

# Machine Learning Models to Predict 2028 Olympic Track and Field Performances

Clare Robson  
Eberly College of Science  
The Pennsylvania State University  
Centre County, PA  
clr5887@psu.edu

**Abstract** – This paper’s aim is to develop and evaluate machine learning models to predict both winning (first place) and medaling (top three) threshold times for the women’s and men’s sprint and mid-distance Track and Field events for the 2028 Olympics. Sprint events are defined to be the 100m, 200m, and 400m races. Mid-distance events are defined to be the 800m and 1500m races. Using a dataset of historic top three Olympic results spanning back to 1896, individualized log-linear regression, k-nearest neighbors (KNN), and autoregressive integrated moving average (ARIMA) models were developed and evaluated via normalized root mean squared error (RMSE) and mean absolute error (MAE). The creation of an optimal model for these predictive purposes furthers the field of Track and Field sports prediction, and provides useable metrics for athletes and coaching staff to consider while training.

The KNN model presented with on average the lowest normalized RMSE and normalized MAE compared to the other considered models, at 0.0154199 and 0.01293649 respectively, demonstrating that on average the KNN model’s predictions were approximately 1.129% away from the actual historical values. 2028 Olympic predictions and respective 95% confidence intervals were generated and plotted using the KNN model.

**Keywords** – Olympic Games, Track and Field, performance prediction, machine learning, trend analysis, modeling

## I. INTRODUCTION

### A. Background

The modern Olympic Games are the world’s foremost sporting events, comprised of summer and winter competitions featuring thousands of athletes from across the globe [7]. This largely commercialized event is viewed by billions, with the 2024 Summer Olympics wracking up approximately 5 billion viewers [9]. Athletes and coaches face huge expectations from their home countries to earn medals, and nations invest heavily in training programs and sports analytics to optimize outcomes.

The Track and Field Athletics at the Olympic Games features a variety of running and field events. This work focuses on sprint and mid-distance events, for women and for men. Sprint events are defined to include the 100m, 200m, and 400m races, and mid-distance events are defined to include the 800m and 1500m races. Many of these races are among the most watched Olympic competitions, and are incredibly competitive, as medal placement can be determined by fractions of seconds.

These events were selected for this project in order to narrow the scope of the project, and to simplify finding the optimal machine learning model. The spread between the 100m sprint and the 1500m race is much narrower than between the 100m and 5000m, the next event in the lineup, and hurdle events, while also classified as sprint events, are different in strategy and final results time spread.

Each of these events have been around since either the first or second Olympics, providing 29 – 30 Olympics worth of data, whereas any other event at the longest has been a part of the Olympics for 26 games (besides the Marathon, which is technically apart of Track and Field at the Olympics). While this is not a huge difference, it would make a difference in the creation of each machine learning model. Therefore, events analyzed are narrowed to sprint and mid-distance events, excluding hurdles.

When training for these events, performance prediction is incredibly important for athletes, coaching staff, and support staff, to dually gain an insight as to what times can be anticipated from medalists (the top 3 finishers), and to develop goal-oriented training plans. Accurate predictive modeling can assist with resource allocation for athletes and teams, allowing those qualified for and competing in multiple events to focus their training on where they’re likely to perform the best. For national teams and sports scientists, modeling supports long term strategic development and to help forecast realistic goals.

Despite advances in sports analytics, existing research into Track and Field performance prediction is limited, and research into Track and Field Olympics performance prediction is even more limited. As it stands, current research tends to focus on single events, oftentimes just the marathon or the 100m sprint, and within that, only on men’s performances rather than women’s.

Additionally, much current research uses environmental conditions or an athlete’s attributes like season best, country of origin, and personal record for prediction, which, while incredibly useful, still leaves a gap for performance prediction accounting for long term historical patterns. This research aims to fill that gap by developing models based upon historical performance data, that can be applied to multiple events and both men and women’s performances. The generalizability of these model’s designs is

purposeful, so that the optimal one can be easily used as a tool moving forward for this range of events.

### *B. Literature Review*

In developing this work, much prior work in sports performance prediction was first examined to gain an understanding of useful methodologies. This research was helpful in the selection and evaluation of this work's machine learning models, because the pros and cons of each solution could be identified. Additionally, this prior work demonstrated where there was room for growth in the field of Track and Field performance prediction.

One paper examined is Szmygin et al. (2023) [1], wherein a multi-layer perceptron (MLP) model is proposed to predict men's 100m sprint times. Researchers analyzed race data from every historical man's 100m performance faster than 10.55 seconds. The model incorporates both race conditions, such as wind speed and location of the race, and athlete specific attributes, such as their birth country, age, and personal best record.

The approach demonstrates how deep learning can capture non-linear relationships in performance data. Additionally though, it also highlights the challenges that come from using detailed individual level features that are not consistently or easily available for all events or historical time periods. With the intent of focusing on women's races as well as men's, and on more events than the 100m, this specific research was designed to focus on aggregated performance data. These detailed level features were not always recorded for women, or for events longer than the 100m, therefore, generalizability is maintained by removing the MLP model from the list of potential models.

Efiong et al (2019) [3], is another paper that focuses on Track and Field performance prediction. Like Szmygin et al, it also focuses on the men's 100m sprint. Using past Olympic results, this paper proposes a MATLAB simulated model to predict likely Olympic gold medalist times for the men's 100m sprint for the next 5 Olympic competitions.

Unlike Szmygin's deep learning model, Eifong's approach is focused on long term trends, demonstrating accurate how predictive modeling can be with historical trend analysis. This finding influenced the inclusion of ARIMA models in this work, as ARIMA models can be used to analyze patterns and trend over time, much like the MATLAB simulated model was designed to do.

Crowley et al. (202) was one of the most important papers in the development of this work. While it doesn't focus on Track and Field, it does focus on predictive modeling in elite swimming, and is still useful for the development of Track and Field performance prediction. Swimming competitions are set up much the same as Track and Field, with events of different lengths, conducted in heats of 8, with the same events for

women and for men, and the paper's overall goal is to develop result predictions based upon historical data, the exact same goal of this work.

The model discussed in the paper predicted winning times for the Tokyo 2024 Olympics using historical trends and event level features, and shows how historical trend modeling can yield robust predictions even with limited input features. Crowley et al. breaks down their predictions by performance categories (1<sup>st</sup> – 3<sup>rd</sup> place, 4<sup>th</sup> – 8<sup>th</sup> place, and 9<sup>th</sup> – 16<sup>th</sup> place), and then uses linear regression and forecasting models to examine trends and predictions.

This paper's work doesn't focus on trends in the same way that Crowley et al. does, but does take influence from the predictions based on performance categories, and creates predictions for first place performances and third place performances. Additionally, it takes influence from how Crowley et al. generated and plotted their 95% confidence intervals.

### *C. Objective Statement*

Methods and insights from each of the above papers informed the modeling framework in this study. In particular, Crowley et al's approach of using previous times to predict performances by categories was useful to examine in the creation of a machine learning model to predict winning and medaling performances in Track and Field, by testing multiple modeling approaches. By using this cross-disciplinary approach, this project introduces a more generalizable framework for predicting Olympic medaling performances. Focus was also put on reproducibility.

Therefore, this paper aims to develop and evaluate several machine learning models to predict gold and bronze medal times for the Los Angeles 2028 Olympic Games, for the men's and women's 100m, 200m, 400m, 800m, and 1500m events, in order to select the most effective one. This model will then be able to generalize performance prediction across multiple events and genders, contributing to the advancement of predictive analytics in Track and Field.

## II. METHODOLOGY

### *A. Data Description and Preprocessing*

For this research, a dataset containing first, second, and third place results for all Track and Field Olympic events was used [6]. This dataset contains the variables gender, event, location, year, medal, athlete name, nationality, and result, and spans from 1896 to 2016. Women's results begin in 1928, when the women's 100m and 800m races were added, with the other women's races being added in later years.

Results from the Tokyo 2020 Olympics [8] and the Paris 2024 [10] were manually added. Due to the Covid-19 Pandemic,

the Tokyo games were postponed until 2021. Therefore, the year listed for those races is 2021, as opposed to 2020.

For the purposes of this work, unnecessary variables were removed from the dataset. Variables retained included Gender, Event, Year, Medal, and Result. Any results from a race resulting in a silver medal were removed from the dataset, as using bronze medal race results provides a medaling threshold. Predicting times needed to win a bronze medal will provide a medaling threshold, which upon hit is predicted to result in an Olympic medal. Events other than the listed sprint and mid-distance events were also removed from the dataset. The final large dataset included 458 observations.

Using this large dataset, smaller datasets for each of the sprint and mid-distance events were created. That way, events could be analyzed separately, as well as all together. From there, those smaller datasets were split even further into datasets containing either gold medal results or bronze medal results. This resulted in the creation of 20 individual datasets.

## B. Model Creation

Using an 80/20 training testing split, training and testing datasets were created from each of the individual datasets to be used in the training of four machine learning models.

### 1. Model 1: General Log-Linear Regression

This model uses the log of race times, but it includes year, gender (male), medal type (gold or bronze), and event as predictors. It's applied to all races together.

$$\log(\text{Result}) \sim \text{Year} + \text{Male} + \text{Medal} + \text{Event}$$

- Model 2: Individualized Log-Linear Regression

This model also uses the log of race times and year as the only predictor. It is applied separately to each race, so each gender and medal specific event is fitted independently from other events.

$$\log(\text{Result}) \sim \text{Year}$$

- Model 3: Individualized K-Nearest Neighbors (KNN),  $k = 3$

This model uses the log of race times and finds the closest neighbors to make predictions. Like the individualized linear regression model, it is applied separately to each race.

$$\log(\text{Result}) \sim \text{Year}$$

- Model 4: Individualized Autoregressive Integrated Moving Average (ARIMA)

This time series model uses race times directly and captures trends over time. It is applied separately to each race to predict future times, and optimal ARIMA parameters are automatically selected based upon Akaike Information Criteria (AIC).

$$\text{Result} \sim \text{ARIMA}(p, d, q)$$

$p$  = number of lagged observations included

$d$  = number of differences to make the series stationary

$q$  = number of past forecast errors included

Upon visually examining histograms of race results, it was determined that the majority were right-skewed and not normally distributed. Examples are shown in Figures A1 and A2 (other histogram examples can be found in Appendix B).

