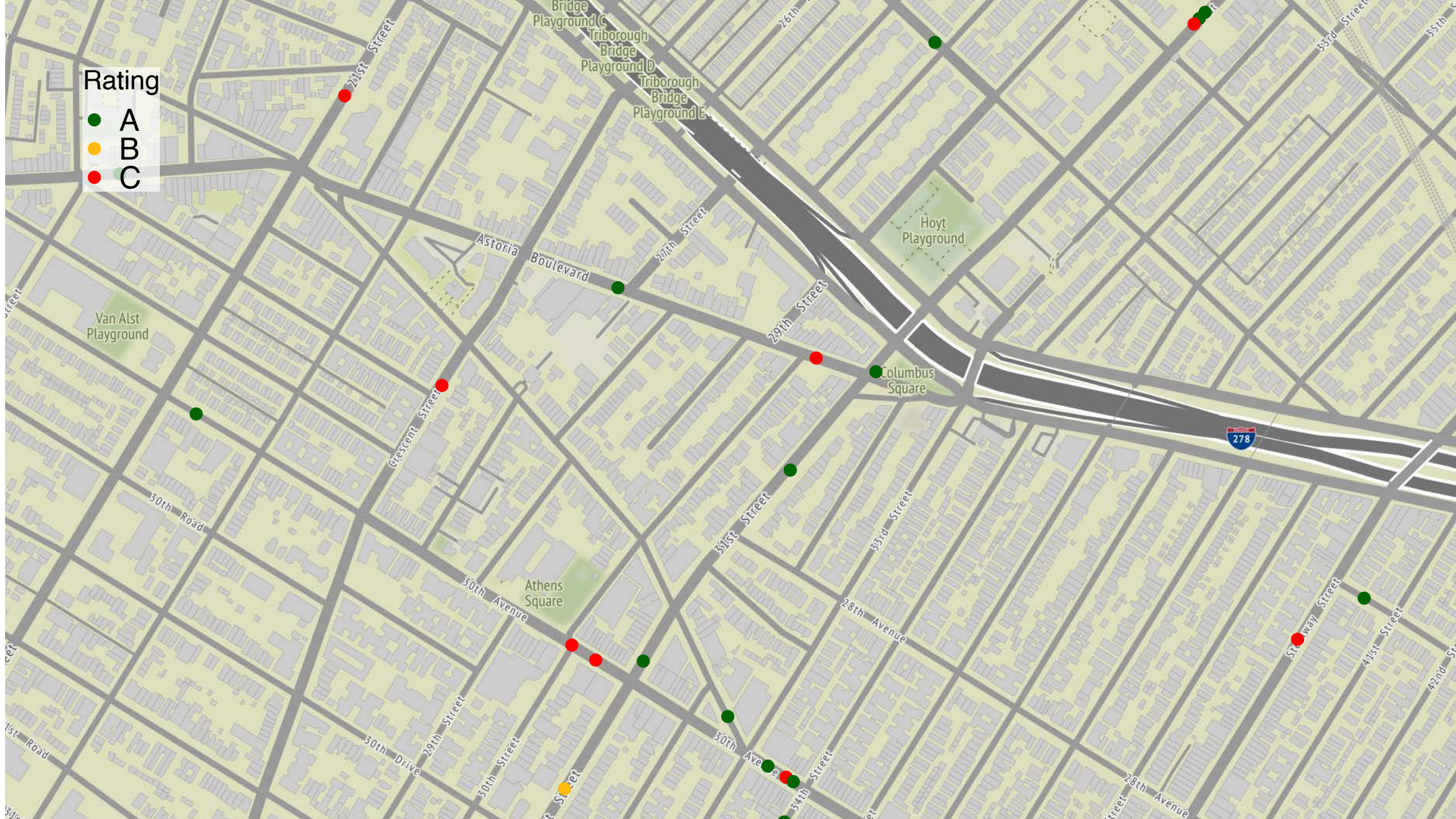
# Focus on New York City

- More covariate data available

- Classes are more equally distributed

# Haversine Formula

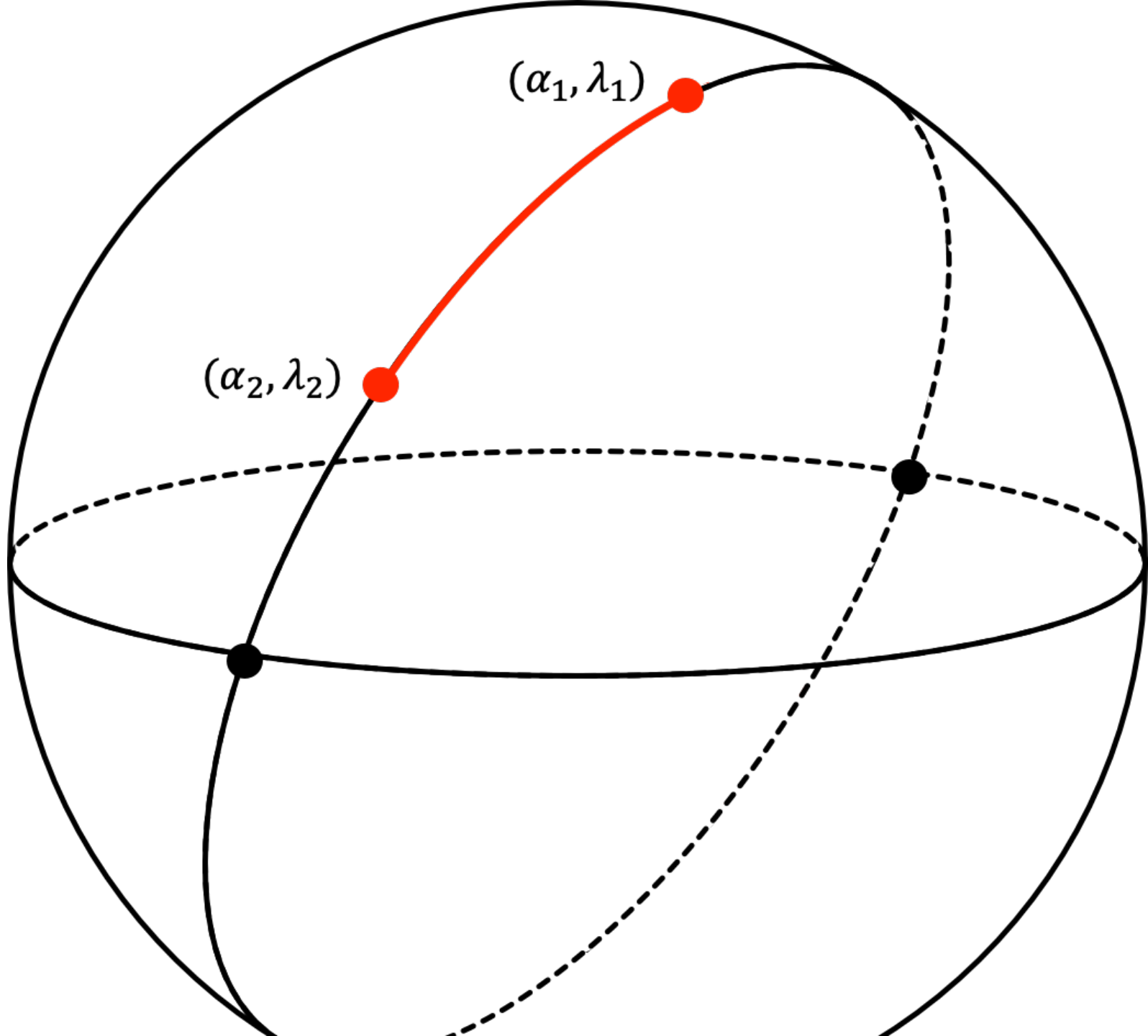$$a = sin^2\left(\frac{\Delta\alpha}{2}\right) + \cos(\alpha_1)\cos(\alpha_2)\,sin^2\left(\frac{\Delta\lambda}{2}\right)$$

$$c = R\left[2\,\text{atan}^2(\sqrt{a}, \sqrt{1-a})\right]$$

With:

