# Analyzing Elo Ratings For NCAA Men's Division 1 Hockey

Alex Tidd

Advisors: Matt Higham and Robin Lock

2025-05-02

## Table of contents

## 1 Things to note

- APA Citation
- Render as docx then pdf
- 1.5 spacing or double space
- add line numbers. On rendered file.

## 2 Abstract

There are 64 NCAA Division 1 Men's Hockey teams. In NCAA Division 1 Men's Hockey, the U.S. College Hockey Online (USCHO) provides the official rankings on the 64 total teams using a system that relies on expert votes. However, this method is not perfect as there is no formal quantitative analytics involved in the voting. To improve on this, a chess ratings system, called Elo, is commonly used in many sports to quantitatively rate players, weighing strength of opponent and recency of match. In this project, we modify the Elo system for NCAA Division 1 Men's Hockey by adding in weights for game goal differential (so that games in which the score differential was large potentially result in a larger bump

in Elo for the winning team) and home-ice advantage (so that the team playing on home-ice has an adjusted probability of winning the game).

## 3 Introduction

- Goal
- Outline of write up
- Background of Elo

NCAA Men's Division 1 Ice Hockey has 64 teams. This inherently makes ranking teams extremely hard. It is impossible to rank teams in the same manner as pro leagues like the NHL and junior leagues like the 3 CHL leagues, where teams are ranked based off of a record-point system.

To rank teams, the current method is an "expert vote" system done by US College Hockey Online (USCHO). This method takes votes by "experts" and ranks teams based on how many votes each team receives. In general, these rankings are pretty accurate, however it lacks any true quantitative analysis, instead relying on the opinions of "experts".

This project aims to solve the issue presented by using a ratings system used in chess to accurately rate and rank players, called Elo, and use in for collegiate hockey. In chess, Elo ratings calculate an "expected outcome" using 2 players' ratings. This expected outcome is essentially a "win probability" for each player. After the expected outcome is calculated, it is compared to the actual outcome. Depending on the outcome, a players new rating will either be increased or decreased. The benefit to using Elo, is that the rating system takes into consideration your strength, the opponents strength and the recency of the match. Big wins/losses garner big adjustments in rating, whereas an expected win/loss won't inflate/tarnish a rating. The equations are shown below:

Expected outcome:

$$E_a = \frac{1}{1 + 10^{\frac{R_b - R_a}{400}}}$$

Rating update:

$$R^`_a = R_a + k(\text{outcome} - E_a)$$

In this case, k is an "update factor" that scales how many points are added or subtracted from a team's rating after the event of a win or loss.

The goal for this project is to use this base Elo ratings system, incorporate factors such as home ice advantage and goal differential, as well as optimize k, to make a ratings model that can accurately predict outcome and rank teams accordingly. This paper will take you through the steps of finding data, creating a ratings function, and optimizing our constants.

## 4 Data

There were two main data sets used for this project. The full 2023-2024 NCAA Men's Division 1 Ice Hockey schedule, which after wrangling, included 1166 games with variables: date, game_type, away_team, away_score, home_team, home_score, overtime, neutral_site, score_diff, and outcome. The score_diff and outcome variables are in reference to the home team, with outcome being either 1: win, 0.5: tie, or 0: loss. The next data frame was the entire 2024-2025 NCAA Men's Division 1 Ice Hockey schedule. The schedule contains 1153 games which had the same variables and setup as the 2023-2024 season. Both data frames were scraped from the website: College Hockey News.

A third data frame for initial rankings was created from scratch using knowledge from the final 2022-2023 season. Initial rankings were created by ranking teams from 2000 to 1300. The top 8 teams from the 2022-2023 season were assigned the rating of 2000, every 8 teams, the assigned rating would decrease by 100. This initial ratings file was used as the initial ratings file to run with the 2023-2024 season to get accurate ratings for the 2024-2025 season.

To get these "accurate" rating for the 2024-2025 schedule. The 2023-2024 schedule was run through the Elo function with k = 100, a value used based off a general exploration which gave relatively accurate rankings compared to the final 2023-2024 season. After, the ratings where scaled using the equation below:

$$rating^` = (rating \cdot 0.7) + 450$$

This adjustment of rankings is used to adjust for player turnover, new coaching hires, and off season improvement. The formula used is based off of the formula used by Fivethirtyeight in their NHL Elo ratings. These final adjusted ratings are what was used for the intial rating in the 2024-2025 season.

## 5 Elo Model

To add goal differential into the Elo system, it was determined that it should be placed in the "rating update" portion of the Elo system, due to goal differential being a "post-game" statistic. On the other hand, home ice advantage affects the predicted probability of a home team winning, therefore, it was added to the "expected outcome calculation" portion of the Elo system. Home ice advantage was also incorporated as a simple point "boost" for the home team instead of a multiplier, based off of FiveThirtyEight's NHL Elo model.

In order to obtain the "best" values for each parameter, a grid search method was used to optimize each of the three parameters to lower the mean absolute residual of the season. In this case the "absolute residual" is calculated by the absolute value of the difference of the game outcome and the expected outcome of the home team. Essentially we want the smallest difference between actual outcome and expected outcome. The smaller the difference, the more accurate our function is predicting actual outcomes of games.

103     The un-optimized Elo function is as follows:

104
$$E_{home} = \frac{1}{1 + 10^{\frac{R_{away} - (R_{home} + \text{homeIce})}{400}}}$$

105
$$E_{away} = \frac{1}{1 + 10^{\frac{(R_{home} + \text{homeIce}) - R_{away}}{400}}}$$

106     $$R\text{`}_{home} = R_{home} + k(d(\text{scoreDiff}) + (\text{outcome} - E_{home}))$$

107     $$R\text{`}_{away} = R_{away} + k(d(-\text{scoreDiff}) + (\text{outcome} - E_{away}))$$

108     Using this model a general exploration to see what possible values of k, d, and home_ice
109     could render a lower mean absolute residual. Based off the exploration, k is found at
110     roughly 100, d at roughly 50, and home_ice at roughly 50. It should be noted that k was
111     explored first with d = 0, and home_ice = 0. Next was home_ice, which was found at varying
112     weights of k and d = 0. d was then explored with k = 100, and home_ice = 0. the only reason
113     k is held in all three cases, is that if k = 0, the function would not be able to update rating as
114     the update portion would always be the initial rating + 0.