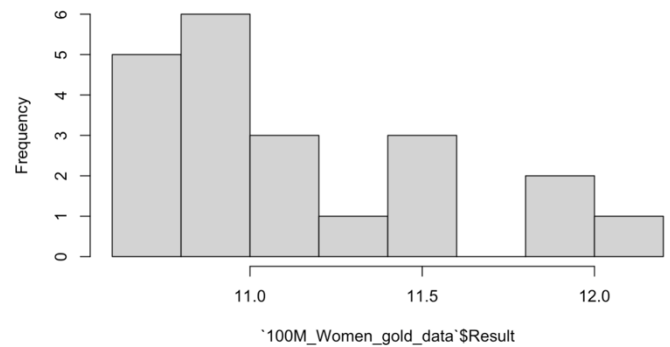


Fig. A1 Histogram of women's 100m gold medal Olympic finals race results

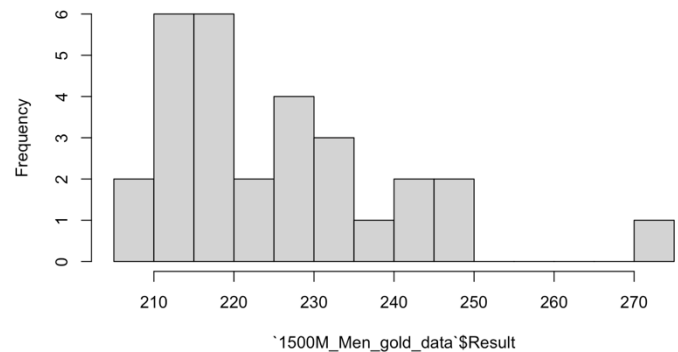


Fig. A2 Histogram of men's 1500m gold medal Olympic finals race results

Additionally, plots of historical gold medal and bronze medal times for each race display the improvement of times over the years at a decreasing rate. The relationship between result and year is non-linear. Plot examples are displayed in Figures A3, A4,

A5, and A6 (remaining figures for each race can be found in Appendix C).

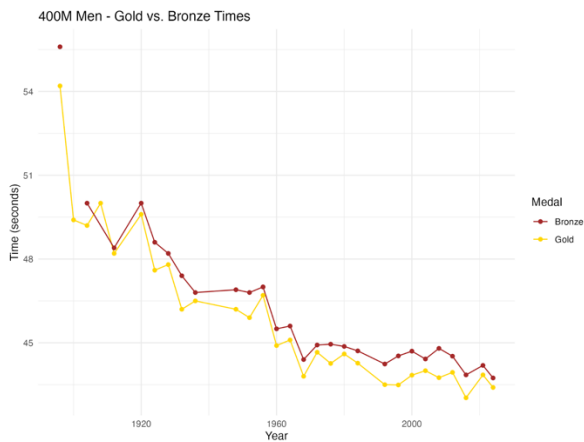


Fig. A3 400m men's Olympic historical trends

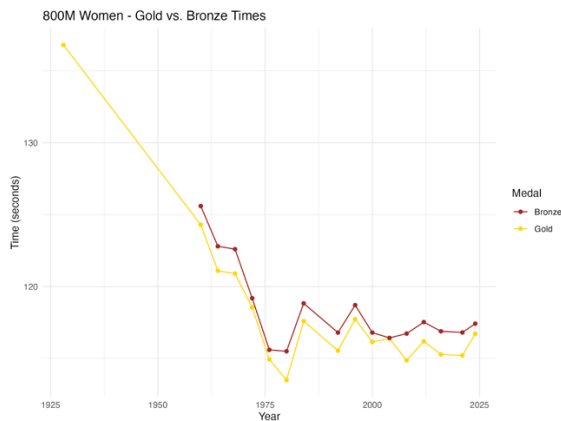


Fig. A4 800m women's Olympic historical trends

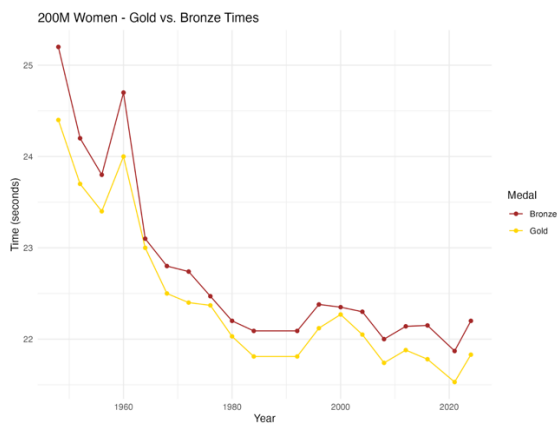


Fig. A5 200m women's Olympic historical trends

These factors limited the machine learning models that could be used, as well as how they could be implemented. Specifically, for the linear regression models and KNN model, a

log transformation was applied to the dependent variable time, so that the impact of extreme values could be reduced, and model performance could be improved. The log transformation ensures that percentage improvements in time are treated consistently, and linearizes displayed trends.

The ARIMA model is built to handle non-normal data, and therefore the log transformation was not applied.

### C. Model Evaluation

These models were evaluated via root mean squared error (RMSE), mean absolute error (MAE), normalized RMSE (NRMSE), and normalized MAE (NMAE). These normalized metrics are created by dividing RMSE and MAE by average race time. For methods 1 through 3, where the regression is based on the log transformation of the results, RMSE and MAE are calculated based on the original result values, rather than the log transformed ones. This way, the normalized values can be compared across all four methods.

This normalization is necessary not just for valid comparison across all four methods, but for comparison within each method from event to event. For example, a MAE of 1.00 seconds might be significant for the men's 100m race, where at an Olympic level this race is now run below 10 seconds, but not for the women's 1500m, which was most recently won in 231 seconds. Each value shows approximately what percentage, relative to average race time, each models' predictions are off by, with NRMSE being more sensitive to large prediction errors than normalized MAE.

Raw results for Models 2, 3, and 4 are displayed in Table A1, Table A2, and Table A3. Model 1's raw results ended up not being comparable, due to the resulting relatively large RMSE and MAE (3.00689 and 1.603443 respectively) and difficulty in normalizing each statistic for comparison. Therefore, these results are not displayed or compared, though the model summary can be found in Appendix E as Figure E1.

### D. Raw Results

Raw results for each method include RMSE, MAE, NRMSE, NMAE, 2028 Olympic time predictions, and values for the lower and upper confidence interval bounds for each of the events' models. In each Medal column, "B" refers to bronze and "G" refers to gold. Tables for models 2, 3, and 4 are displayed in Tables A1 – 3.

TABLE A1  
MODEL 2 RAW RESULTS (INDIVIDUALIZED LOG-LINEAR REGRESSION MODEL)

Race	Medal	LM_RMSE	LM_MAE	Normalized_LM_RMSE	Normalized_LM_MAE	LM_2025_Predictions	Lower_95_CI	Upper_95_CI
100M_Men	B	0.257167924852734	0.210026666079081	0.024647107981535	0.02100800333437	9.5818983397633	9.2672156523184	9.8565810272082
100M_Men	G	0.144756460116676	0.115014396259586	0.01417595884831	0.011230308644943	9.5402854283743	9.2656740491832	9.8269030364803
100M_Women	B	0.193514319558586	0.160804713677737	0.01515560949238	0.0144381336635454	10.5145560949238	10.187532055335	10.841620729314
100M_Women	G	0.24904028103519	0.21828457169658	0.022030124434697	0.0192864972759903	10.5592347966405	10.1049484534308	11.1034746388902
1500M_Men	B	4.4212027640097	3.86033779645916	0.020214835101374	0.0176503515120783	206.694348627407	198.300916184413	217.08710740242
1500M_Men	G	5.22094155653914	3.6826284180747	0.0238586177439202	0.016828820956187	206.602941617946	196.8456541328	216.341371822093
1500M_Women	B	5.46721570483724	5.23263837348013	0.0228107966483912	0.0218320725427874	242.618119245031	231.595315980245	253.279922508817
1500M_Women	G	4.89462166506772	4.5062318875836	0.020076106418018	0.0184830610219645	236.594236841733	232.007337181082	241.1811366502384
200M_Men	B	0.210726813650311	0.191439851015374	0.0103216503551289	0.0093769519502044	19.3939115651676	18.96303929018	19.8247839013172
200M_Men	G	0.378423025619622	0.319723350216517	0.0184281970109385	0.015596786080602	19.1337272939573	18.3000884960994	19.9257560910353
200M_Women	B	0.406423878313809	0.389407395070521	0.0180934392126348	0.0173358884839408	21.3448518954495	20.4613760355579	22.2283277561412
200M_Women	G	0.55472507890003	0.41051006244335	0.0244911734616169	0.0181240660503525	21.335453391164	20.1106879503887	22.560218767844
400M_Men	B	0.765039326476235	0.6663097802381411	0.017022971593862	0.0149826807942393	42.582577028888	40.920413566912	44.2447404908545
400M_Men	G	0.563439673275903	0.47730200903267	0.012602230290026	0.010341209033086	42.2307160056962	41.048808226813	43.4126271482554
400M_Women	B	1.08461875348858	0.941772666233535	0.0216626611911677	0.0188153836219406	48.745321341077	47.057120065469	50.511852204658
400M_Women	G	0.528913879979373	0.42726399903489	0.0108458075182596	0.0087613947538225	48.2402631353463	47.491879152356	48.9886471183366
800M_Men	B	2.9959721391005	2.8399962194232	0.0277524754794525	0.023676287852373	99.055518181616	94.1890909058937	103.92122485473
800M_Men	G	2.97378388446059	2.7109696322828	0.0282435547959027	0.02574782923951	99.085402838978	92.730756893387	105.4004669236
800M_Women	B	1.18392302400218	1.1306356790869	0.01008108995627	0.00962750067344767	115.06447845894	112.22955412364	117.899402796217
800M_Women	G	2.2510190393789	2.11324995906544	0.019342390451634	0.0181585784113376	110.886450422546	109.131521472064	114.641369418028

TABLE A2  
MODEL 3 RAW RESULTS (INDIVIDUALIZED KNN MODEL, K = 3)