R = (Mean) radius of the earth (6,371km)
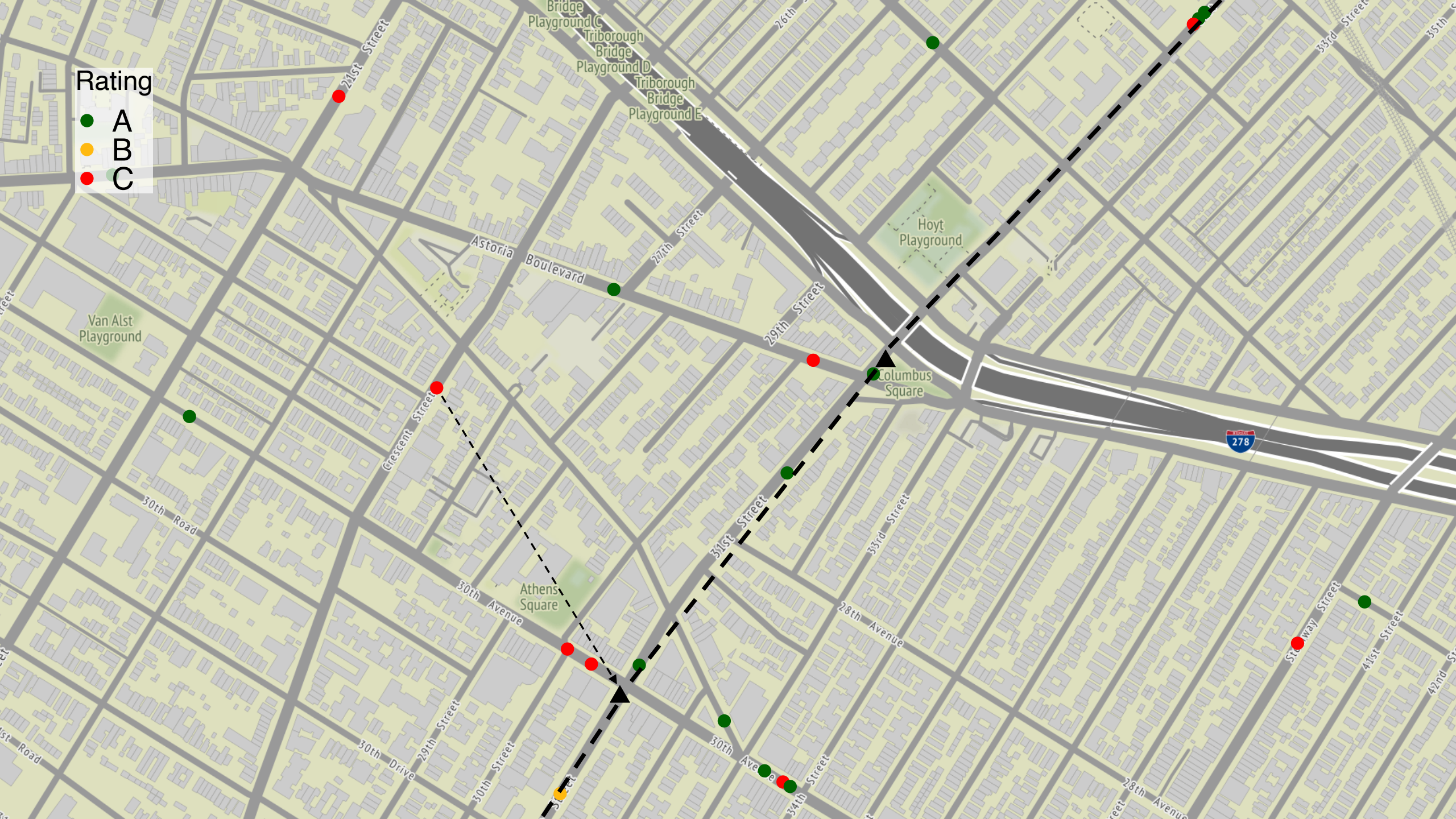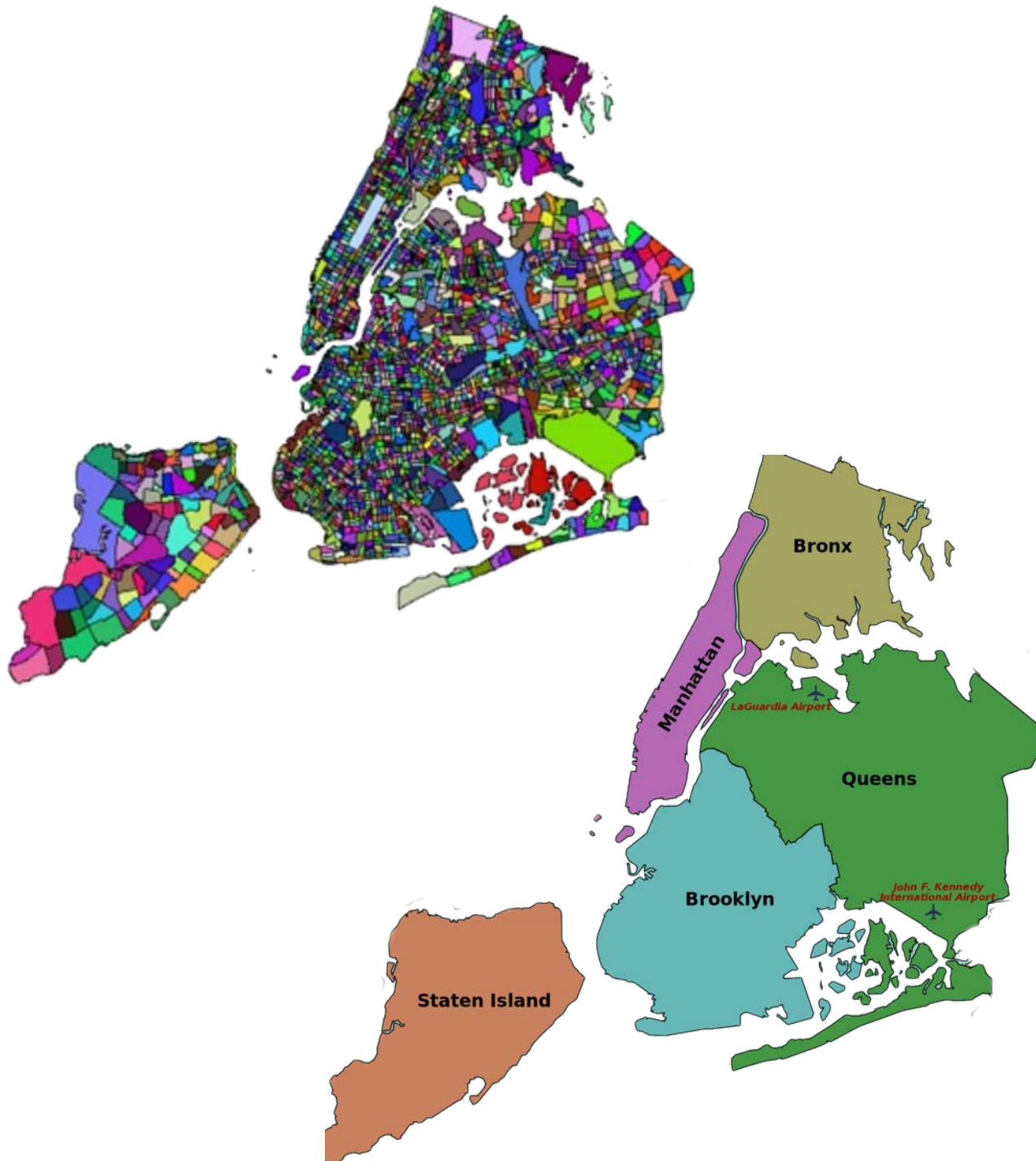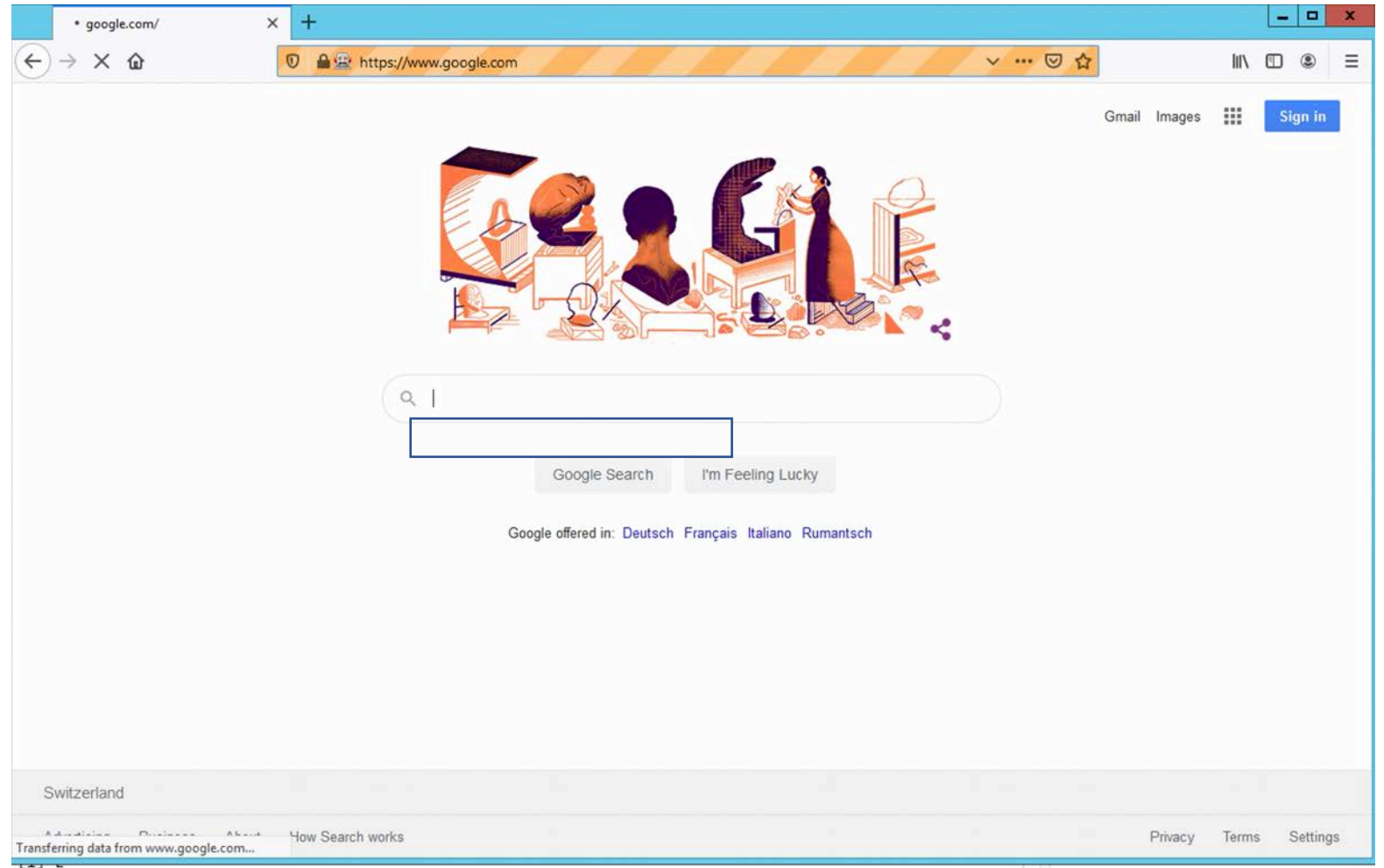
$\alpha_{1,2} = Latitude$  $\lambda_{1,2} = Longitude$