115     These explorations provided essential knowledge for the grid search as initial optimization
116     could be found using smaller ranges of values, significantly decreasing the run time of
117     each optimization. This background knowledge cut run time of one optimization dawn from
118     16.3 days to roughly 30 minutes as each optimization grid search used 1000 combinations
119     of the three variables instead of 4.5 million.

120     Using this method k was found to be optimized at 37.22222, home_ice = 53.33333, and d =
121     40.55556, this provided a mean absolute residual of 0.3515844. This looked very promising
122     at first glance. However, upon deeper inspection, two major concerns were found. First,
123     the bottom four teams had end of season ratings in the negatives. Second, there were
124     many residuals of 1 and 0. This raised the most concern as expected outcome is bounded
125     by 0 < expected outcome < 1, meaning we should never see residuals of exactly 1 or 0.
126     Upon deeper inspection, it was found that there were many games were expected
127     outcome was to the effect of 0.9 repeating or 1e-20. Which whilst technically within our
128     bounds, is highly unrealistic. A good model that predicts win probability should never give
129     a pregame win probability of nearly 100% or 0%. Something needed to be changed with the
130     model. Goal differential was looked at first.

131     Upon inspecting and researching different ways to incorporate goal differential in an Elo
132     function, FiveThirtyEight was looked at again. In their NHL model, goal differential was
133     incorporated as shown below:

134
$$\left(0.6686 \cdot ln(\text{scoreDiff})\right) + 0.8048)$$

135     The new update function now looks as follows:

136  $$R^`\_{home} = R\_{home} + k(((0.6686\cdot ln(\text{scoreDiff}))+0.8048)(\text{outcome} -
137  E\_{home}))$$

138  $$R^`\_{away} = R\_{away} + k(((0.6686\cdot ln(\text{-scoreDiff}))+0.8048)(\text{outcome} -
139  E\_{away}))$$

140  The addition of 0.8048 is put in as the factor for score_diff = 1 as the natural logarithm of 1
141  is 0, which would then render the entire update portion of the function useless, running
142  into the same issue if k = 0.

143  Now that goal factor is optimized, all that was left was to rerun the grid search to optimize k
144  and home_ice. Using the grid search method, k is optimized at 88 and home_ice optimized
145  at 40, which gives a mean absolute residual of 0.3876531. Whilst this may not have a mean
146  absolute residual lower than the first model, it maintains that all teams keep ratings above
147  zero and keep predicted outcomes in a more realistic zone.

## 6 Results

149  Results from this optimization show that k = 88, home_ice = 40, and the calculation to
150  utilize score differential $= \left(0.6686 \cdot ln(\text{scoreDiff})\right) + 0.8048)$. The plots below were used to
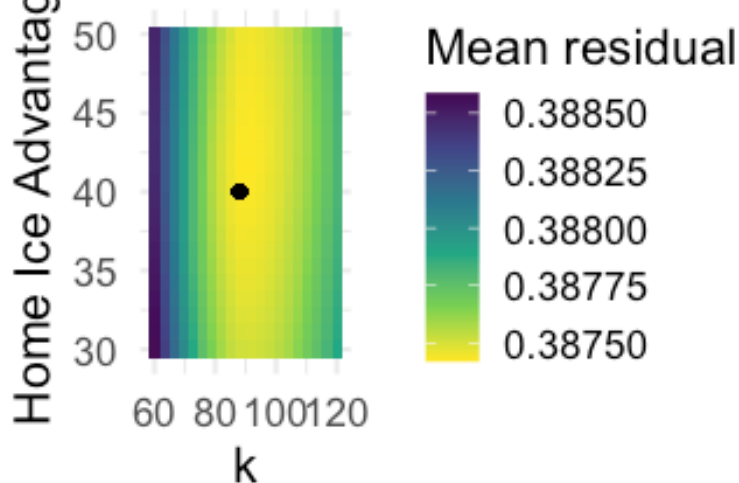151  further imply this.

152  Grid search for k and home_ice:

```
library(furrr)
library(progressr)
library(scico)
plan(multisession)
handlers("progress")
options(progressr.enable = TRUE)

grid_100 = expand.grid(k = seq(60, 120, length.out = 20), home_ice = seq(30,
50, length.out = 20), d = seq(0, 100, length.out = 1))

mean_residuals_100 = with_progress({future_pmap_dbl(grid_100, \ (k, home_ice,
d) update_rankings_residuals(season = schedule_reg, end_date = "2025-03-25",
ratings = rankings2324, k = k, home_ice = home_ice, d = d), .progress =
TRUE)})

residual_100_df <- grid_100 |> mutate(mean_residual = mean_residuals_100)

ggplot(data = residual_100_df, aes(x = k,
                                   y = home_ice)) +
  geom_tile(aes(fill = mean_residual)) +
  geom_point(aes(x = 88, y = 40), color = "black", fill = "black") +
  scale_fill_viridis_c(option = "D",
                       ##limits = c(0.36, 0.41),
                       oob = scales::squish,
                       name = "Mean residual",
```

```
178                         direction = -1) +
179     labs(x = "k",
180          y = "Home Ice Advantage",
181          title = "Optimization of k and Home Ice Advantage",
182          caption = "k is optimized at 88, home ice advantage at 40, \nand goal
183  differential factor \nof  0.6686 * log(abs(score_diff)) + 0.8048") +
184     theme_minimal(base_size = 16) +
185     theme(legend.position = "right",
186           plot.title = element_text(hjust = 0.5),
187           plot.caption = element_text(hjust = 0.5),
188           plot.margin = margin(t = 40, b = 20, l = 60, r = 60))
```



k is optimized at 88, home ice advantage at 40,
and goal differential factor
of  0.6686 * log(abs(score_diff)) + 0.8048