Race	Medal	KNN_RMSE	KNN_MAE	Normalized_KNN_RMSE	Normalized_KNN_MAE	KNN_2025_Predictions	Lower_95_CI	Upper_95_CI
100M_Men	B	0.44580848858361	0.332469711991185	0.042726518524063	0.0318640705377789	9.8699052720055	9.06341959189828	10.678309524328
100M_Men	G	0.20545325733249	0.161897597252625	0.0201194910135383	0.0158541796271217	9.79999658653828	9.49084771376397	10.109145835126
100M_Women	B	0.110636345341451	0.090308761443084	0.00835789946668314	0.005204455473234	10.85024045547324	10.787677576741	11.1164115337906
100M_Women	G	0.23791485241841	0.206236368145297	0.0102200299595853	0.0182219798679358	10.6798842726171	10.169459559131	11.1902729539242
1500M_Men	B	1.9571238624085	1.7388088063034	0.008948402015549017	0.007905078211936725	215.456117963496	212.267707271272	218.644528655719
1500M_Men	G	4.47307948738427	4.20707913269571	0.0218119810843398	0.019225477186664	18.928643941599	208.89911588753	228.9581711356468
1500M_Women	B	5.69337001899834	5.54387867244651	0.023735065673008	0.0241529119229064	229.83354397033	253.224884801495	250.1002854509
1500M_Women	G	4.562560784739606	5.13780989032414	0.022524570340397	0.0201735795938016	237.643451255749	233.457940140708	241.82896210079
200M_Men	B	0.164710644196389	0.137829874006633	0.00607510714180172	0.005108815386272224	19.93025568272224	18.6701047998858	20.1902854509
200M_Men	G	0.255874928140973	0.196247971424366	0.014604296070254	0.00950576336519804	19.466280461534	18.9310291403671	20.00125746839
200M_Women	B	0.37138935865556	0.290368024854269	0.0120268169105963	0.0105208169105963	22.1165031606561	21.365107629133	22.869835718339
200M_Women	G	0.370798124108981	0.277117169918987	0.0168378251262243	0.0123427707065381	21.7129357831871	21.000483417742	22.423875306001
400M_Men	B	0.25315511012148	0.196749413979556	0.00404211886084629	0.00404211886084629	44.3510578312523	43.854562980426	44.867569644081
400M_Men	G	0.528645950619055	0.541387000169811	0.0117287045097446	0.0117287045097446	43.7293603515	42.2396237158806	45.2191973544194
400M_Women	B	0.7906519478508	0.524890811307212	0.010689961900051	0.010488428569004	49.42871859414189	49.1806751652898	49.676782002308
400M_Women	G	0.41112703613861	0.295734789192441	0.008430323883515	0.00606428157050827	49.0493122671432	48.144120426723	49.954039821431
800M_Men	B	1.07303013633557	0.99931888145097	0.0094232677410214	0.00925695252378469	103.260876357972	100.97123009641	105.50522760302
800M_Women	B	0.7888718272713	1.65363162862639	0.017939701495881	0.0151054965750802	101.451792920292	97.476133099801	105.35442548405
800M_Men	G	0.889746095077839	0.79322282172702	0.0075617584777064	0.00675427982585747	117.279666646518	115.1615502769848	119.397782092188
800M_Women	G	0.8904554998214	0.83324788461841	0.00760512321933547	0.0071596751789536	115.620618750214	114.166613850382	117.05462860047

TABLE A3  
MODEL 4 RAW RESULTS (INDIVIDUALIZED ARIMA MODEL)

Race	Medal	ARIMA_RMSE	ARIMA_MAE	Normalized_ARIMA_RMSE	Normalized_ARIMA_MAE	ARIMA_2025_Predictions	ARIMA_95CI_Lower	ARIMA_95CI_Upper
100M_Men	B	0.527603240277224	0.507471666666778	0.050565789127932	0.0486370778864077	10.5070833333333	9.0404029534084	11.9737573713184
100M_Men	G	0.239818837843808	0.19138711594251	0.0234847808050442	0.018740967289683	10.3091304347628	9.1723265094255	11.440962100491
100M_Women	B	0.484020031222904	0.47406250000028	0.043458866618981	0.042564532312034	11.3231200000001	10.479331700967	12.1669183290104
100M_Women	G	0.53602625477374	0.383875	0.0474046806511986	0.03971715214703	11.100975	10.3038808695088	11.9113931304612
1500M_Men	B	9.62790303890388	8.99623061383817	0.0440211109249038	0.0411419781622905	226.919828651377	196.624874860707	257.214782424047
1500M_Men	G	14.1557432368821	13.548550724638	0.046488809180949	0.0605719306677553	227.224782606096	198.10480228428	256.344762933111
1500M_Women	B	5.58253729648887	5.34800000000001	0.0232919515033595	0.0223133944341683	242.894	232.253743207615	253.534256792885
1500M_Women	G	5.55451179383892	5.13888888888889	0.022782757350201	0.0210780081421595	238.064444444444	227.986934132636	249.3491495756453
200M_Men	B	0.717181482992643	0.624857142857132	0.0351284029992782	0.030906247102742	20.8042857142857	19.0115713435132	22.597000805881
200M_Men	G	0.91517578396674	0.848833333333332	0.045687189171226	0.0413115818521223	20.55	18.641857889305	22.458145106365
200M_Women	B	0.62746851862036	0.571849999999949	0.0279518509999942	0.0230929198715042	22.8181959999999	20.7797917395489	24.842062040491
200M_Women	G	1.0846836773759	0.977995054535385	0.047888648241487	0.042954914787675	22.514480255987	21.4348085999999	23.5941985811953
400M_Men	B	2.73759027605999	2.23640269596554	0.051101739384338	0.04570780683871	46.707806985665	41.488744346296	51.907077396271
400M_Men	G	3.6294668593073	2.30660666666667	0.056961875139607	0.049871149624495	45.8656217391305	40.659101236059	51.0739424221
400M_Women	B	1.6166164630051	1.12333333333333	0.02290680504521	0.022442727573591	50.25	47.7619613509693	52.738306407537
400M_Women	G	1.992790679515658	0.943333333333343	0.020357977023561	0.0193483148080543	49.7099999999999	47.1542053244004	52.657940755195
800M_Men	B	4.74331234667723	4.57968253668255	0.043938545553678	0.0424227987990108	100.399047619048	90.144822121668	126.65373205958
800M_Men	G	4.7014874625333	4.34283080060431	0.0446848777020628	0.040296689721181	111.499667575102	96.1577508187458	124.841584331458
800M_Women	B	5.23756304655768	5.15999999999999	0.04459777984939	0.0438373297002724	122.8	116.902397275788	128.29760274212
800M_Women	G	2.717553896147	2.5648078923112	0.0235102947736	0.0220388040245	118.492307692511	106.833935422523	131.050678962369

### III. RESULTS

#### A. Metric Comparisons

Mean NRMSE, minimum NRMSE, maximum NRMSE, NRMSE range, mean NMAE, minimum NMAE, maximum NMAE, and NMAE range are displayed in Table A4.

TABLE A4  
SUMMARY OF NORMALIZED RMSE AND MAE METRICS ACROSS MODELS 2, 3 AND 4

	Mean Normalized RMSE	Min. Normalized RMSE	Max. Normalized RMSE	Range Normalized RMSE	Mean Normalized MAE	Min. Normalized MAE	Max. Normalized MAE	Range Normalized MAE
Model 2	0.0192	0.01008	0.02824	0.01816	0.0167	0.00876	0.02631	0.01755
Model 3	0.01542	0.00569	0.04273	0.03703	0.01294	0.00442	0.03186	0.02744
Model 4	0.03916	0.02036	0.06469	0.04433	0.03585	0.01874	0.06057	0.04183

Model 3 has both the lowest mean NRMSE and mean NMAE values, at 0.01542 and 0.01294 respectively. Model 2 has the smallest ranges of values, with a 0.01816 range for NRMSE, and a 0.01755 range for NMAE.

Mean values can also be compared for each race, between each model. These results are displayed in Table A5, with the highlighted values being the lowest per race between each model.

TABLE A5  
NORMALIZED RMSE AND MAE COMPARISON FOR EACH RACE

Medal	Race	Normalized_LM_RMSE	Normalized_KNN_RMSE	Normalized_ARIMA_RMSE	Normalized_LM_MAE	Normalized_KNN_MAE	Normalized_ARIMA_MAE	Lowest_RMSE	Lowest_MAE
B	100M Men	0.024647179981535	0.0427265189340483	0.050565789127932	0.02100800333437	0.0318640705377789	0.0486370778864077	LM	LM
G	100M Men	0.01417595884831	0.0201194910135383	0.0234847808050442	0.011230308644943	0.0158541796271217	0.0187409672896863	LM	LM
B	100M Women	0.0175545609503082	0.010005997776275	0.043458866618981	0.0144381336635454	0.0083578994668314	0.042564532312034	KNN	KNN
G	100M Women	0.0202039124434697	0.02102009595853	0.047404860571986	0.0192864972759903	0.0182219798679358	0.039172115210743	KNN	KNN
B	1500M Men	0.0202148351013714	0.00894840201554907	0.040211102049038	0.0073852175903	0.0141918792905	0.0411918792905	KNN	KNN
G	1500M Men	0.023686174738202	0.02181181654336	0.04689018801949	0.0162830527681	0.0182547718664	0.0251933267535	LM	LM
B	1500M Women	0.020219796643912	0.02375437740407	0.0234951915035495	0.02175427254784	0.0233606673088	0.022133964341683	LM	LM
G	1500M Women	0.02007816810418	0.022254034704397	0.027827557355201	0.0184306130144	0.02130736959304	0.021078001421595	LM	LM
B	200M Men	0.010216603501589	0.00687723586172	0.05152402699872	0.005124951950544	0.00675101449672	0.0306604271595	LM	LM
G	200M Men	0.016284910193045	0.00870690790354	0.045486911871226	0.00595678606002	0.00595678606002	0.01515818521237	LM	LM
B	200M Women	0.019105641959942	0.01195959995942	0.01735588480351	0.01195959995942	0.01195959995942	0.01195959995942	KNN	KNN
G	200M Women	0.01941737616349	0.016378751292045	0.01814969999999	0.01294777000000	0.01294777000000	0.042564532312034	KNN	KNN
B	400M Men	0.01702779918930	0.005696496023897	0.010712208484	0.01496806290000	0.0045421884466	0.050277491606287	KNN	KNN
G	400M Men	0.01292069039326	0.01578524203336	0.05696481757967	0.0103113003618	0.011287245050789	0.0191714660249	LM	LM
B	400M Women	0.021660251911677	0.016989619100351	0.01881538219461	0.015488438290000	0.0142277737950	0.024277737950	KNN	KNN
G	400M Women	0.019486818758596	0.020430432838316	0.020357877023561	0.023347343793255	0.009064215760567	0.019538314080565	KNN	KNN
B	800M Men	0.023752474945262	0.009942326741214	0.043835645563876	0.026307828782373	0.0201569525278489	0.042277737950	KNN	KNN
G	800M Men	0.02834534769027	0.021793802194681	0.04648857720026	0.0251783002194681	0.015704550178300	0.042277737950	KNN	KNN
B	1500M Women	0.01910809350502	0.045871777767504	0.045871777767504	0.006754279000000	0.006754279000000	0.04397327002724	KNN	KNN
G	1500M Women	0.019342934025404	0.007665102193547	0.020351002917336	0.018158734113376	0.01159275178637	0.022086060421595	LM	LM



the KNN model's ability to show local patterns for each race, so it has greater flexibility in examining trends for individual races.