1km

# Demographic Data

- U.S. Census Bureau: Counties
  - merging via counties

- U.S. Census Bureau: Census Tracts
  - merging via census tract
  - Translation necessary: AddTrac
  - Geocoding service: Addresses
    - State FIPS
    - County FIPS
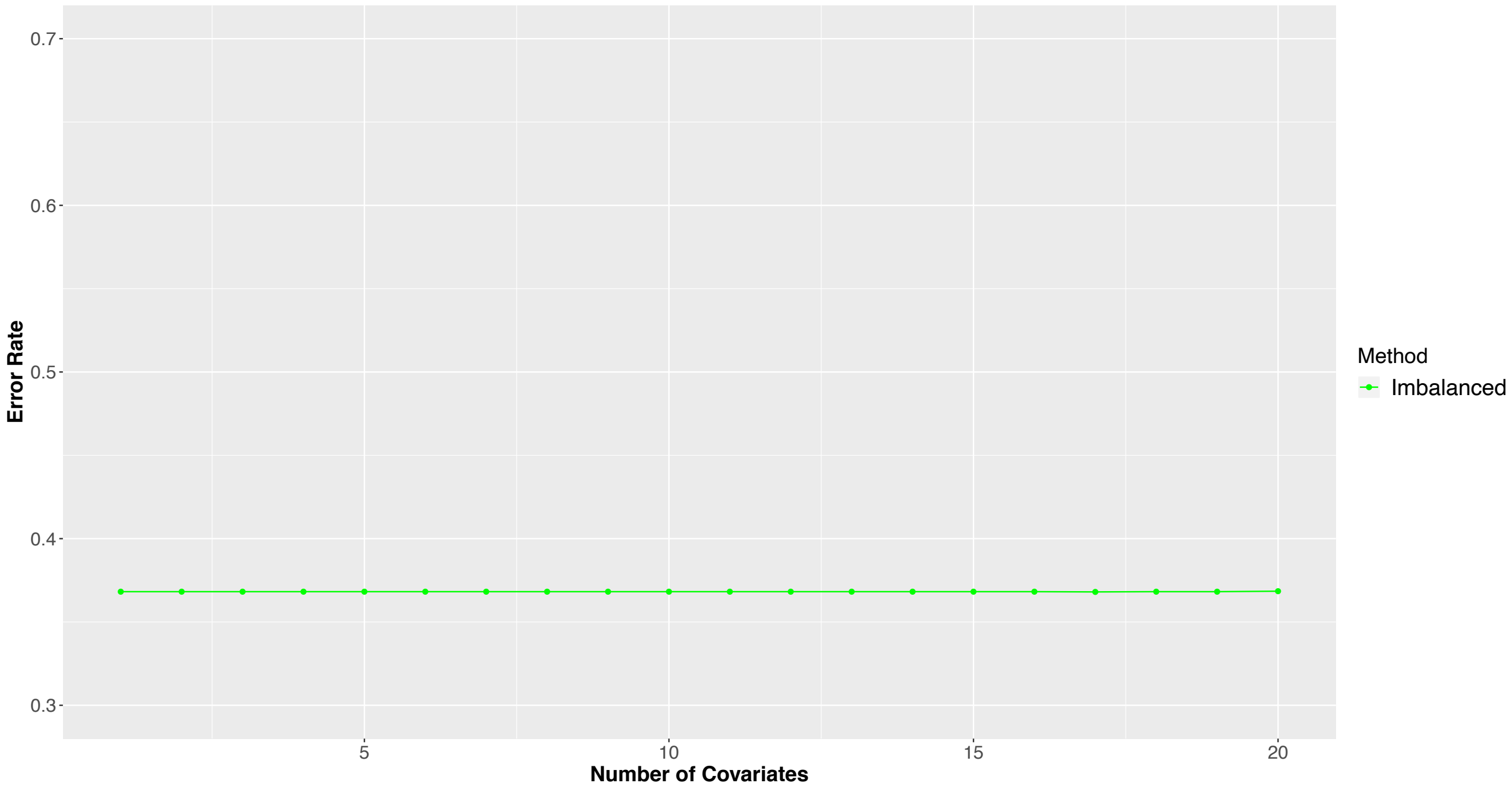    - Census Tract Id

# Google web scraper

- Automated Googles search

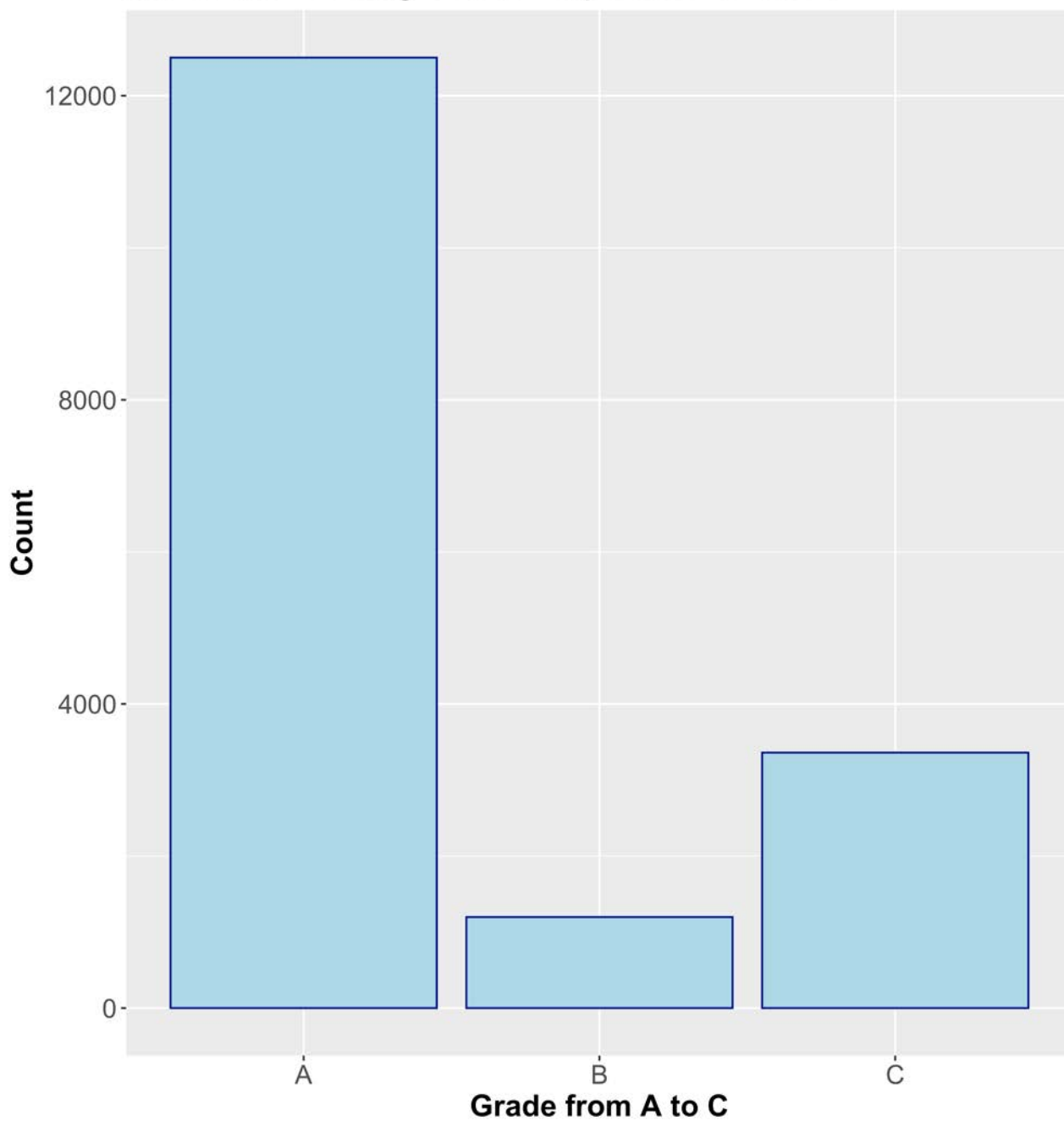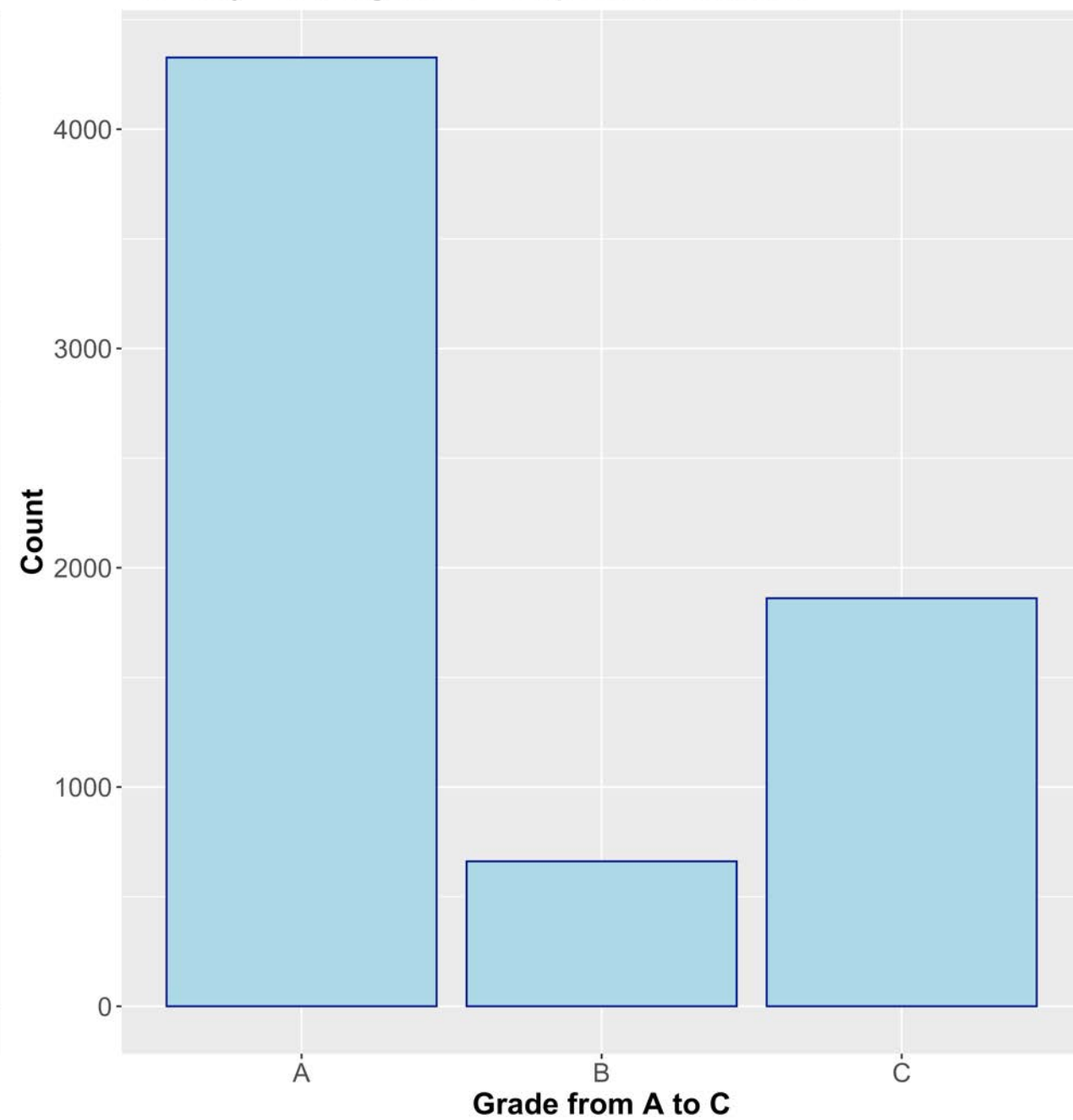- Gather Google star ratings and number of reviews