189

190    Heat map showing the grid search to optimize update factor, k, and home ice advantage,
191    home_ice. k and home_ice are optimized to lower the mean absolute value of the
192    difference between expected outcome and actual outcome. Goal differential factor was
193    held at $(0.6686 \cdot ln(\text{scoreDiff})) + 0.8048$. k = 88, home_ice = 40.
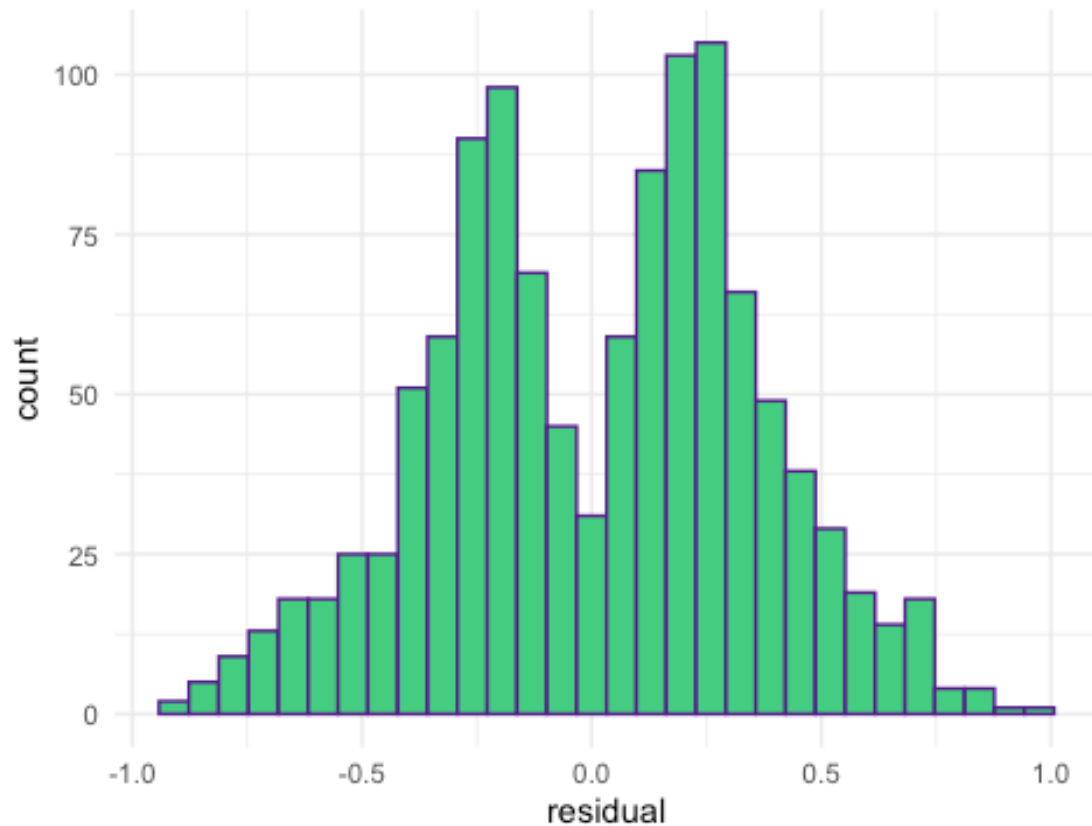
194    Plot of residuals:

```
195  apr15 = apr_15_ranking |> bind_rows()
196
197  schedule_elo = schedule |>
198    mutate(home_elo = NA) |>
```

```r
199     mutate(away_elo = NA)
200
201   schedule_apr15 = left_join(schedule_elo, apr15,
202                           by = join_by(date == date, home_team == Team)) |>
203     mutate(home_elo = rating) |>
204     select(-rating)
205
206   merged_sched_apr15 = left_join(schedule_apr15, apr15,
207                           by = join_by(date == date, away_team == Team)) |>
208     mutate(away_elo = rating) |>
209     select(-rating)
210
211   schedule_full_apr15 = merged_sched_apr15 |>
212     mutate(outcome_away = abs(outcome - 1)) |>
213     ## Calculating expected outcome variable for home and away team
214     mutate(exp_home = 1/(1 + 10^((away_elo - home_elo)/400))) |>
215     mutate(exp_away = 1/(1 + 10^((home_elo - away_elo)/400))) |>
216     ## Using expected outcome variable to generate new Elo ratings based on
217   actual outcome and expected outcome
218     mutate(elo_new_home = home_elo + 100 * (outcome - exp_home)) |>
219     mutate(elo_new_away = away_elo + 100 * (outcome_away - exp_away))
220
221   ## Making a residual column
222   schedule_full_apr15 <- schedule_full_apr15 |>
223     mutate(residual = outcome - exp_home) |>
224     mutate(abs_residual = abs(residual))
225
226   ggplot(data = schedule_full_apr15, aes(x = residual)) +
227     geom_histogram(color = "purple4", fill = "seagreen3") +
228     labs(caption = "Positive residual past 0.5 indicates model predicting a
229   home loss when actual result is a home win. Negative residual beyond -0.5
230   indicate a predicted home wins with an observed home loss") +
231     theme_minimal()
```

Negative residual beyond -0.5 indicate a predicted home wins with an observed home loss

232

233    Plot of residuals shows all the differences between expected outcome and actual
234    outcome. Positive residual past 0.5 indicates model predicting a home loss when actual
235    result is a home win. Negative residual beyond -0.5 indicate a predicted home wins with an
236    observed home loss. We see that the model has large amounts of residuals around the
237    |0.25| which means that the expected outcome is still within the range of the actual
238    outcome. For example, actual outcome is 0, and expected is 0.25, the model still predicted
239    a loss. The only grey area is around ties, since ties are 0.5, a |0.25| difference could be
240    construed as tie or a win/loss.