In terms of stability, Model 2 tends to be more stable, with less variation in its performance across different events. This is reflected in its smaller range of NRMSE and NMAE values. Model 3, while more accurate, shows slightly higher variability, which suggests it may be more sensitive to the specific data points in each event. Despite this, Model 3's overall accuracy makes it a more reliable choice for prediction, though Model 2's stability provides an advantage in scenarios where consistency across different events is more important.

Using NMAE, as MAE is less prone to influence from outliers than RMSE, Model 3's predictions on average reflect a 1.294% error rate.

### C. 2028 Olympic Predictions

Time predictions along with 95% confidence intervals generated using Model 3, the KNN model, are shown below in Table A6.

TABLE A6

2028 OLYMPIC TIME PREDICTIONS AND 95% CONFIDENCE INTERVALS

Race	Medal	Predictions (secs)	95% Confidence Interval
100M Men	B	9.8699	[9.0634, 10.6764]
	G	9.8000	[9.4908, 10.1091]
100M Women	B	10.9520	[10.7877, 11.1164]
	G	10.6799	[10.1695, 11.1903]
1500M Men	B	215.4561	[212.2677, 218.6445]
	G	218.9286	[208.8991, 228.9582]
1500M Women	B	241.5291	[229.8334, 253.2249]
	G	237.6435	[233.4579, 241.8290]
200M Men	B	19.9326	[19.6701, 20.1950]
	G	19.4663	[18.9310, 20.0015]
200M Women	B	22.1165	[21.3632, 22.8699]
	G	21.7129	[21.0005, 22.4254]
400M Men	B	44.3511	[43.8345, 44.8676]
	G	43.7294	[42.2395, 45.2192]
400M Women	B	49.4287	[49.1807, 49.6768]
	G	49.0493	[48.1441, 49.9545]
800M Men	B	103.2609	[100.9712, 105.5505]
	G	101.4153	[97.4761, 105.3544]
800M Women	B	117.2797	[115.1616, 119.3978]
	G	115.6206	[114.1866, 117.0546]

Confidence intervals were calculated using the following formula:

$$\hat{y}_{2028} \pm 1.96 \times \hat{\sigma}_{\text{residual}}$$

Individual plots were generated for each race, gender, and medal type, displaying historical results, 2028 predictions, and 2028 95% confidence intervals in the same line graph. One example is shown in Figure A6 (the remaining plots can be found in Appendix D).

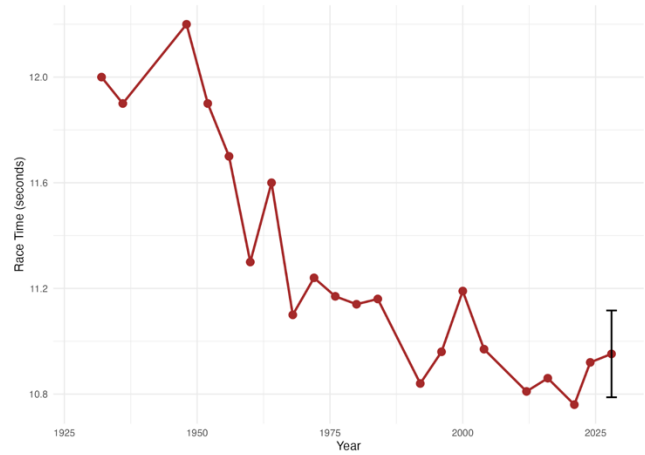


Fig. A6 Women's 100m bronze medalist race times vs year, including 2028 predictions and 95% CI's

### III. RESULTS

Even though Model 3, the KNN model, demonstrated strong predicted performance based on NRMSE and NMAE, there are still several limitations to consider. One limitation is the distribution of residuals. Even though  $\log(\text{time})$  was used in the creation of the model, errors still skewed positive as shown by the QQ plots in Figure A7.

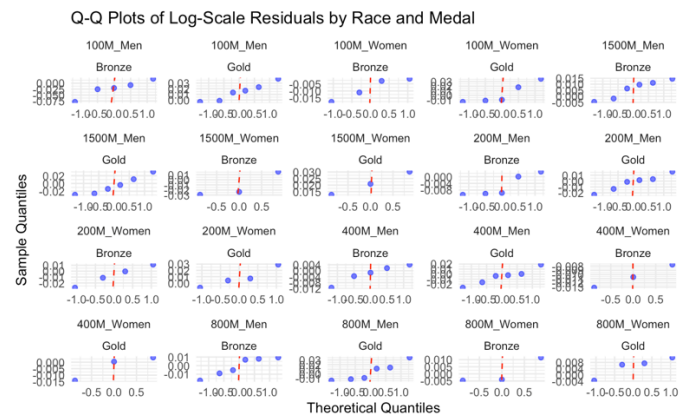


Fig. A7 QQ plots of log-scale residuals by race and medal

The positive skew suggests that the model tends to underpredict fast race times and overpredict slow race times. Therefore, the 95% confidence intervals generated and reported in Table A5 may not be reliable. To improve error distribution, a different transformation and/or additional transformations may be necessary in the future.

Another limitation of the KNN model is its reliance on local similarity, which can make it sensitive to variations in training data. Its predictions are influenced by nearby historical observations, which could lead to instability when predicting future results, like in the case of predicting times for the 2028 Olympics. Additionally, the choice of  $k = 3$  introduces bias-variance tradeoff. While  $k = 3$  was determined to be the optimal choice for this model, adjusting model parameters and transformations in future work may mean a different  $k$  value should be used to create a model that isn't too sensitive to noise, but also isn't reducing responsiveness.

#### IV. CONCLUSIONS

This paper describes the development of 4 machine learning models, designed to predict winning and medaling times for the sprint and mid-distance events at the 2028 LA Olympics. Using historical Olympic results, log-linear regression models, a KNN model, and an ARIMA model were built. After assessing each model via NRMSE, NMAE, and range of NRMSE and NMAE values, Model 3, the KNN model, was assessed to be optimal.

While this work is valuable, there are limitations. Model 3 relies solely on historical results, and does not include any outside factors such as environment or an athlete's training. Future work could expand upon these models by incorporating outside factors, and/or applying additional data transformations to limit positively skewed errors.

Taking into consideration the limitations of this work, future work could also include customizing the chosen  $k$  value per each race category. Building this into the model would be a complex solution to help improve accuracy with more regard to the differences for each smaller data set.

Overall methodology provides a framework by which to predict Olympic race performances, for multiple events and for both men and women. This framework can aid athletes, coaches, and analysts in understanding expected medaling and winning thresholds for the 2028 Olympic Games.

#### REFERENCES

- [1] A. Szmygin, M. Wojtowicz, Ż. Świdarska-Chadaj and R. Roszczyk, "Prediction of athletes' performance results using machine learning algorithms," 2023 24th International Conference on Computational Problems of Electrical Engineering (CPEE), Grybów, Poland, 2023, pp. 1-5, doi: 10.1109/CPEE59623.2023.10285142.
- [2] Crowley, Emmet, Kwok Ng, Iñigo Mujika, and Cormac Powell. "Speeding up or Slowing Down? Analysis of Race Results in Elite-

Level Swimming from 2011-2019 to Predict Future Olympic Games Performances." *Measurement in Physical Education and Exercise Science* 26,no. 2 (2022): 130–40. doi:10.1080/1091367X.2021.1952592.