Methodology illustrated with LDA

Methodology illsutrated with LDA
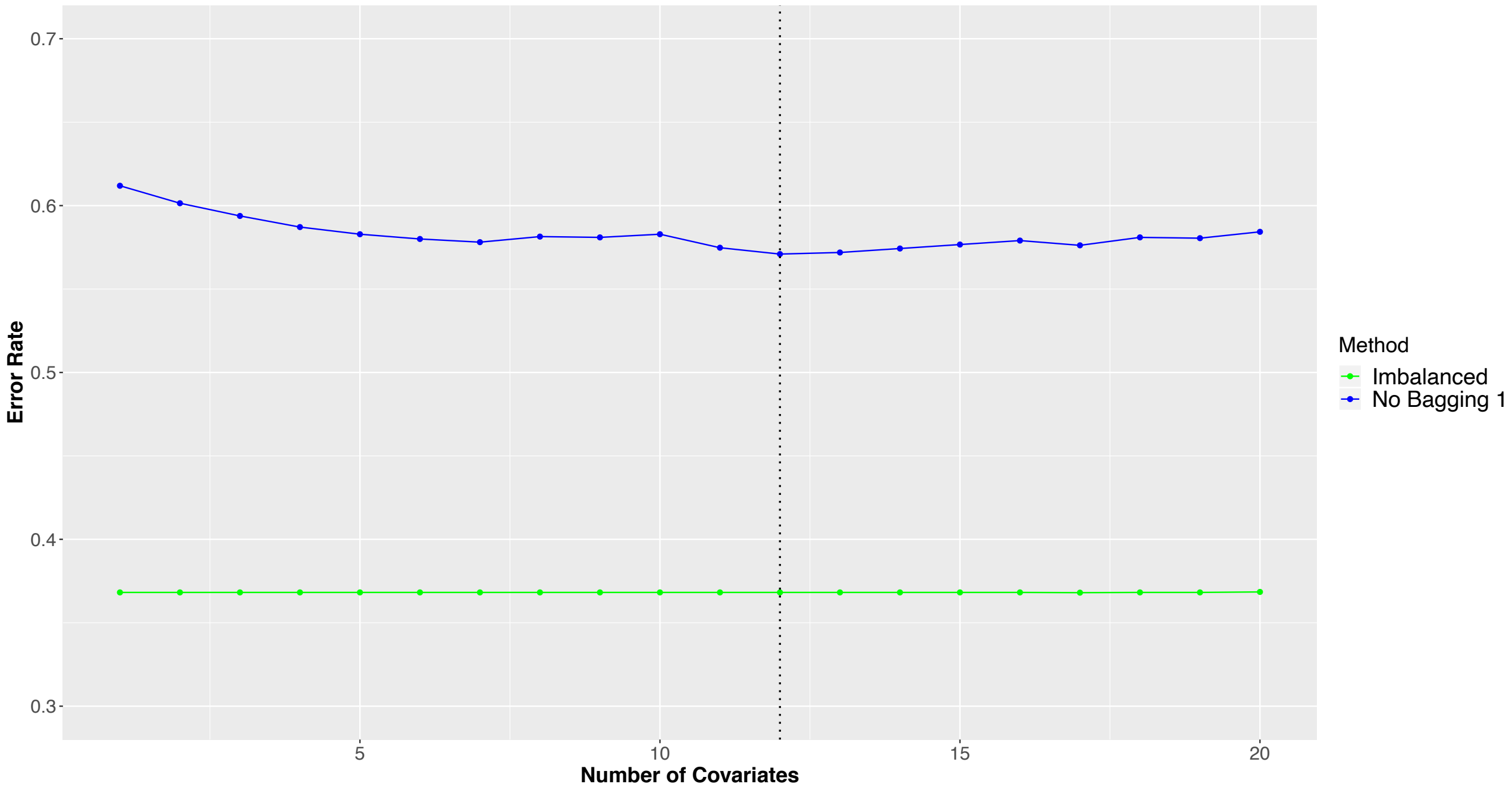