241    Binned expected value graph:

```
242    schedule_full_apr15 |> summarise(mean_resid = mean(abs_residual, na.rm =
243    TRUE))

244    # A tibble: 1 × 1
245      mean_resid
246           <dbl>
247    1      0.294

248    prop_wins15 <- schedule_full_apr15 |>
249      mutate(binned_exp = floor(exp_home / 0.1) * 0.1 + 0.05) |>
250      group_by(binned_exp) |>
251      summarise(win_prop = mean(outcome, na.rm = TRUE),
```

```
252                    totalgames = n()) |>
253      filter(!is.na(binned_exp))
254
255    modgdha = lm(win_prop ~ binned_exp, data = prop_wins15, weights = totalgames)
256    summary(modgdha)
257
258    Call:
259    lm(formula = win_prop ~ binned_exp, data = prop_wins15, weights = totalgames)
260
261    Weighted Residuals:
262        Min      1Q   Median      3Q      Max
263    -1.03772 -0.41495  0.01982  0.49843  0.79635
264
265    Coefficients:
266                Estimate Std. Error t value Pr(>|t|)
267    (Intercept) -0.11495    0.04086  -2.813   0.0227 *
268    binned_exp   1.24411    0.06996  17.782 1.02e-07 ***
269    ---
270    Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
271
272    Residual standard error: 0.6344 on 8 degrees of freedom
273    Multiple R-squared:  0.9753,    Adjusted R-squared:  0.9722
274    F-statistic: 316.2 on 1 and 8 DF,  p-value: 1.024e-07
275    ggplot(data = prop_wins15, aes(x = binned_exp,
276                                   y = win_prop,
277                                   size = totalgames)) +
278      geom_point(color = "black", shape = 16) +
279      geom_smooth(aes(color = "Fitted Model",
280                  weight = totalgames), method = "lm", se = FALSE, size =
281    1.2) +
282      geom_abline(data = data.frame(1),
283                  aes(color = "Expected Linear Model",
284                      linetype = "Expected Linear Model"),
285                  slope = 1, intercept = 0, linetype = 2, size = 1) +
286      scale_color_manual(name = "Model",
287                         values = c("Fitted Model" = "gold",
288                                    "Expected Linear Model" = "black")) +
289      labs(title = "Proportion of Home Team Wins \nfrom Home Expected Outcome",
290           x = "Expected Outcome",
291           y = "Proportion of \nGames Won",
292           caption = "Size of points indicate more games played. \nModel weighs
293    point based off of amount of games played") +
294      guides(size = "none") +
295      theme_minimal(base_size = 16) +
296      theme(legend.position = "right",
297            legend.background = element_rect(fill = "white", color = NA),
298            legend.title = element_text(size = 12),
299            legend.text = element_text(size = 10),
```
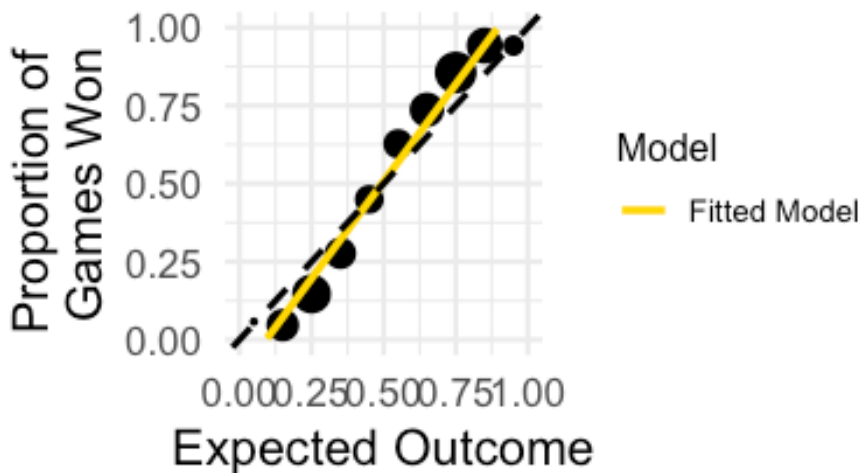
```
300         plot.title = element_text(hjust = 0.5),
301         plot.caption = element_text(hjust = 0.5),
302         plot.margin = margin(t = 40, b = 20, l = 40, r = 40)) +
303    xlim(c(0, 1)) +
304    ylim(c(0, 1))
```



Proportion of Home Team Wins from Home Expected Outcome

Size of points indicate more games played.
Model weighs point based off of amount of games played

305

306   This plot bins every game played in the 2024-2025 season by every 0.1. Then uses the
307   amount of games played in each bin to weigh each point. The y-axis is the proportion of
308   wins for each binned proportion. We expect a slope of 1. This means that for example, at
309   expected outcome of 0.9 we expect 90% of the games are won. In this case the slope of our
310   fitted model is 1.24 with an intercept of -0.11. This further shows the accuracy of the
311   optimized Elo function.

312   Final rankings are shown below:

```
313   apr15 = apr_15_ranking |> bind_rows()
314
315   apr15_lagged = apr15 |> group_by(date) |>
316     summarise(last_date = last(date)) |>
317     mutate(lag_date = lag(last_date)) |>
318     select(-last_date)
319
320   apr15_full = left_join(apr15, apr15_lagged, join_by(date == lag_date)) |>
```
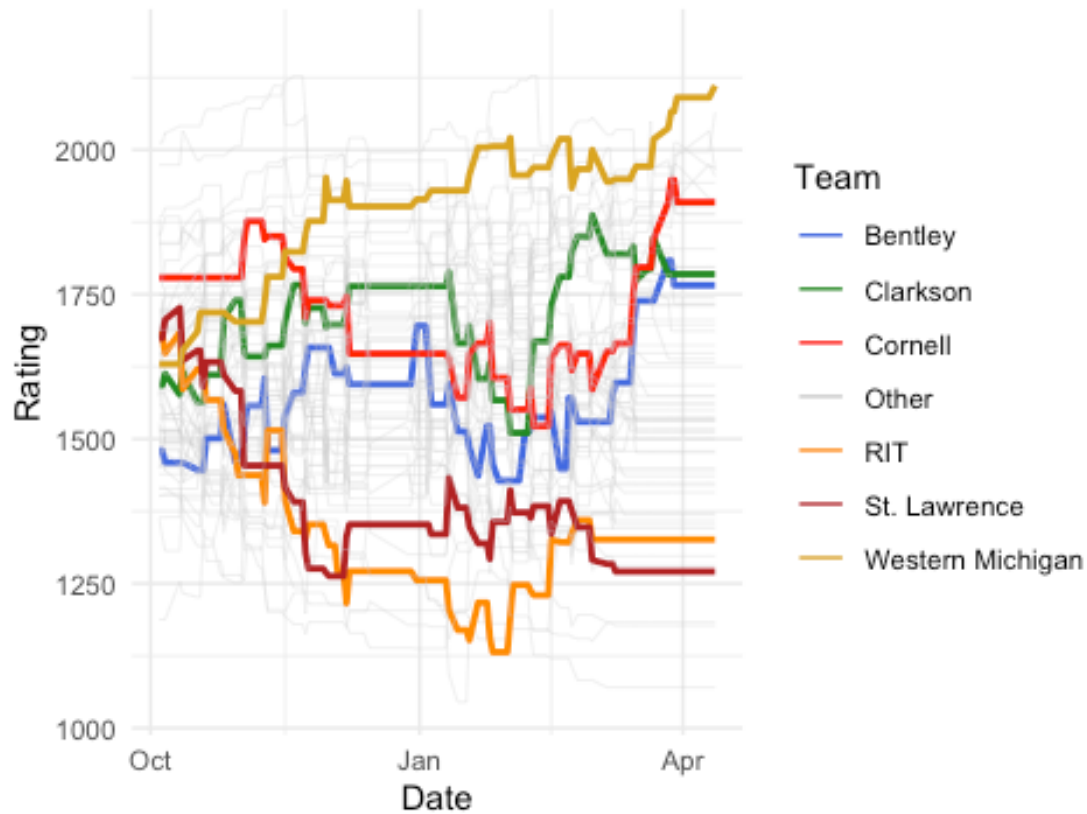
```
321    select(-date) |>
322    rename(date = date.y)
323
324  highlight = c("St. Lawrence", "Western Michigan", "Clarkson", "Bentley",
325  "RIT", "Cornell")
326
327  highlighted_color = c(
328    "St. Lawrence" = "firebrick",
329    "Western Michigan" = "goldenrod",
330    "Clarkson" = "forestgreen",
331    "Bentley" = "royalblue",
332    "RIT" = "darkorange",
333    "Cornell" = "red")
334
335  apr15_color = apr15_full |> mutate(highlight = if_else(Team %in% highlight,
336  Team, "Other"))
337
338  ggplot(data = apr15_color, aes(x = date,
339                                 y = rating,
340                                 group = Team)) +
341    geom_line(aes(color = highlight,
342                  alpha = highlight,
343                  linewidth = highlight)) +
344    scale_color_manual(values = c("Other" = "grey80",
345                                    highlighted_color),
346                         name = "Team") +
347    scale_alpha_manual(values = c("Other" = 0.5,
348                                    setNames(rep(1, length(highlighted_color)),
349                                             names(highlighted_color))),
350                         guide = "none") +
351    scale_linewidth_manual(values = c("Other" = 0.2,
352                                      setNames(rep(1,
353  length(highlighted_color)),
354                                                 names(highlighted_color))),
355                           guide = "none") +
356    theme_minimal() +
357    labs(color = "Team",
358         title = "Full Season Rankings",
359         x = "Date",
360         y = "Rating") +
361    theme(legend.position = "right",
362          plot.title = element_text(hjust = 0.5))
```