- [3] Efiog, John & Olajubu, Emmanuel & Dr.Aranuwa, Felix. (2019). Formulation of Sprint Time Predictive Model for Olympic Athletic Games. *International Journal of Information Technology and Computer Science*. 4. 33-43. 10.5815/ijitcs.2019.04.04.
- [4] Wu, P. P. Y., Garufi, L., Drovandi, C., Mengersen, K., Mitchell, L. J. G., Osborne, M. A., & Pyne, D. B. (2022). Bayesian prediction of winning times for elite swimming events. *Journal of Sports Sciences*, 40(1), 24–31. <https://doi.org/10.1080/02640414.2021.1976485>
- [5] Mujika, I., Pyne, D. B., Wu, P. P., Ng, K., Crowley, E., & Powell, C. (2023). Next-Generation Models for Predicting Winning Times in Elite Swimming Events: Updated Predictions for the Paris 2024 Olympic Games. *International journal of sports physiology and performance*, 18(11), 1269–1274. <https://doi.org/10.1123/ijspp.2023-0174>
- [6] Ravaliya, Jay (2017). Olympic Track and Field Results. Retrieved from <https://www.kaggle.com/datasets/jayrav13/olympic-track-field-results?resource=download>
- [7] NBC Olympics (2024). Track and field 101: Olympic history, records and results. Retrieved from <https://www.nbcolympics.com/news/track-and-field-101-olympic-history-records-and-results>
- [8] Olympics (2024). Tokyo 2020 Athletics Results. Retrieved from <https://olympics.com/en/olympic-games/tokyo-2020/results/athletics>
- [9] Olympics (2024). Around 5 billion people - 84 percent of the potential global audience - followed the Olympic Games Paris 2024. Retrieved from <https://olympics.com/ioc/news/around-5-billion-people-84-percent-of-the-potential-global-audience-followed-the-olympic-games-paris-2024>
- [10] Runner's World (2024). Results from Track and Field at the 2024 Paris Olympics. Retrieved from <https://www.runnersworld.com/races-places/a61776840/2024-paris-olympics-results-track-and-field/>

#### APPENDIX A ASSETS FROM MAIN BODY

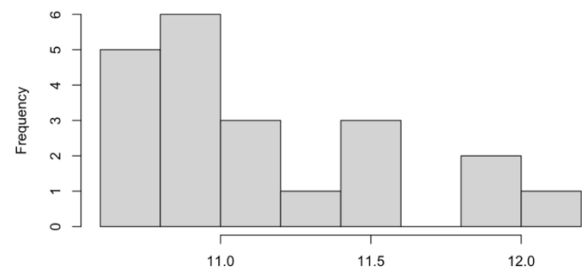


Fig. A1 Histogram of women's 100m gold medal Olympic finals race results

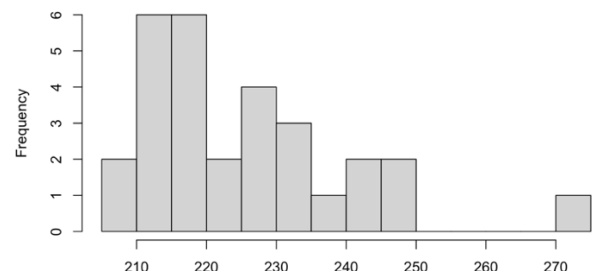


Fig. A2 Histogram of men's 1500m gold medal Olympic finals race results

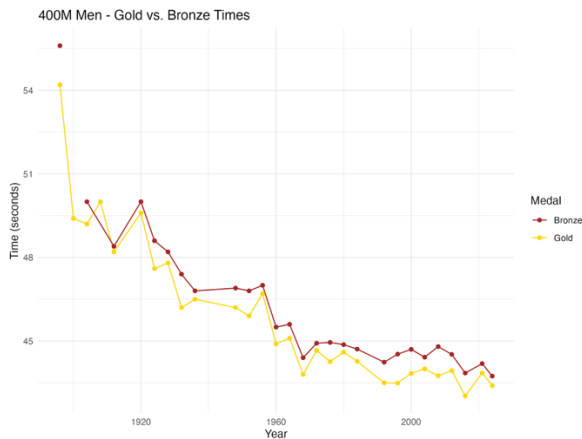


Fig. A3 400m men's Olympic historical trends

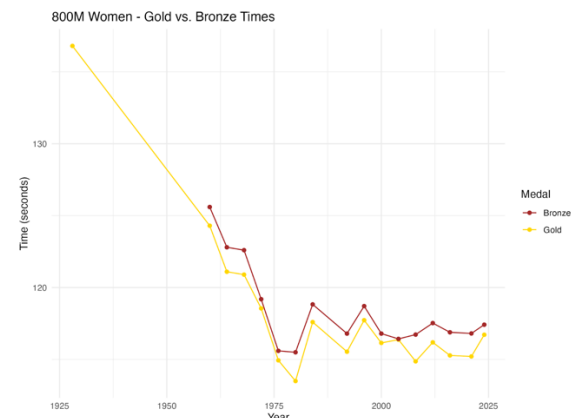


Fig. A4 800m women's Olympic historical trends

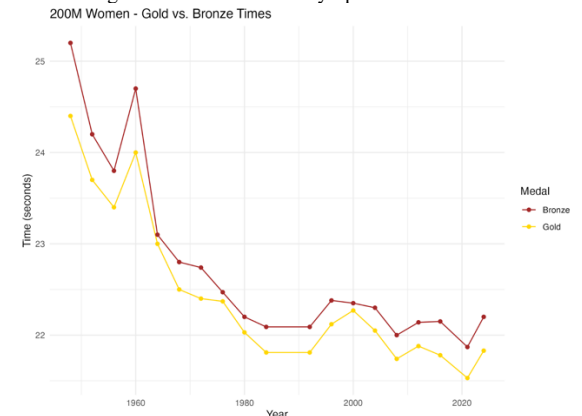


Fig. A5 200m women's Olympic historical trends

TABLE A1  
MODEL 2 RAW RESULTS (INDIVIDUALIZED LOG-LINEAR REGRESSION MODEL)

Race	Medal	LM_RMSE	LM_MAE	Normalized LM_RMSE	Normalized LM_MAE	LM_2020_Predictions	Lower_95_CI	Upper_95_CI
100M_Men	B	0.25167924852734	0.218208866079081	0.0246471079981535	0.021089003333437	9.5618963397633	9.2672156523184	9.856810272082
100M_Men	G	0.144759640116676	0.115014396295996	0.01475959594581	0.011263266449443	9.5462854263743	9.2657404974832	9.8290320349853
100M_Women	B	0.1951431955659	0.16864713877737	0.0179149586532062	0.014181133662544	10.514556946228	10.1875132665335	10.84160729214
100M_Women	G	0.249494028103519	0.218208866079081	0.022023914344687	0.018266497759903	10.5592347966405	10.148946523663	11.1034745898902
1500M_Men	B	4.4212027646097	3.86033779645916	0.02027148381013174	0.0176503515120783	206.894386274027	198.300916184413	215.387781079402
1500M_Men	G	5.2209415563914	3.68262284180747	0.0238646177439202	0.016828620656187	206.802941617946	196.8646641376	216.347317623003
1500M_Women	B	5.46721570493724	5.23263637348013	0.0228707964901812	0.0218320725427874	242.618119240331	231.950315980245	253.279922509517
1500M_Women	G	4.89462166506772	4.5062318735836	0.0200761064180186	0.0184630610219645	236.594236847733	232.207337181082	241.18113662304
200M_Men	B	0.0217026813650311	0.01439851015374	0.0103216505351289	0.0093789519502044	19.3839115651676	18.363028229018	20.4027789013172
200M_Men	G	0.037482025619622	0.019123350216517	0.0184281970103065	0.0155696780605002	18.1133722935673	18.300984960994	19.9257609103053
200M_Women	B	0.404243878313809	0.389407395070521	0.0180934382126348	0.017338864834948	21.3448518958495	20.461376035579	22.283277561412
200M_Women	G	0.554752707806003	0.410510060244335	0.0244911734616349	0.0181240660593525	21.3354533591164	20.1106879503887	22.560218767844
400M_Men	B	0.765020326476235	0.66630078281411	0.01702027191598362	0.01498268078942393	42.582577020888	40.800413669218	44.2447049008545
400M_Men	G	0.563439673275903	0.477350209093267	0.012206232003026	0.010341200335816	42.2307160056692	41.04800226813	43.4126227845254
400M_Women	B	1.08461875348858	0.941727668223535	0.0218682611911677	0.018813383619406	48.7845321341077	47.0572120665469	50.5118522016885
400M_Women	G	0.538913879873703	0.42726399038499	0.010845807182596	0.0087613943736225	48.740263135463	47.491979152356	48.9886471183366
800M_Men	B	2.9959721391005	2.8399621946232	0.027752474579425	0.0263078287852373	96.055181816166	94.188989068937	103.92122458734
800M_Men	G	2.9376388466059	2.7108968322828	0.0284435547596027	0.0257478362923951	96.085402838078	92.7307936883537	105.4004689236
800M_Women	B	1.1839232003718	1.106536790869	0.0100810899550227	0.00962750667344767	115.06447845884	112.229554123464	117.896402796217
800M_Women	G	2.2510190393769	2.113249695905544	0.019342230451634	0.0181585784113376	110.886450422546	109.131531427064	112.641369418028

TABLE A2  
MODEL 3 RAW RESULTS (INDIVIDUALIZED KNN MODEL, K = 3)

Race	Medal	KNN_RMSE	KNN_MAE	Normalized KNN_RMSE	Normalized KNN_MAE	KNN_2020_Predictions	Lower_95_CI	Upper_95_CI
100M_Men	B	0.44380848858361	0.332489711991185	0.042726189340963	0.031864075377789	9.8699057220055	9.0341059196828	10.674380542328
100M_Men	G	0.20545353733249	0.161897587625625	0.0201194910135383	0.0158541796271217	9.7999905865328	9.4908471736973	10.1460061061895
100M_Women	B	0.118063845341451	0.09030838781443884	0.0106005662776275	0.0083576994066314	10.3620445447334	10.78767756741	11.1164113337208
100M_Women	G	0.23791488521841	0.206236368144297	0.0210220695508653	0.018221798679358	10.6798842726171	10.18494959131	10.92122458734
1500M_Men	B	4.73701238624085	1.7389288063034	0.0084842015548017	0.0079507821193675	215.456117963496	212.26770271272	218.64452805719
1500M_Men	G	4.73701238624085	1.7389288063034	0.0218119810835971	0.01922547186664	208.8991158675	208.8991158675	208.8991158675
1500M_Women	B	5.8933701899634	5.54387867824481	0.022754377504487	0.02313065673308	241.529119229264	239.833545397023	253.2248840611895
1500M_Women	G	5.45540784729686	5.13782885922414	0.022245347424987	0.021073769938616	237.8434015527149	235.45784541078	241.85995210079
200M_Men	B	0.164707464119639	0.13782885922414	0.002607723597716	0.00571071741400172	19.3025695627224	19.3025695627224	19.3025695627224
200M_Men	G	0.255874629140973	0.196247874104356	0.013460429070354	0.0095657556519024	19.466286546513	19.301091453671	20.001592746939
200M_Women	B	0.3713851863556	0.290368624854269	0.0216535642486616	0.0129826169105963	22.1165201659661	21.263152796278	22.869537818339
200M_Women	G	0.370387182108881	0.277117899198087	0.016378521262243	0.012234777006381	21.7129357381871	21.0004839417742	22.4253873300001
400M_Men	B	0.25315511012148	0.196748413975556	0.00560246605212897	0.004421188084629	44.3510578312253	43.834540980426	44.867583644081
400M_Men	G	0.278645850610955	0.5413970016981	0.01578523233336	0.0117287045097446	43.72906353515	42.2395237158806	45.2191975544194
400M_Women	B	0.530553194766558	0.524880811307212	0.010889615900351	0.010488428569004	49.4287185544189	49.180675168298	49.677602302308
400M_Women	G	0.41112735013861	0.295734798192441	0.0084304923883515	0.0060428157605827	49.3493122671432	48.144126482733	49.8543038914131
800M_Men	B	1.07307316335357	0.99931888145097	0.0094523267141014	0.00925695252378469	103.2608763257972	100.97123006411	105.5502270302
800M_Men	G	1.8887121627213	1.65363162862639	0.017939701945681	0.0157054955705802	101.512792329296	97.476133096081	105.35425484805
800M_Women	B	0.88974695077739	0.793222862172702	0.0075761758777064	0.0075427892585747	117.27966846518	115.16150570948	119.400778292318
800M_Women	G	0.89204554998214	0.833247288461841	0.0076651032193547	0.0071596751798536	115.62081795214	114.186613863382	117.05482680047