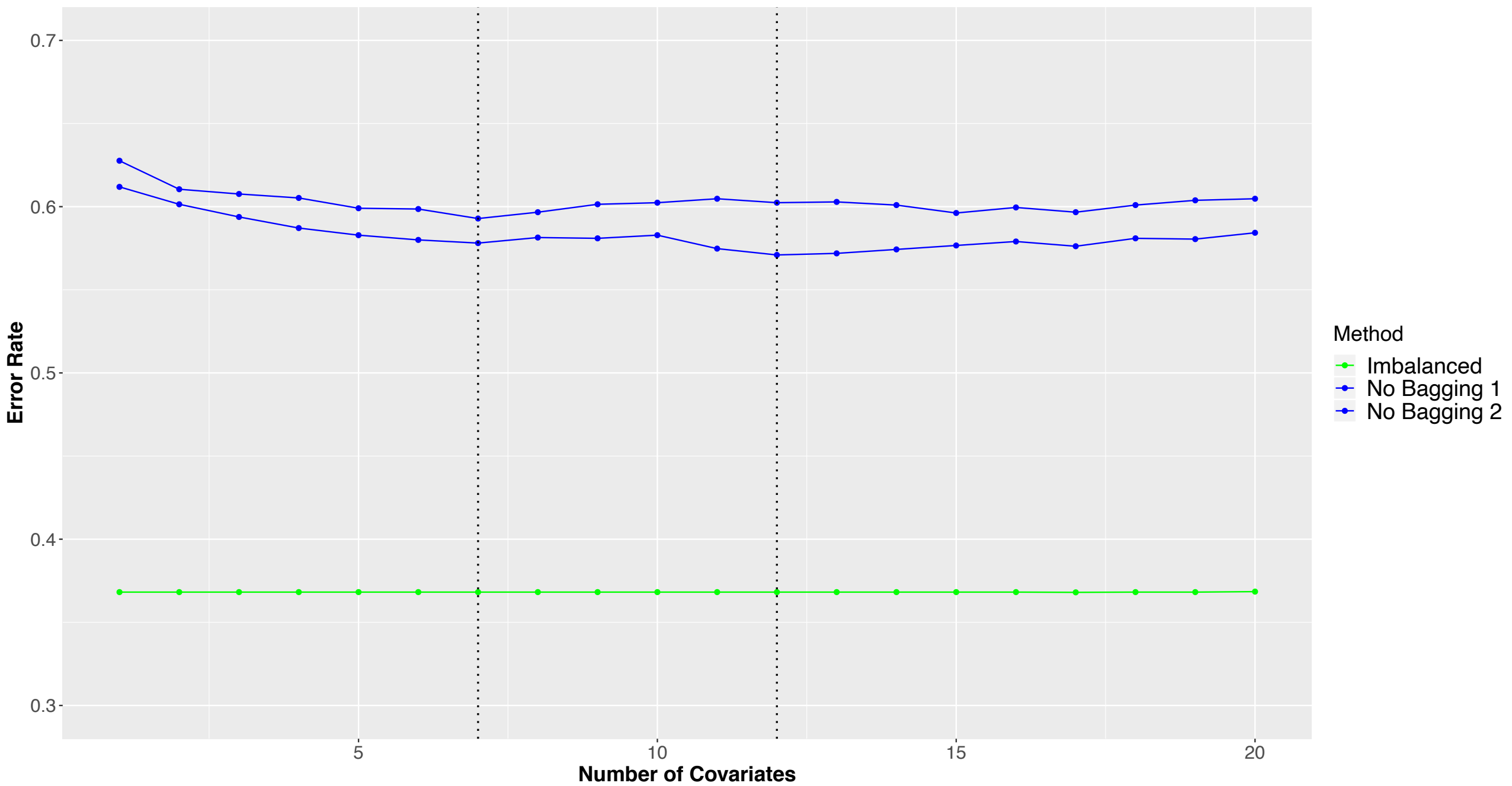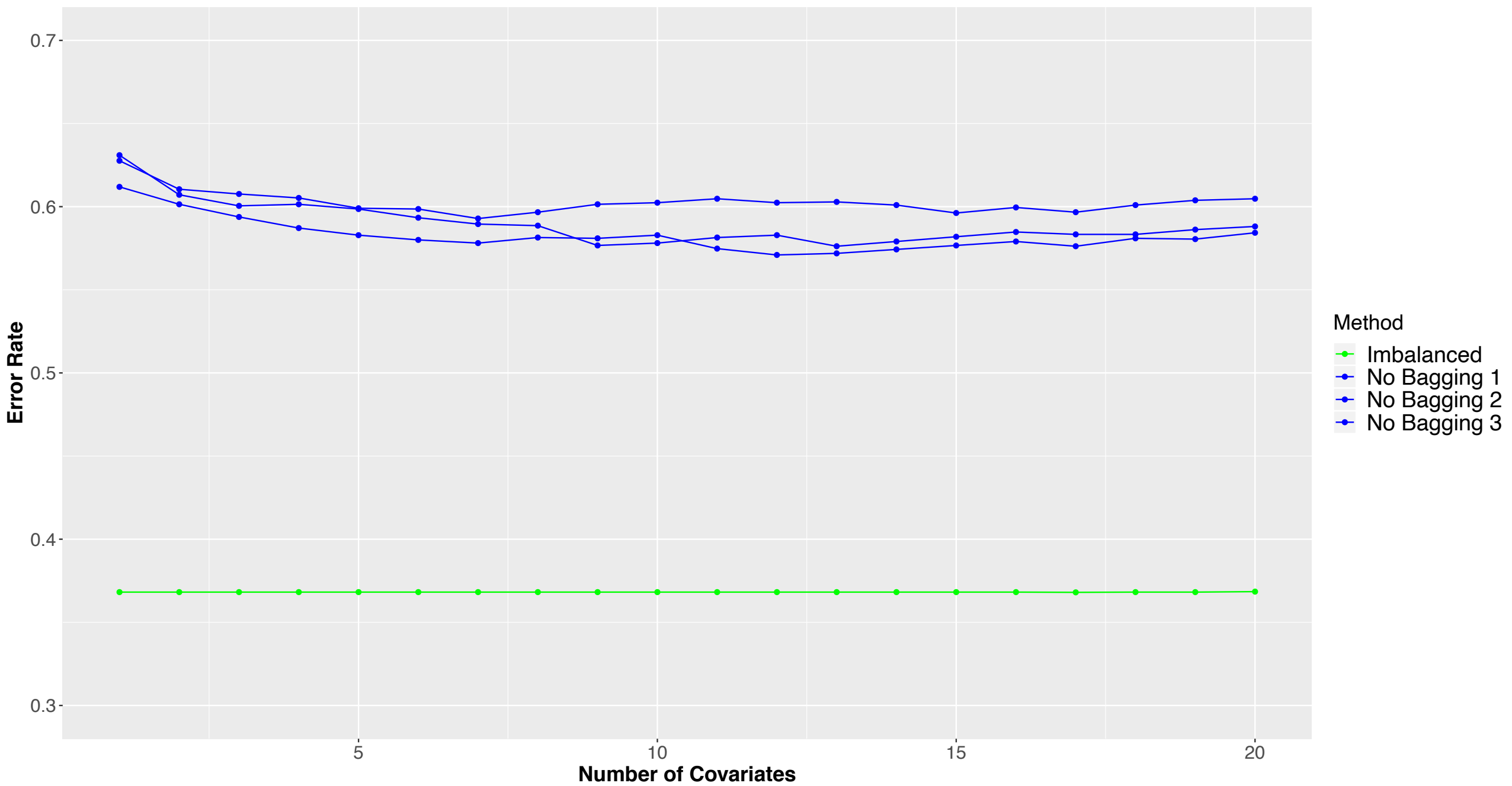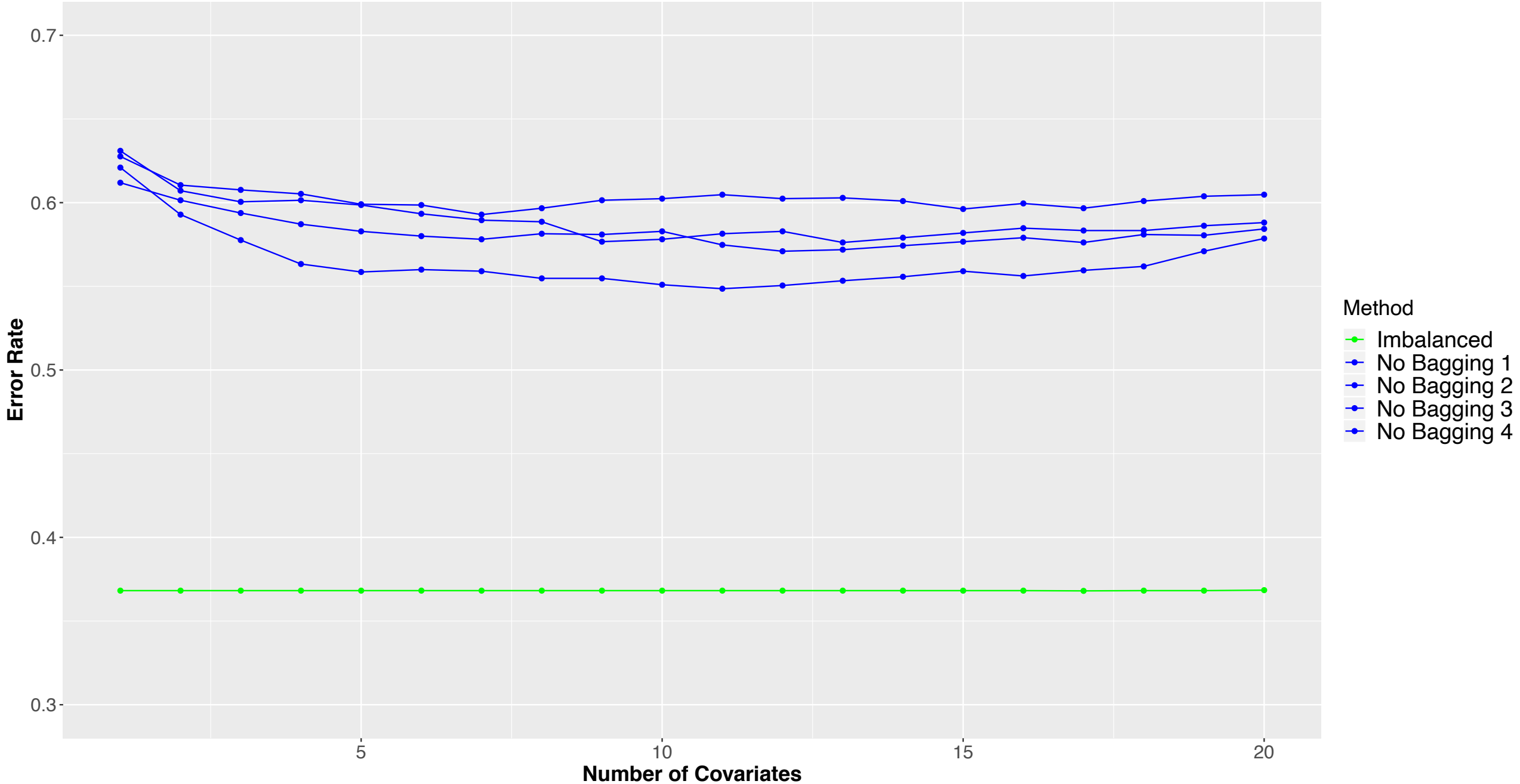
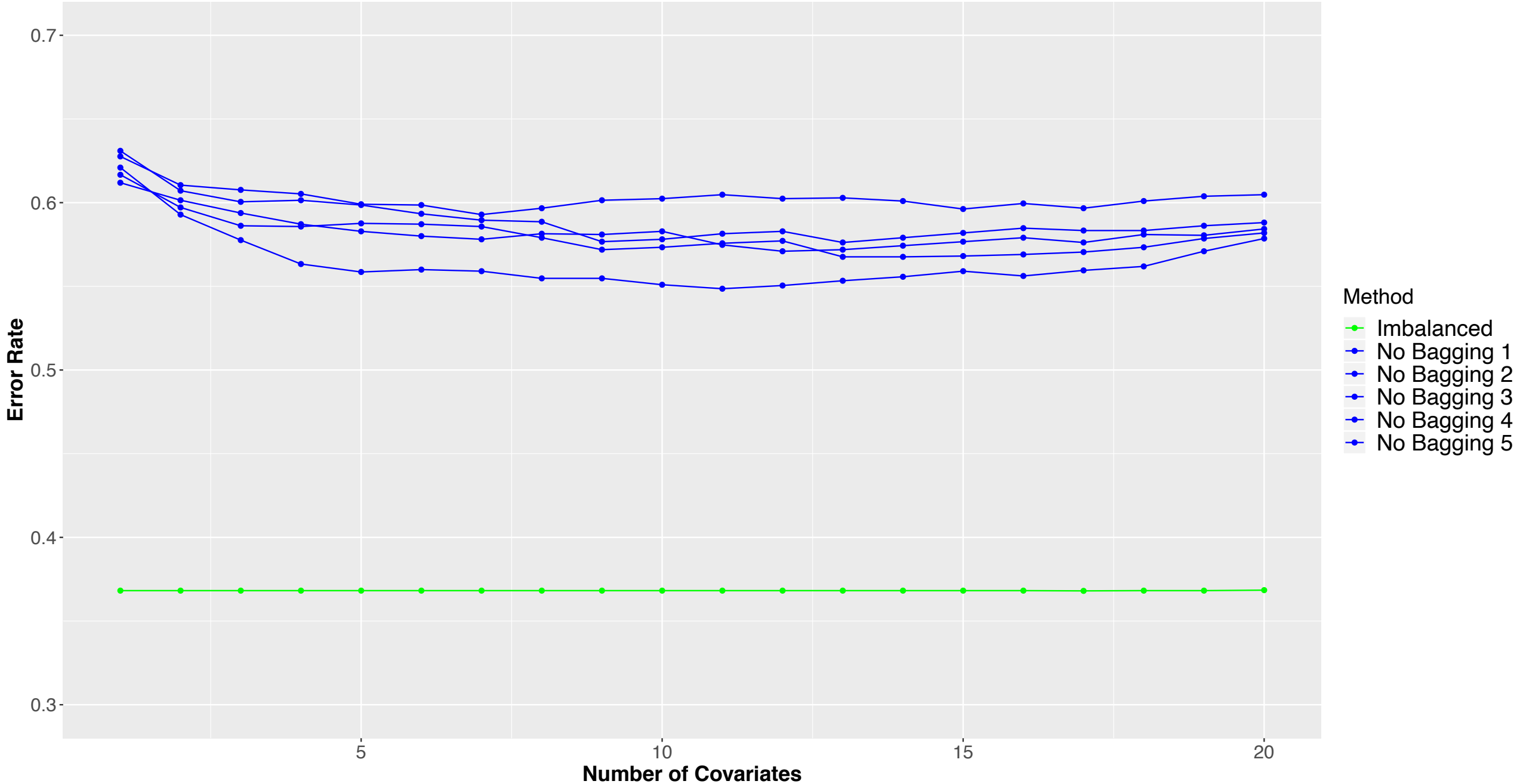