## Full Season Rankings



363

364
```
print(n = 64, apr_15_ranking[[110]] |> arrange(desc(rating)))
```

365  # A tibble: 64 × 3
366     Team              rating date
367     <chr>              <dbl> <date>
368   1 Western Michigan   2186. 2025-04-12
369   2 Boston University  1990. 2025-04-12
370   3 Connecticut        1970. 2025-04-12
371   4 Penn State         1964. 2025-04-12
372   5 Denver             1954. 2025-04-12
373   6 Michigan State     1943. 2025-04-12
374   7 Boston College     1915. 2025-04-12
375   8 Cornell            1909. 2025-04-12
376   9 Minnesota State    1902. 2025-04-12
377  10 Maine              1890. 2025-04-12
378  11 Massachusetts      1888. 2025-04-12
379  12 North Dakota       1829. 2025-04-12
380  13 Minnesota          1798. 2025-04-12
381  14 Quinnipiac         1796. 2025-04-12
382  15 Arizona State      1786. 2025-04-12
383  16 Clarkson           1785. 2025-04-12
384  17 Omaha              1774. 2025-04-12
385  18 Bentley            1766. 2025-04-12

```
386    19 St. Thomas          1744. 2025-04-12
387    20 Providence         1734. 2025-04-12
388    21 Michigan           1733. 2025-04-12
389    22 Northeastern       1730. 2025-04-12
390    23 Dartmouth          1706. 2025-04-12
391    24 Notre Dame         1697. 2025-04-12
392    25 Ohio State         1687. 2025-04-12
393    26 Long Island        1662. 2025-04-12
394    27 New Hampshire      1648. 2025-04-12
395    28 Holy Cross         1638. 2025-04-12
396    29 Harvard            1632. 2025-04-12
397    30 Bowling Green      1593. 2025-04-12
398    31 Colorado College   1575. 2025-04-12
399    32 Bemidji State      1574. 2025-04-12
400    33 Colgate            1563. 2025-04-12
401    34 Minnesota-Duluth   1558. 2025-04-12
402    35 Wisconsin          1549. 2025-04-12
403    36 Alaska             1538. 2025-04-12
404    37 Vermont            1536. 2025-04-12
405    38 Augustana          1533. 2025-04-12
406    39 St. Cloud State    1533. 2025-04-12
407    40 Brown              1530. 2025-04-12
408    41 Sacred Heart       1518. 2025-04-12
409    42 Merrimack          1514. 2025-04-12
410    43 Mass.-Lowell       1511. 2025-04-12
411    44 Union              1500. 2025-04-12
412    45 Princeton          1479. 2025-04-12
413    46 Army               1477. 2025-04-12
414    47 Lake Superior      1449. 2025-04-12
415    48 Ferris State       1444. 2025-04-12
416    49 Air Force          1435. 2025-04-12
417    50 Michigan Tech      1424. 2025-04-12
418    51 Lindenwood         1420. 2025-04-12
419    52 American Int'l     1407. 2025-04-12
420    53 Stonehill          1372. 2025-04-12
421    54 Niagara            1369. 2025-04-12
422    55 Rensselaer         1359. 2025-04-12
423    56 Canisius           1352. 2025-04-12
424    57 Alaska-Anchorage   1345. 2025-04-12
425    58 RIT                1326. 2025-04-12
426    59 Northern Michigan  1298. 2025-04-12
427    60 Yale               1281. 2025-04-12
428    61 St. Lawrence       1271. 2025-04-12
429    62 Robert Morris      1183. 2025-04-12
430    63 Miami              1177. 2025-04-12
431    64 Mercyhurst         1070. 2025-04-12
```

432    1.  In these final season ranking our Elo system ranked Western Michigan as the top
433        team, with Boston University second. In fact, of the top 4 teams in our model, 3 of

434    them made the Frozen Four, with Denver finishing in 5th. Western Michigan
435    ultimately ended up winning the National Championship.

## 7 Conclusion

437 In order to make a better ratings system in NCAA Men's Division 1 Ice Hockey, an Elo style
438 system was created. Using goal differential, an update factor, and home ice advantage, an
439 Elo model was optimized using a grid search method to get the mean absolute value of
440 game residuals down to 0.388, and successfully ranked the four teams to make the Frozen
441 Four in the top five and successfully ranked the national champions, Western Michigan, in
442 the top spot for the end of 2024-2025 season rankings.

443 This model allows teams to be ranked quantitatively and use factors such as strength of
444 schedule and quality of win/loss to accurately scale the effect of each win or loss. In the
445 future, I would like to try to add more parameters in to the model to see how big of an effect
446 things like, OT and neutral site have on expected outcome. Mostly I would like to break
447 down score differential into goals for and goals against to see if having good defense is
448 more important than good offense and vice versa. My biggest regret is not being able to
449 optimize the score differential parameter myself, and if given more time would do this.