TABLE A3  
MODEL 4 RAW RESULTS (INDIVIDUALIZED ARIMA MODEL)

Race	Medal	ARIMA_RMSE	ARIMA_MAE	Normalized ARIMA_RMSE	Normalized ARIMA_MAE	ARIMA_2020_Predictions	ARIMA_95CI_Lower	ARIMA_95CI_Upper
100M_Men	B	0.52703025027224	0.50741666666778	0.0505657691273932	0.0488310778864077	10.5078333333339	9.8404026934394	11.373737373184
100M_Men	G	0.238818837842808	0.191378811584251	0.023848780802442	0.018740907286883	10.3091304347628	9.17235028942155	11.44580210144
100M_Women	B	0.484820031222904	0.47406250000028	0.043458868619891	0.040564543212024	11.3231200000001	10.4793318739887	12.166918320014
100M_Women	G	0.53626252477374	0.383875	0.047404850271986	0.033917215214703	11.103675	10.330880665588	11.9148691304612
1500M_Men	B	8.9720393808388	8.9823061383817	0.044221110248036	0.04114197828205	226.919823651377	186.62487486077	257.21478242407
1500M_Men	G	1.1557423326261	13.254650774638	0.06488801980949	0.050571832671583	227.2470280696	198.1048228428	256.34742930111
1500M_Women	B	5.8933701899634	5.54387867824481	0.022754377504487	0.02313065673308	241.529119229264	239.833545397023	253.2248840611895
1500M_Women	G	5.45540784729686	5.13782885922414	0.022245347424987	0.021073769938616	237.8434015527149	235.45784541078	241.85995210079
200M_Men	B	0.164707464119639	0.13782885922414	0.002607723597716	0.00571071741400172	19.3025695627224	19.3025695627224	19.3025695627224
200M_Men	G	0.255874629140973	0.196247874104356	0.013460429070354	0.0095657556519024	19.466286546513	19.301091453671	20.001592746939
200M_Women	B	0.3713851863556	0.290368624854269	0.0216535642486616	0.0129826169105963	22.1165201659661	21.263152796278	22.869537818339
200M_Women	G	0.370387182108881	0.277117899198087	0.016378521262243	0.012234777006381	21.7129357381871	21.0004839417742	22.4253873300001
400M_Men	B	0.25315511012148	0.196748413975556	0.00560246605212897	0.004421188084629	44.3510578312253	43.834540980426	44.867583644081
400M_Men	G	0.278645850610955	0.5413970016981	0.01578523233336	0.0117287045097446	43.72906353515	42.2395237158806	45.2191975544194
400M_Women	B	0.530553194766558	0.524880811307212	0.010889615900351	0.010488428569004	49.4287185544189	49.180675168298	49.677602302308
400M_Women	G	0.41112735013861	0.295734798192441	0.0084304923883515	0.0060428157605827	49.3493122671432	48.144126482733	49.8543038914131
800M_Men	B	1.07307316335357	0.99931888145097	0.0094523267141014	0.00925695252378469	103.2608763257972	100.97123006411	105.5502270302
800M_Men	G	1.8887121627213	1.65363162862639	0.017939701945681	0.0157054955705802	101.512792329296	97.476133096081	105.35425484805
800M_Women	B	0.88974695077739	0.793222862172702	0.0075761758777064	0.0075427892585747	117.27966846518	115.16150570948	119.400778292318
800M_Women	G	0.89204554998214	0.833247288461841	0.0076651032193547	0.0071596751798536	115.62081795214	114.186613863382	117.05482680047

TABLE A4  
SUMMARY OF NORMALIZED RMSE AND MAE METRICS ACROSS MODELS 2, 3, 4

	Mean Normalized RMSE	Min. Normalized RMSE	Max. Normalized RMSE	Range Normalized RMSE	Mean Normalized MAE	Min. Normalized MAE	Max. Normalized MAE	Range Normalized MAE
Model2	0.0192	0.01008	0.02824	0.01816	0.0167	0.00876	0.02631	0.01755
Model3	0.01542	0.00569	0.04273	0.03703	0.01294	0.00442	0.03186	0.02744
Model4	0.03916	0.02036	0.06469	0.04433	0.03585	0.01874	0.06057	0.04183

TABLE A5  
NORMALIZED RMSE AND MAE COMPARISON FOR EACH RACE

Model	Race	Normalized LM_RMSE	Normalized KNN_RMSE	Normalized ARIMA_RMSE	Normalized LM_MAE	Normalized KNN_MAE	Normalized ARIMA_MAE	Invest RMSE	Invest MAE
B	100M_Men	0.0246471079981535	0.042726519834063	0.006656781273932	0.0210080033337	0.031884075377789	0.0488310778864077	LM	LM
G	100M_Men	0.01475959594581	0.001194910333333	0.035484780082442	0.011263266449443	0.015841779269614	0.047690276668618	LM	LM
B	100M_Women	0.017354965632082	0.003458692787275	0.034358813064554	0.012453306345454	0.0194536974668314	0.045645342310284	KNN	KNN
B	100M_Women	0.023508724348897	0.004746808701768	0.016298487739903	0.016298487739903	0.0238172115014703	0.0238172115014703	LM	LM
B	100M_Women	0.020271455311653	0.008494240000000	0.017650511111111	0.017650511111111	0.017650511111111	0.017650511111111	LM	LM
B	100M_Men	0.021871677433032	0.021811859583333	0.064888081880493	0.021871677433032	0.019250477188864	0.021871677433032	LM	LM
B	100M_Men	0.020271455311653	0.008494240000000	0.017650511111111	0.017650511111111	0.017650511111111	0.017650511111111	LM	LM
B	150CM_Women	0.0228196649312	0.023754377040877	0.023291551653395	0.023291551653395	0.0231366057427874	0.0225133044341683	LM	LM
B	100M_Women	0.0228196649312	0.023754377040877	0.023291551653395	0.023291551653395	0.0231366057427874	0.0225133044341683	LM	LM
B	200M_Men	0.021021165325159	0.008067235597812	0.009376951920044	0.008067235597812	0.007875314140212	0.030604210214072	LM	LM
B	200M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	200M_Men	0.018094392793438	0.016350635641686	0.073915859299542	0.012608195849308	0.0230289159715042	0.0230289159715042	KNN	KNN
B	200M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.015668780080002	0.005000000000000	0.04131181821223	KNN	KNN
B	400M_Men	0.018428197103835	0.004458799703534	0.034561887118126	0.				



TABLE A6

2028 OLYMPIC TIME PREDICTIONS AND 95% CONFIDENCE INTERVALS

Race	Medal	Predictions (secs)	95% Confidence Interval
100M Men	B	9.8699	[9.0634, 10.6764]
	G	9.8000	[9.4908, 10.1091]
100M Women	B	10.9520	[10.7877, 11.1164]
	G	10.6799	[10.1695, 11.1903]
1500M Men	B	215.4561	[212.2677, 218.6445]
	G	218.9286	[208.8991, 228.9582]
1500M Women	B	241.5291	[229.8334, 253.2249]
	G	237.6435	[233.4579, 241.8290]
200M Men	B	19.9326	[19.6701, 20.1950]
	G	19.4663	[18.9310, 20.0015]
200M Women	B	22.1165	[21.3632, 22.8699]
	G	21.7129	[21.0005, 22.4254]
400M Men	B	44.3511	[43.8345, 44.8676]
	G	43.7294	[42.2395, 45.2192]
400M Women	B	49.4287	[49.1807, 49.6768]
	G	49.0493	[48.1441, 49.9545]
800M Men	B	103.2609	[100.9712, 105.5505]
	G	101.4153	[97.4761, 105.3544]
800M Women	B	117.2797	[115.1616, 119.3978]
	G	115.6206	[114.1866, 117.0546]



Fig. A6 Women's 100m bronze medalist race times vs year, including 2028 predictions and 95% CI's