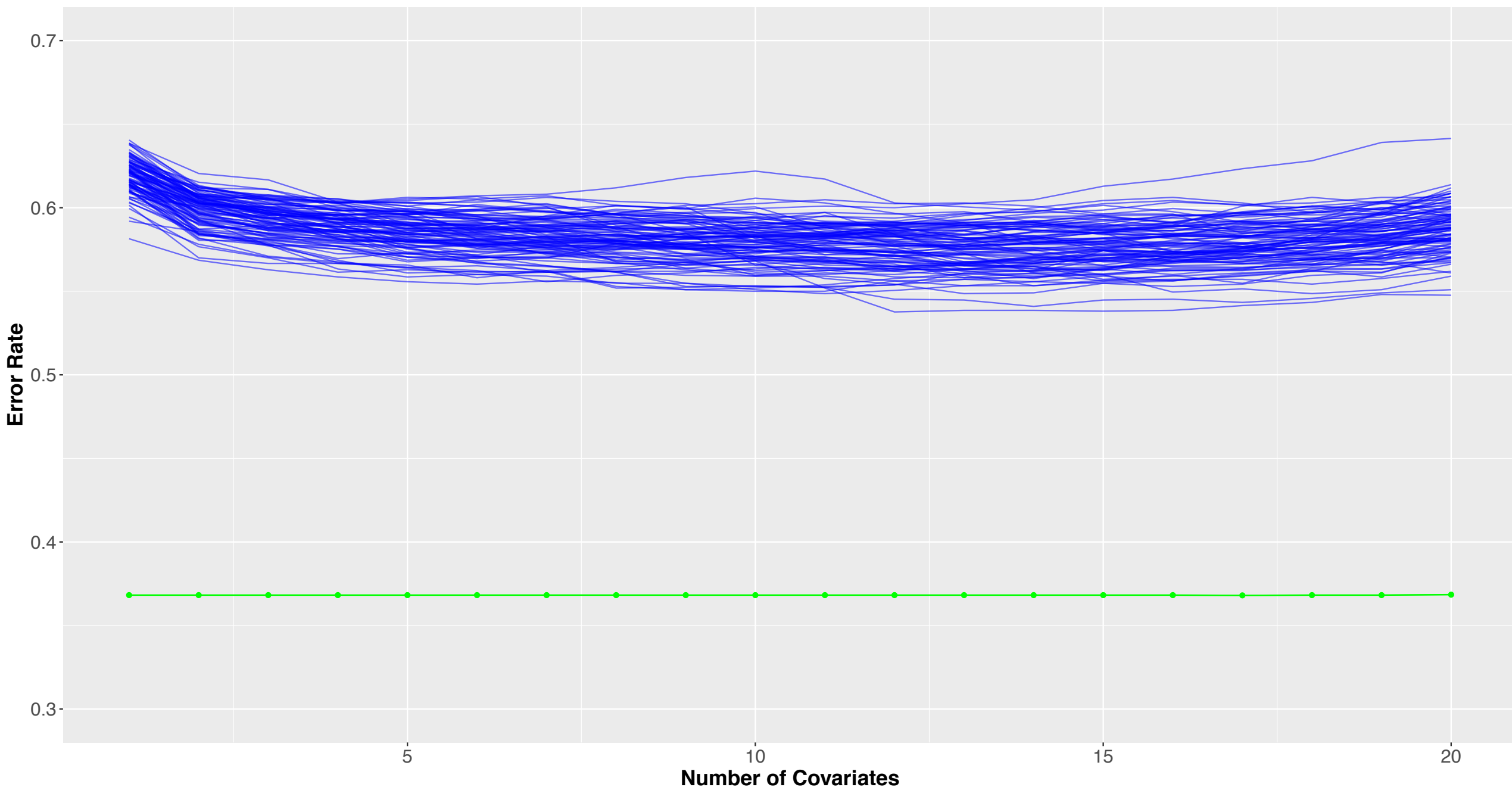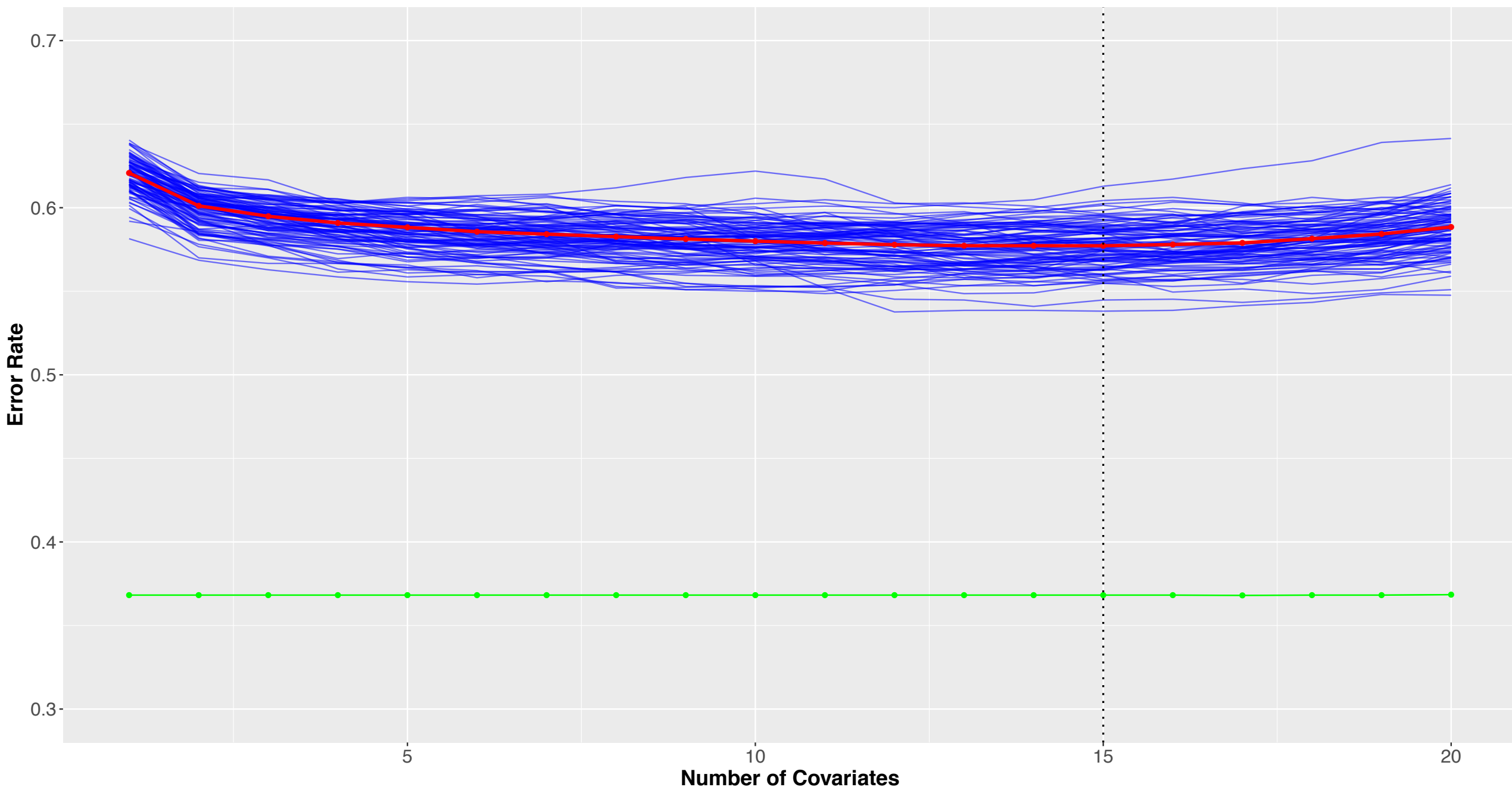Methodology illustrated with LDA