## 8 Code Appendix

451 Libraries:

```
452 library(elo)
453 library(dplyr)
454 library(tidyverse)
455 library(cowplot)
456 library(ggrepel)
457 library(lubridate)
458 library(rvest)
459 library(here)
460 library(forcats)
461 library(progressr)
462 library(furrr)
463 library(vctrs)
464 library(purrr)
```

465 Scraping function:

```
466 ##Function to load in schedule
467 scrape_men <- function(season = "20232024"){
468   ## URL for schedule data frame
469   url_hockey <-paste("https://www.collegehockeynews.com/schedules/?season=",
470 season, sep = "")
471   ## Selecting which schedule table to grab
472   tab_hockey <- read_html(url_hockey) |>
473     html_nodes("table")
474
```

```r
475     ## The website likes to switch which table it uses. If function doesn't
476   work try changing which table number you select
477     stats_dirty <- tab_hockey[[1]] |> html_table()
478
479     ## Creating regex for date, and conference to make date and conference
480   columns in dataframe
481     regex_date <- "October|November|December|January|February|March|April"
482     regex_conference <- "Atlantic Hockey|Big Ten|CCHA|ECAC|Hockey
483   East|NCHC|Ind|Exhibition|Non-Conference"
484     ## Combining regexs with original table so that the original scraped
485   dataframe has date and conferene as variables
486     stats_regex <- stats_dirty |> mutate(date = if_else(str_detect(X1,
487   regex_date),
488                                                     true = X1, false =
489   NA_character_),
490                                         conference = if_else(str_detect(X1,
491   regex_conference),
492                                                     true = X1, false
493   = NA_character_))
494
495     ## Filling in respective dates and conferences
496     stats_filled <-stats_regex |> fill(date, .direction = "down") |>
497       ## Selecting date and congference so they show up as X1 and X2 in the
498   dataframe
499       fill(conference, .direction = "down") |> select(date, conference,
500   everything())
501     ##filtering out anywhere that a conference value is undetected (Game
502   category, not an actual game played)
503     stats_filled_cleaner <- stats_filled |> filter(!str_detect(X1, regex_date)
504   &
505                                                     !str_detect(X1,
506   regex_conference))
507     print(head(stats_filled_cleaner))
508
509     ## Dataframe is now in a format that is able to be worked on. Now creating
510   specific variables that we want to look at
511     ## Selecting first 8 columns
512     schedule_new <- stats_filled_cleaner |> select(date, conference, X1, X2,
513   X3, X4, X5, X6) |>
514       ## Taking out first two rows (no data in them). Renaming columns to match
515   what their variable is.
516       slice(-1 , -2) |> rename(game_type = conference, away_team = X1,
517   away_score = X2, location_marker = X3, home_team = X4, home_score = X5,
518   overtime = X6) |>
519       ## Take out the day of the week in our date columns as we don't need to
520   know if a game was played on a Monday per-se.
521       separate(col = date, into = c("weekday", "dm", "y"),
522               sep = ", ") |>
523     unite("new_date", c(dm, y),
```

```
524             sep = " ") |>
525        select(-weekday) |>
526        ##making date column into a <date> variable
527        mutate(date = mdy(new_date)) |>
528        ## Taking out the <chr> date variable
529        select(-new_date) |>
530        select(date, everything()) |>
531        ## Filtering out where there is no away team since that means no game was
532    played
533        filter(away_team != "") |>
534        ## Filtering out exhibition games since we aren't looking at exhibition
535    games
536        filter(game_type != "Exhibition") |>
537        ## Turning scores from <chr> to <dbl> variables
538        mutate(away_score = as.double(away_score)) |>
539        mutate(home_score = as.double(home_score)) |>
540        ## creating a variable to indicate if a game was played at a neutral site
541        mutate(neutral_site = case_when(location_marker == "vs." ~ 1,
542                                        location_marker == "at" ~ 0)) |>
543        ## Making the neutral_site variable as <lgl>
544        mutate(neutral_site = as.logical(neutral_site)) |>
545        ## taking out location_marker
546        select(-location_marker) |>
547        ## Making a logical overtime variable. Note we are not differentiating
548    between OT and 2OT
549        mutate(overtime = case_when(overtime == "" ~ 0,
550                                     overtime == "ot" ~ 1,
551                                     overtime == "2ot" ~ 1)) |>
552        mutate(overtime = as.logical(overtime)) |>
553        ##Filtering out NA "overtime" values as this indicates no game played,
554    since overtime will either be TRUE or FALSE
555        filter(!is.na(overtime))|>
556        ## Creating a score differential variable to indicate a win, loss, or tie
557    for the home team. If we know the outcome for the home team, we know the
558    outcome for the away team.
559        mutate(score_diff = home_score - away_score) |>
560        ## making an outcome variable for home team so ties get input as 0.5,
561    wins get input as 1, and loss get input as 0.
562        mutate(outcome =
563                case_when(score_diff == 0 ~ "0.5",
564                          score_diff > 0 ~ "1",
565                          score_diff < 0 ~ "0")) |>
566        ## turning score_diff from <chr> to <dbl>
567        mutate(outcome = as.double(outcome)) |>
568        ## Filtering out games where D1 team played against D3 teams as these are
569    exhibition as well
570        filter(game_type != "Non-Conference v. D3")
571
572      ## Tidy schedule is returned
573      return(schedule_new)
```

```
574   }
575
576   ##Load in Schedule
577   schedule <- scrape_men("20242025")
578
579   ## Load in my arbitrary initial elo ranking
580   X22Rankings <- read_csv(here("datasets_dataframes/22Rankings.csv"))
```