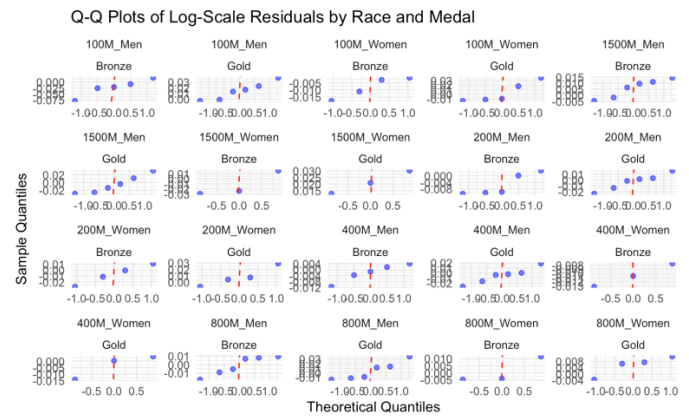


Fig. A7 QQ plots of log-scale residuals by race and medal

## APPENDIX B

### ADDITIONAL HISTOGRAMS OF HISTORICAL RACE TIMES, BY RACE, GENDER, AND MEDAL

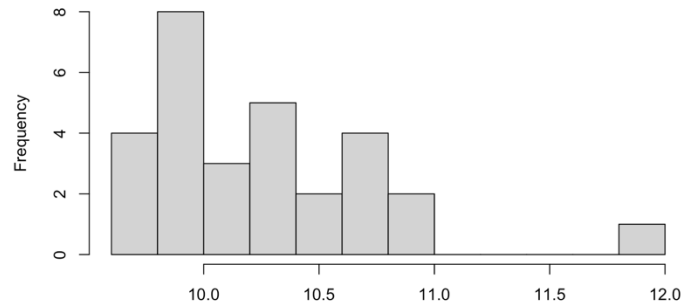


Fig. B1 Histogram of men's 100m gold medal Olympic finals race results

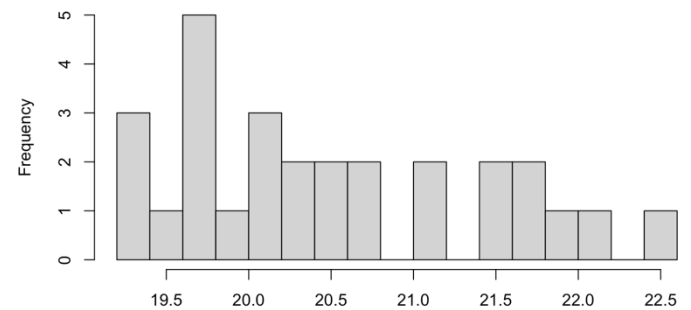


Fig. B2 Histogram of men's 200m gold medal Olympic finals race results

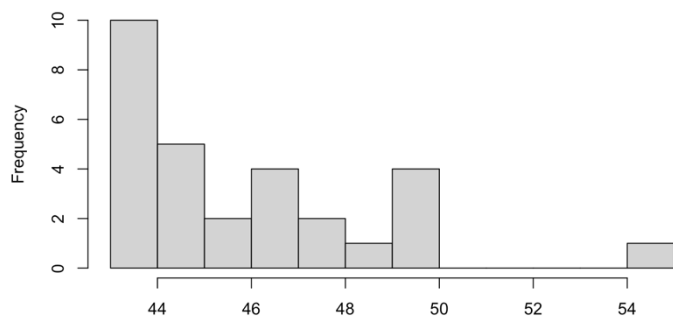


Fig. B3 Histogram of men's 400m gold medal Olympic finals race results

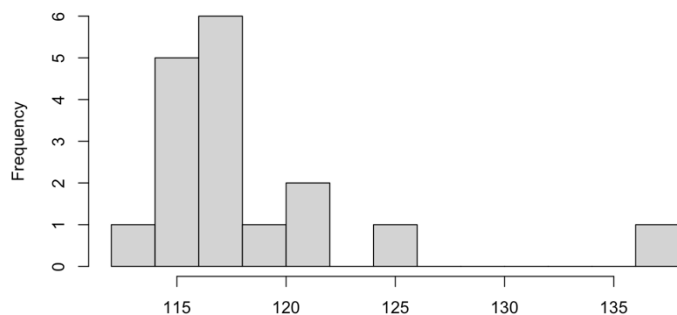


Fig. B4 Histogram of women's 800m gold medal Olympic finals race results

### APPENDIX C ADDITIONAL FIGURES DISPLAYING HISTORICAL PERFORMANCE PLOTS, FOR GOLD AND BRONZE MEDAL TRENDS

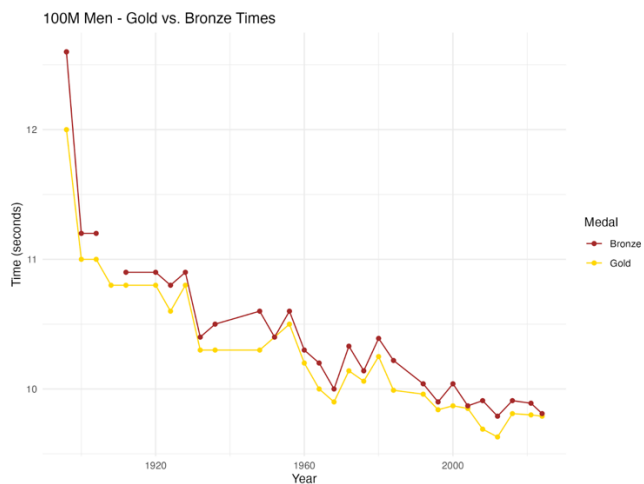


Fig. C1 100m men's Olympic historical trends



Fig. C2 100m women's Olympic historical trends

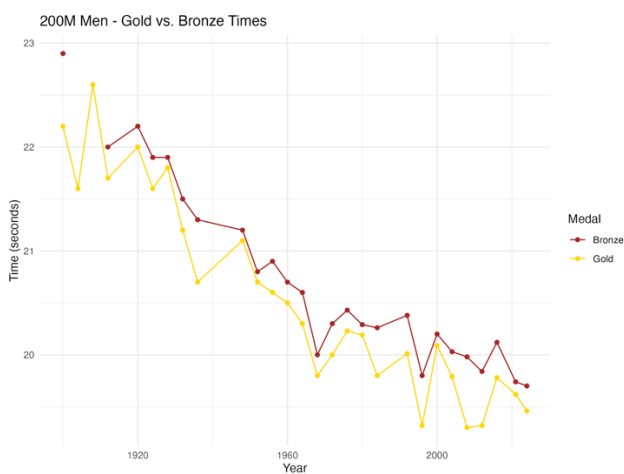


Fig. C3 200m men's Olympic historical trends

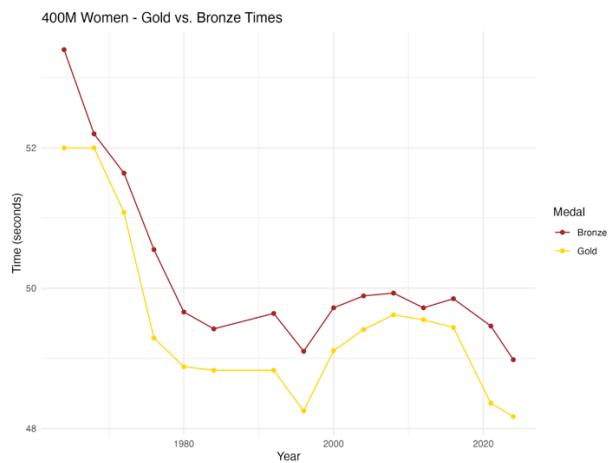


Fig. C4 400m women's Olympic historical trends

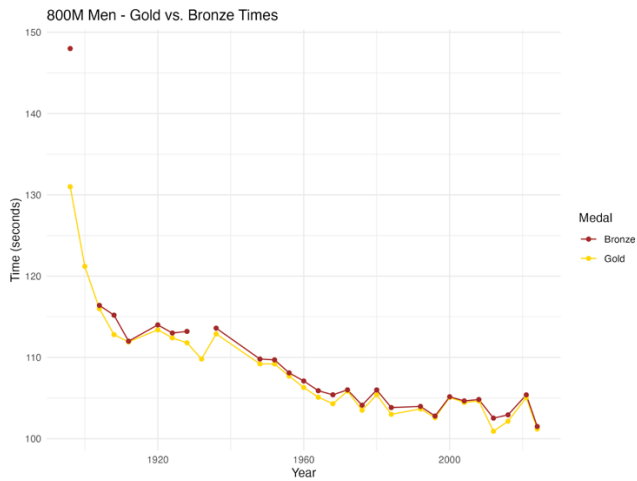


Fig. C5 800m men's Olympic historical trends

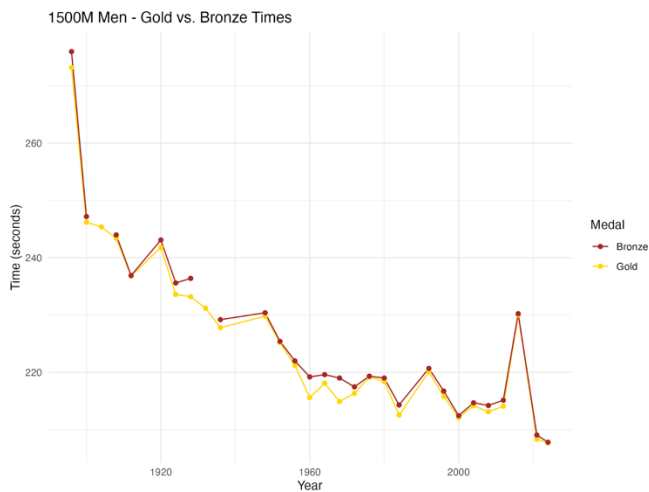


Fig. C6 1500m men's Olympic historical trends

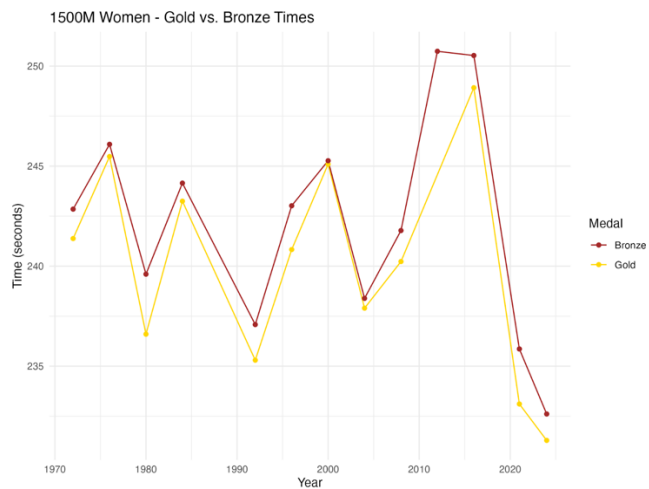


Fig. C7 1500m women's Olympic historical trends

## APPENDIX D ADDITIONAL FIGURES DISPLAYING HISTORICAL RESULTS, 2028 PREDICTIONS, AND 2028 95% CONFIDENCE INTERVALS

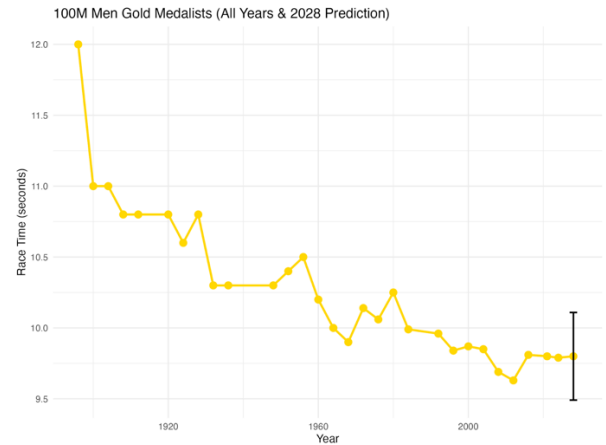


Fig. D1 Men's 100m gold medalist race times vs year, including 2028 predictions and 95% CI's