Methodology illustrated with LDA

Methodology illustrated with LDA

Methodology illustrated with LDA

Methodology illustrated with LDA

# Evaluation of Over- and Under-Bagging

- Under-bagging:
  - No loss of information[1]
- Over-bagging:
  - Reduce the risk for over-fitting[1]

- Computationally Intensive[1]

[1]According to Fernandez et al. (2018, pp. 82 - 83, 175 - 176)

# 1st Model Selection Approach

- **1)** 100 Bagged Samples **2)** 10-Fold-CV **3)** Best Subset Selection[1]

- Three Functions:
  - Best_subset_selection <− function(df train , df test , Y, FUN){ ... }
  - K_fold_CV<− function(df, Y, K, FUN){ ... }
  - Over_under_bagging <− function(df, Y, B, sample size, FUN){ ...}

- Number of estimated models:
  - $B \times K \times 2^p \approx 1\ Billion$      with $B = 100, K = 10, p = 100$

[1]According to James et al. (2017, p. 205)

# 2$^{nd}$ Model Selection Approach

- Increase the function's efficiency

- Instead of CV, use out-of-bag errors:
  - $B \times K \times 2^p \approx 100\ Million$     with $B = 100, K = 1, p = 100$

# 3[rd] Model Selection Approach

- Instead of best subset selection, new function for forward stepwise selection[1]:
  - forward_stepwise_selection <− function(df train , df test , Y, FUN){ … }
  - $B \times K \times \left( \frac{p(p+1)}{2} - 1 \right) \approx 10\ Million$ with $B = 100, K = 10, p = 100$
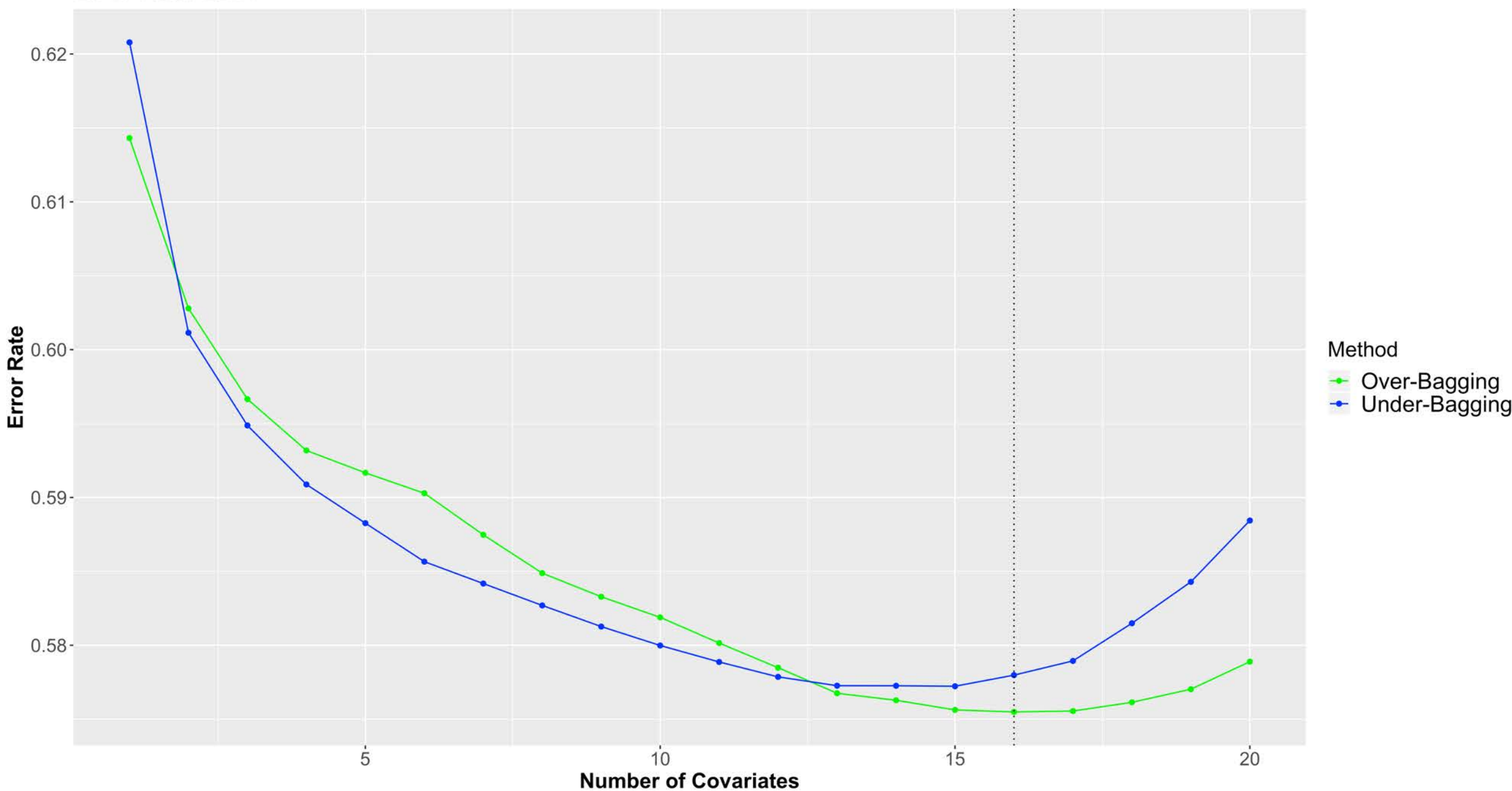
[1]According to James et al. (2017, p. 207)

# 4th Model Selection Approach

- Reducing the number of covariates to 20
  - Some demographic could be eliminated due to perfect multicollinearity (e.g. Ethnicity per census track)
  - Highest correlation to Inspection Grades
    (Not an optimal approach[1])

- Number of estimated models:
  - $B \times K \times \left( \frac{p(p+1)}{2} - 1 \right) \approx 200'000$      with $B = 100, K = 10, p = 20$

- **1)** Under- and Over-bagging **2)** 100 Bagged Samples **3)** 10-Fold-CV **4)** Best Subset Selection for only 20 Covariates