Function to do single day ratings:

```
582   ##Function to update rankings
583   ##rating is the variable, ratings is the df.
584   update_rankings <- function(season, game_date, ratings, k = 20){
585     ## Filters schedule to a specific date
586     elo_ratings_update <- season |> filter(date == game_date) |>
587       ## Joins the Elo ratings from our rating file to the schedule file. Puts
588   updated ratings in the schedule
589       left_join(ratings, by = join_by(away_team == Team)) |>
590       rename(away_elo = rating) |>
591       ## Updates ratings for home team in the schedule file
592       left_join(ratings, by = join_by(home_team == Team)) |>
593       rename(home_elo = rating) |>
594       ## Creating an away team outcome variable. Opposite of home team or same
595   if tie.
596       mutate(outcome_away = abs(outcome - 1)) |>
597       ## Calculating expected outcome variable for home and away team
598       mutate(exp_home = 1/(1 + 10^((away_elo - home_elo)/400))) |>
599       mutate(exp_away = 1/(1 + 10^((home_elo - away_elo)/400))) |>
600       ## Using expected outcome variable to generate new Elo ratings based on
601   actual outcome and expected outcome
602       mutate(elo_new_home = home_elo + k*(outcome - exp_home)) |>
603       mutate(elo_new_away = away_elo + k*(outcome_away - exp_away)) |>
604       rename(date_update_rankings = date)
605
606     ranked_home <- left_join(ratings, elo_ratings_update, by = join_by(Team ==
607   home_team)) |>
608       relocate(elo_new_home) |>
609       mutate(rating = if_else(!is.na(elo_new_home),
610                             true = elo_new_home,
611                             false = rating)) |>
612       select(Team, rating, date_update_rankings)
613
614     ratings_new <- left_join(ranked_home, elo_ratings_update, by = join_by(Team
615   == away_team)) |>
616       relocate(elo_new_away) |>
617       mutate(rating = if_else(!is.na(elo_new_away),
618                             true = elo_new_away,
619                             false = rating)) |>
620       select(Team, rating, date_update_rankings.x) |>
621       mutate(date_update_rankings.x = game_date) |>
```

```r
        rename(date = date_update_rankings.x)

   return(ratings_new)
}
```

Creating a vector of dates to be iterated to automate ratings:

```r
dates_vec <- unique(schedule$date)
## defining what new_rankings is going to be
new_rankings = rankings
## Creating a for loop with the update_rankings function to update rankings
up to a specified date
for (i in dates_vec) {
  new_rankings <- update_rankings(season = schedule, game_date = i, ratings =
new_rankings, k = 100)

}
```

Function to iterate through dates to automate ratings

```r
update_rankings_iter <- function(season, end_date, ratings, k){

  new_rankings <- list()

  season_cut <- season |>
    ## Filter by a specified end date to deal with NA values (gmaes that have
yet to be played)
    filter(date <= ymd(end_date))
  ## Creating a vector for unique dates in a season
  dates_vector <- unique(season_cut$date)
  ## Defining our rankings within the function
  new_rankings[[1]] <- ratings
  ## Creating a for loop for the function to generate new ratings with the
updtae_rankings function
  for (i in 1:length(dates_vector)) {
    new_rankings[[i + 1]] <- update_rankings(season = season_cut, game_date =
dates_vector[i], ratings = new_rankings[[i]], k = k)
  }
  return(new_rankings)
}
```

Function to update single day ratings with all three parameters, k, score differential, home ice advantage:

```r
update_rankings_gd_ha <- function(season, game_date, ratings, k = 100,
home_ice = 50, d = 0.5){
  ## Filters schedule to a specific date
  elo_ratings_update <- season |> filter(date == game_date) |>
    ## Joins the Elo ratings from our rating file to the schedule file. Puts
updated ratings in the schedule
```

```
666      left_join(ratings, by = join_by(away_team == Team)) |>
667      rename(away_elo = rating) |>
668      ## Updates ratings for home team in the schedule file
669      left_join(ratings, by = join_by(home_team == Team)) |>
670      rename(home_elo = rating) |>
671      ## Creating an away team outcome variable. Opposite of home team or same
672  if tie.
673      mutate(outcome_away = abs(outcome - 1)) |>
674      ## Calculating expected outcome variable for home and away team
675      mutate(exp_home = 1/(1 + 10^((away_elo - (home_elo + home_ice))/400))) |>
676      mutate(exp_away = 1/(1 + 10^(((home_elo + home_ice) - away_elo)/400))) |>
677      ## Using expected outcome variable to generate new Elo ratings based on
678  actual outcome and expected outcome
679      ## fivethirtyeights parameters
680      ## 0.6686 * log(abs(score_diff)) + 0.8048
681      mutate(score_mult = if_else(score_diff == 0,
682                                  true = 0.8048,
683                                  false = 0.6686 * log(abs(score_diff)) +
684  0.8048)) |>
685      mutate(elo_new_home = home_elo + k * score_mult * (outcome - exp_home))
686  |>
687      mutate(elo_new_away = away_elo + k * score_mult * (outcome_away -
688  exp_away)) |>
689      rename(date_update_rankings = date)
690
691    ranked_home_gd_ha <- left_join(ratings, elo_ratings_update, by =
692  join_by(Team == home_team)) |>
693      relocate(elo_new_home) |>
694      mutate(rating = if_else(!is.na(elo_new_home),
695                              true = elo_new_home,
696                              false = rating)) |>
697      select(Team, rating, date_update_rankings)
698
699    ratings_new_gd_ha <- left_join(ranked_home_gd_ha, elo_ratings_update, by =
700  join_by(Team == away_team)) |>
701      relocate(elo_new_away) |>
702      mutate(rating = if_else(!is.na(elo_new_away),
703                              true = elo_new_away,
704                              false = rating)) |>
705      select(Team, rating, date_update_rankings.x) |>
706      mutate(date_update_rankings.x = game_date) |>
707      rename(date = date_update_rankings.x)
708
709    return(ratings_new_gd_ha)
710  }
```