[1]According to Hastie, Tibshirani & Friedman (2013, p. 245)
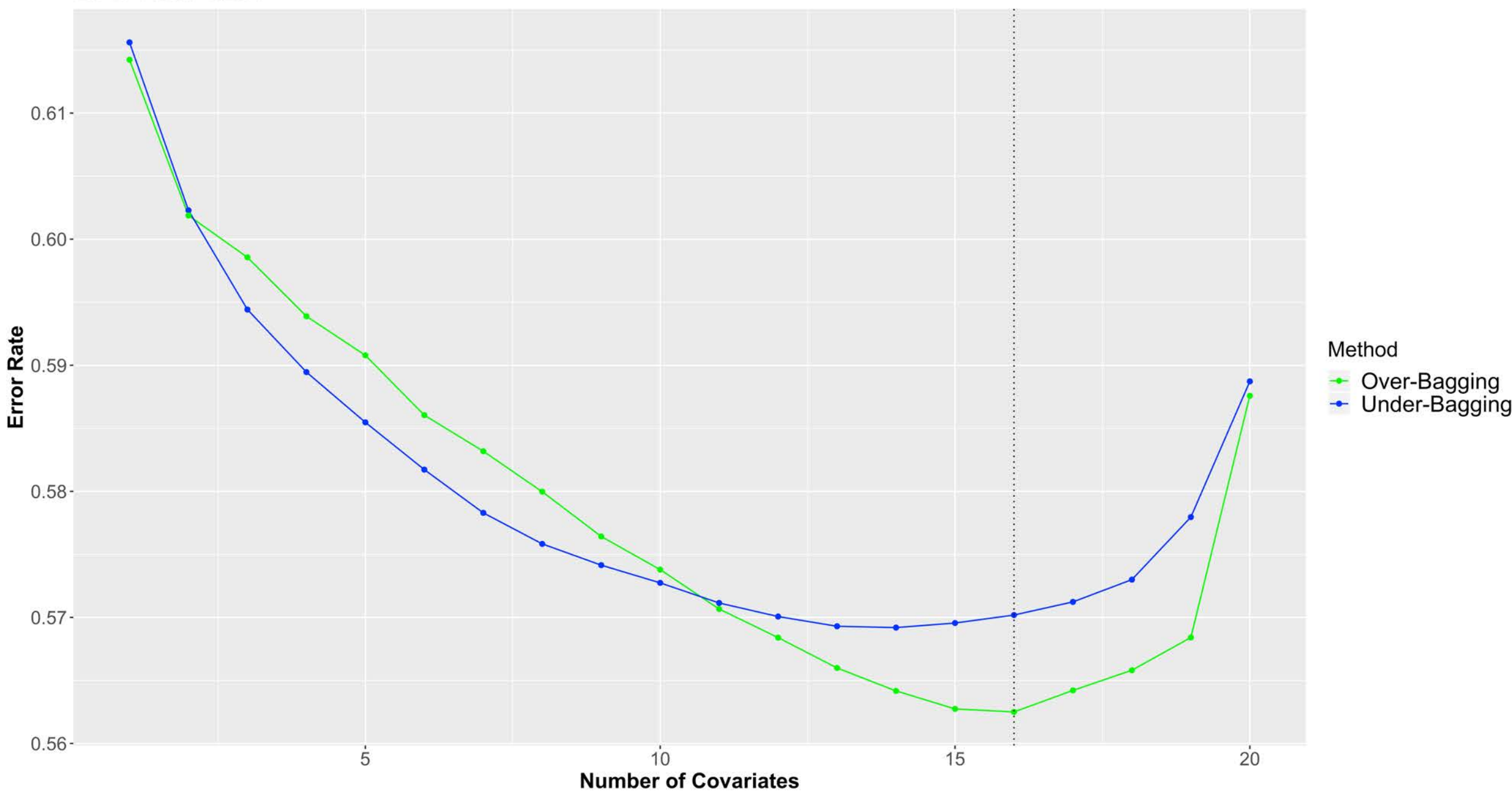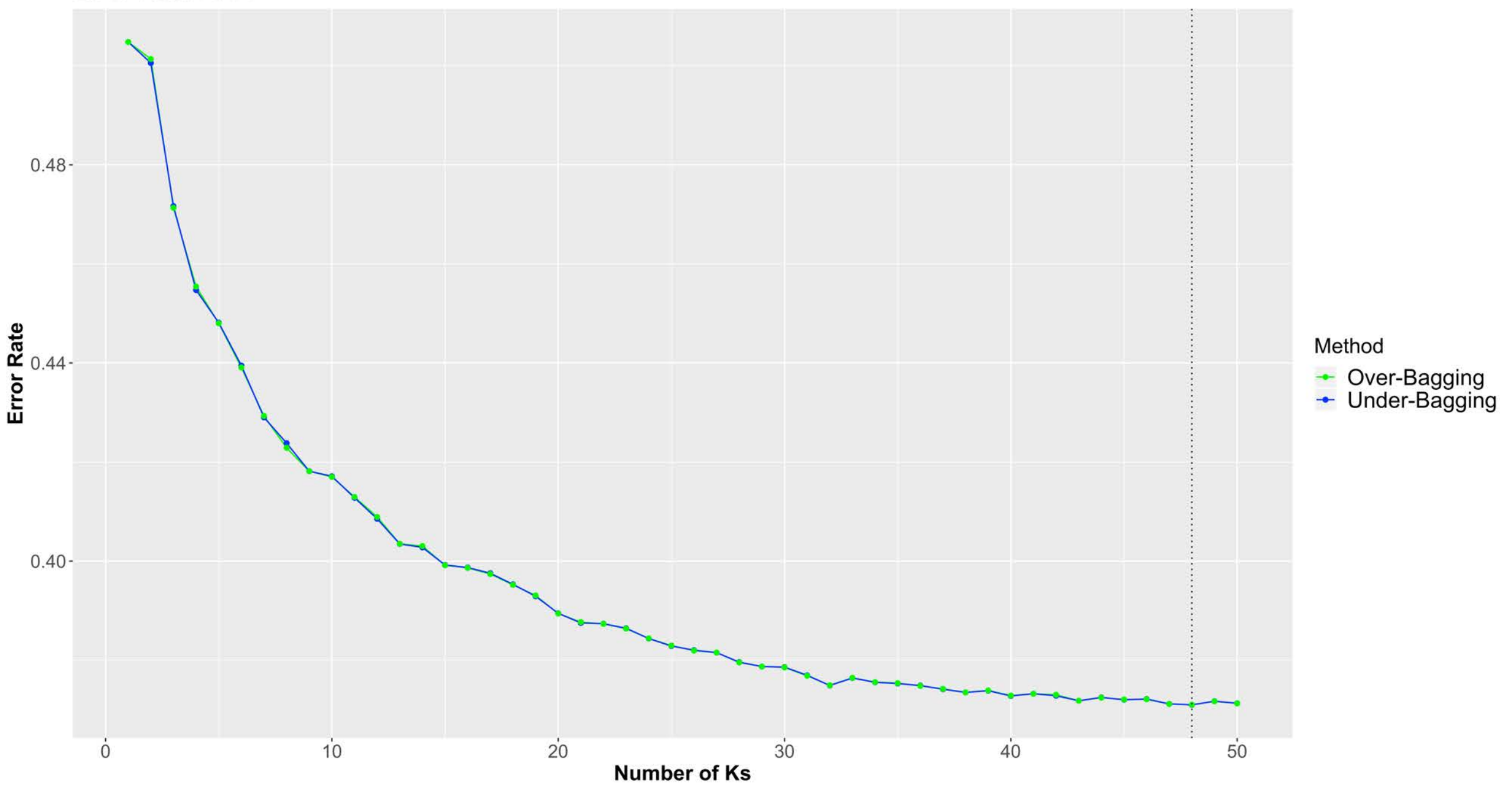
Error Rate LDA

Decision Boundaries LDA

Error Rate QDA

Decision Boundaries QDA

Error Rate KNN

Decision Boundaries KNN

# Prediction Trade-Off

## QDA

Error Rate: 54%

| Obs / Pred | A | B | C |
|---|---|---|---|
| A | 2839 | 462 | 1164 |
| B | 1260 | 173 | 594 |
| C | 228 | 26 | 103 |

## KNN

Error Rate: 66%

| Obs / Pred | A | B | C |
|---|---|---|---|
| A | 1464 | 204 | 470 |
| B | 1928 | 335 | 885 |
| C | 935 | 122 | 506 |

# Random Forest and Boosting

- Methodology needs very high computational power
  - Fewer iterations lead to higher variance

- Tendence to over-fitting with over-bagging
  - Training error of 5%

| RF over-bagging | RF under-bagging | Boosting under-bagging | Boosting over-bagging |
|---|---|---|---|
| 0.533 | 0.51 | 0.68 | 0.63 |

# Results



Final error rates

# Conclusion

- Correlation of best 20 covariates: between 0.09 – to 0.03

- More sophisticated approach to imbalance problem (e.g. Synthetic Resampling Technique)

- Opportunities for agency:
  - Use of internal data
  - Unique data from every food store

# Appendix 1: OOB Testing Sample

- Probability for a not picking observation

$$\frac{N-1}{N}$$

- Probability for a not picking N observations (with replacement)

$$\left(\frac{N-1}{N}\right)^N$$

- Probability for a not picking N observations (with replacement)

$$\lim_{N \to \infty} \left(\frac{N-1}{N}\right)^N = e^{-1} = 0.368$$