Function to iterate through entire season automatically:

```
711  rankings = X22Rankings
712
713
```

```r
##Function to update rankings
##rating is the variable, ratings is the df.
update_rankings_gd_ha <- function(season, game_date, ratings, k = 100,
home_ice = 50, d = 0.5){
  ## Filters schedule to a specific date
  elo_ratings_update <- season |> filter(date == game_date) |>
    ## Joins the Elo ratings from our rating file to the schedule file. Puts
updated ratings in the schedule
    left_join(ratings, by = join_by(away_team == Team)) |>
    rename(away_elo = rating) |>
    ## Updates ratings for home team in the schedule file
    left_join(ratings, by = join_by(home_team == Team)) |>
    rename(home_elo = rating) |>
    ## Creating an away team outcome variable. Opposite of home team or same
if tie.
    mutate(outcome_away = abs(outcome - 1)) |>
    ## Calculating expected outcome variable for home and away team
    mutate(exp_home = 1/(1 + 10^((away_elo - (home_elo + home_ice))/400))) |>
    mutate(exp_away = 1/(1 + 10^(((home_elo + home_ice) - away_elo)/400))) |>
    ## Using expected outcome variable to generate new Elo ratings based on
actual outcome and expected outcome
    ## 0.6686 * log(abs(score_diff)) + 0.8048
    mutate(score_mult = if_else(score_diff == 0,
                                true = 0.8048,
                                false = 0.6686 * log(abs(score_diff)) +
0.8048)) |>
    mutate(elo_new_home = home_elo + k * score_mult * (outcome - exp_home))
|>
    mutate(elo_new_away = away_elo + k * score_mult * (outcome_away -
exp_away)) |>
    ##mutate(elo_new_home = home_elo + k*(outcome - exp_home) + d *
score_diff) |>
    mutate(elo_new_home = if_else(elo_new_home < 100,
                                true = 100,
                                false = elo_new_home)) |>
    ##mutate(elo_new_away = away_elo + k*(outcome_away - exp_away) + d * -1 *
(score_diff)) |>
    mutate(elo_new_away = if_else(elo_new_away < 100,
                                true = 100,
                                false = elo_new_away)) |>
    rename(date_update_rankings = date)

  ## Find out which is the right date, select() and relocate(), **DO THIS
FIRST**: try renaming date in one of the df to make less confusing (at
start).
  ranked_home_gd_ha <- left_join(ratings, elo_ratings_update, by =
join_by(Team == home_team)) |>
    relocate(elo_new_home) |>
    mutate(rating = if_else(!is.na(elo_new_home),
                            true = elo_new_home,
```

```
764                              false = rating)) |>
765        select(Team, rating, date_update_rankings)
766
767    ratings_new_gd_ha <- left_join(ranked_home_gd_ha, elo_ratings_update, by =
768    join_by(Team == away_team)) |>
769        relocate(elo_new_away) |>
770        mutate(rating = if_else(!is.na(elo_new_away),
771                                true = elo_new_away,
772                                false = rating)) |>
773        select(Team, rating, date_update_rankings.x) |>
774        mutate(date_update_rankings.x = game_date) |>
775        rename(date = date_update_rankings.x)
776
777    return(ratings_new_gd_ha)
778    }
```

779    Creating a vector of dates to use again to automate:

```
780    dates_vec <- unique(schedule$date)
781    ## defining what new_rankings is going to be
782    new_rankings = rankings
783    ## Creating a for loop with the update_rankings function to update rankings
784    up to a specified date
785    for (i in dates_vec) {
786      new_rankings <- update_rankings_gd_ha(season = schedule, game_date = i,
787    ratings = new_rankings, k = 100, home_ice = 50, d = 0.5)
788
789    }
```

790    Function to return a completed set of rankings:

791    Function to run ratings and return the mean residual

```
792    update_rankings_residuals = function(season, end_date, ratings, k, home_ice,
793    d){
794
795      new_rankings = update_rankings_iter_gd_ha(season = season, end_date =
796    end_date, ratings = ratings, k = k, home_ice = home_ice, d = d)
797
798      full_rankings = new_rankings |> bind_rows()
799
800      lagged_dates = full_rankings |> group_by(date) |>
801        summarise(last_date = last(date)) |>
802        mutate(lag_date = lag(last_date)) |>
803        select(-last_date)
804
805      lagged_rankings <- left_join(full_rankings, lagged_dates, join_by(date ==
806    lag_date)) |>
807        select(-date) |>
808        rename(date = date.y)
```

```
809
810    season = season |>
811      mutate(home_elo = NA) |>
812      mutate(away_elo = NA)
813
814    merged_season_home = left_join(season, lagged_rankings,
815                                    by = join_by(date == date, home_team ==
816  Team)) |>
817      mutate(home_elo = rating) |>
818      select(-rating)
819
820    merged_season =  left_join(merged_season_home, lagged_rankings,
821                                by = join_by(date == date, away_team == Team))
822  |>
823      mutate(away_elo = rating) |>
824      select(-rating)
825
826    full_season = merged_season |>
827      mutate(outcome_away = abs(outcome - 1)) |>
828      ## Calculating expected outcome variable for home and away team
829      mutate(exp_home = 1/(1 + 10^((away_elo - (home_elo + home_ice))/400))) |>
830      mutate(exp_away = 1/(1 + 10^(((home_elo + home_ice) - away_elo)/400))) |>
831      ## Using expected outcome variable to generate new Elo ratings based on
832  actual outcome and expected outcome
833      ##
834      mutate(score_mult = if_else(score_diff == 0,
835                                  true = 0.8048,
836                                  false = 0.6686 * log(abs(score_diff)) +
837  0.8048)) |>
838      mutate(elo_new_home = home_elo + k * score_mult * (outcome - exp_home))
839  |>
840      mutate(elo_new_away = away_elo + k * score_mult * (outcome_away -
841  exp_away)) |>
842      mutate(elo_new_home = if_else(elo_new_home < 100,
843                                    true = 100,
844                                    false = elo_new_home)) |>
845      mutate(elo_new_away = if_else(elo_new_away < 100,
846                                    true = 100,
847                                    false = elo_new_away)) |>
848      mutate(residual = outcome - exp_home) |>
849      mutate(abs_residual = abs(residual))
850
851    mean_residual = full_season |> summarise(avg = mean(abs_residual, na.rm =
852  TRUE))
853
854    return(pull(mean_residual))
855  }
```

Additional data frames:

857 Optimization process:

858 Optimized ratings:

```
apr_15_ranking = update_rankings_iter_gd_ha(schedule, "2025-04-15",
                                            rankings2324, optimal$k,
optimal$home_ice,
                                            0)
```

863 Data frame for mapping plot:

```
plan(multisession)
handlers("progress")
options(progressr.enable = TRUE)

grid_100 = expand.grid(k = seq(60, 120, length.out = 20), home_ice = seq(30,
50, length.out = 20), d = seq(0, 100, length.out = 1))

mean_residuals_100 = with_progress({future_pmap_dbl(grid_100, \ (k, home_ice,
d) update_rankings_residuals(season = schedule_reg, end_date = "2025-03-25",
ratings = rankings2324, k = k, home_ice = home_ice, d = d), .progress =
TRUE)})

residual_100_df <- grid_100 |> mutate(mean_residual = mean_residuals_100)
```

## 9 Works Cited

878 College Hockey News. (2024). *Men's Division I Hockey Schedule* – 2023–2024 Season.
879 https://www.collegehockeynews.com/schedules/

880 College Hockey News. (2025). *Men's Division I Hockey Schedule* – 2024–2025 Season.
881 https://www.collegehockeynews.com/schedules/

882 **Silver, N.** (n.d.). *How our NHL predictions work*. FiveThirtyEight.
883 https://fivethirtyeight.com/methodology/how-our-nhl-predictions-work/

### 9.1