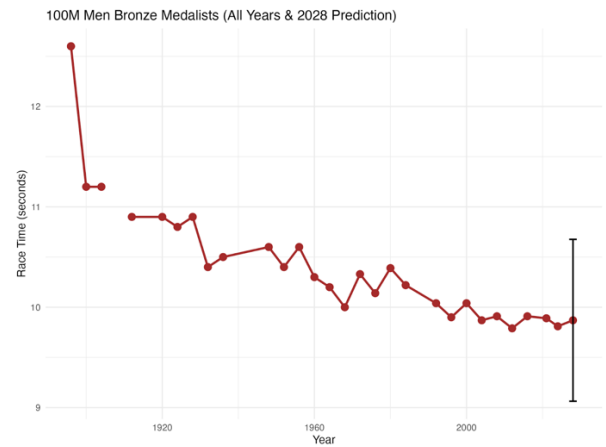


Fig. D2 Men's 100m bronze medalist race times vs year, including 2028 predictions and 95% CI's

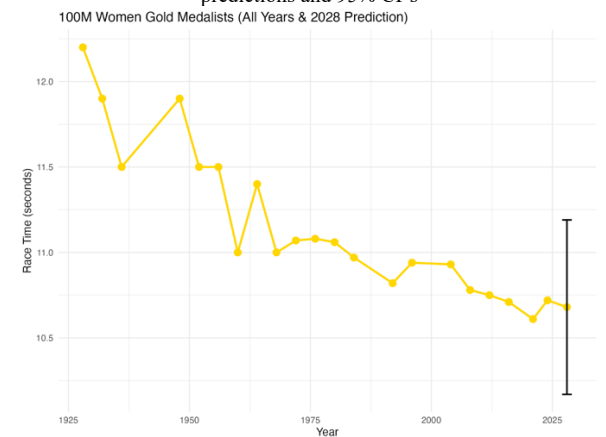


Fig. D3 Women's 100m gold medalist race times vs year, including 2028 predictions and 95% CI's

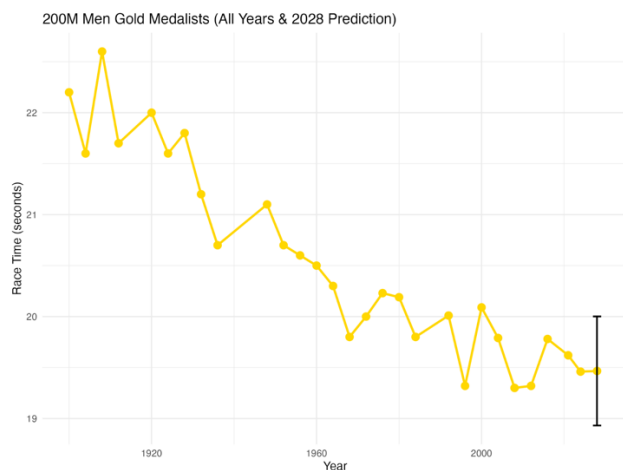


Fig. D4 Men's 200m gold medalist race times vs year, including 2028 predictions and 95% CI's

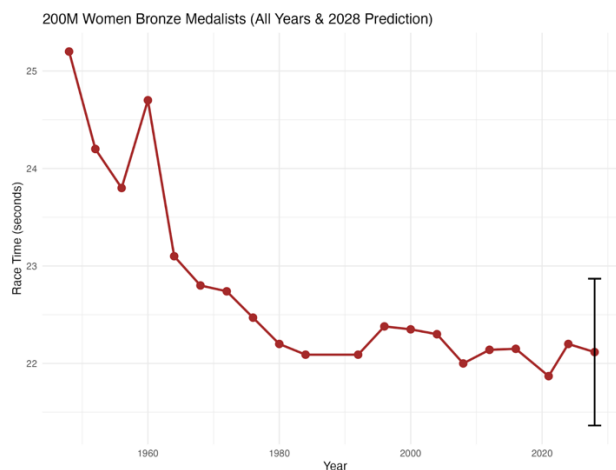


Fig. D7 Women's 200m bronze medalist race times vs year, including 2028 predictions and 95% CI's

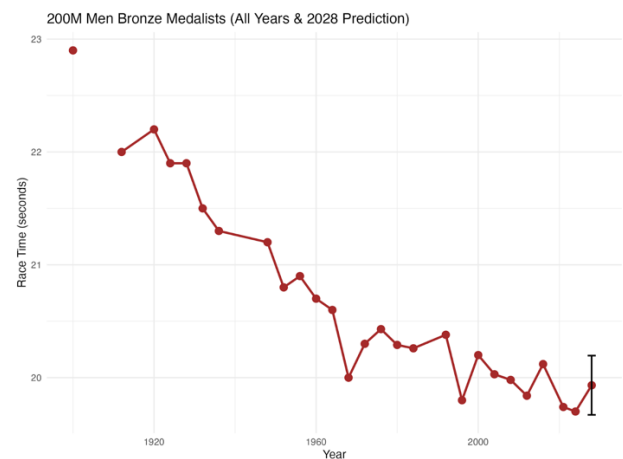


Fig. D5 Men's 200m bronze medalist race times vs year, including 2028 predictions and 95% CI's

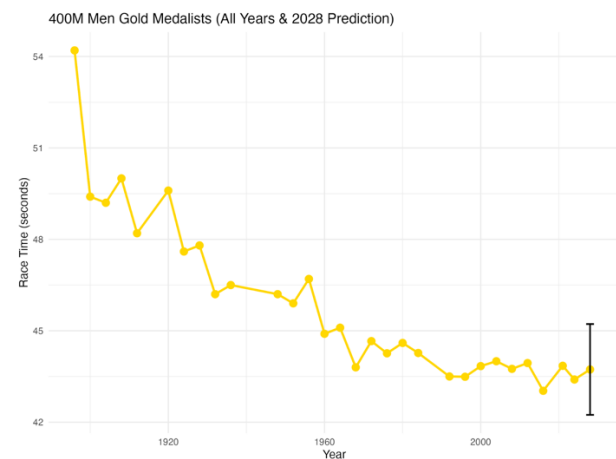


Fig. D8 Men's 400m gold medalist race times vs year, including 2028 predictions and 95% CI's

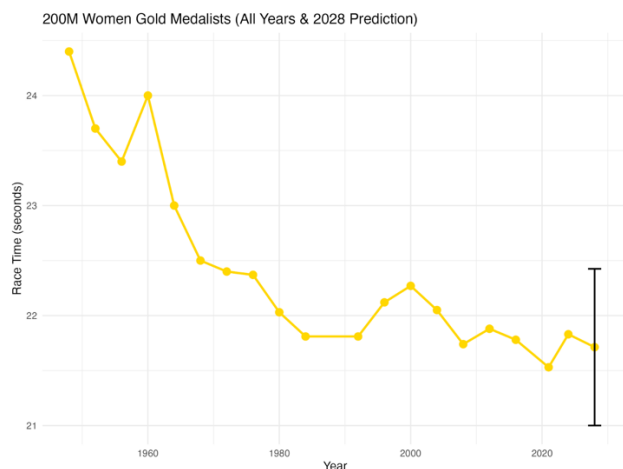


Fig. D6 Women's 200m gold medalist race times vs year, including 2028 predictions and 95% CI's

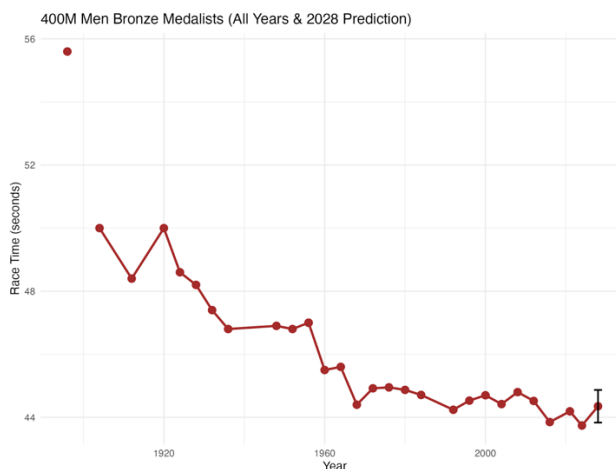


Fig. D9 Men's 400m bronze medalist race times vs year, including 2028 predictions and 95% CI's

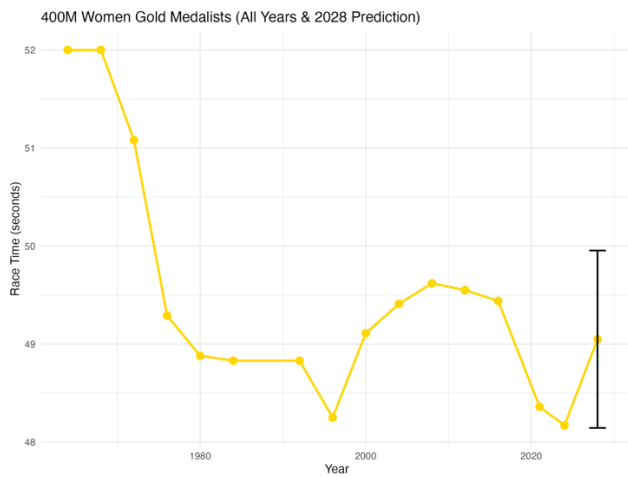


Fig. D10 Women's 400m gold medalist race times vs year, including 2028 predictions and 95% CI's

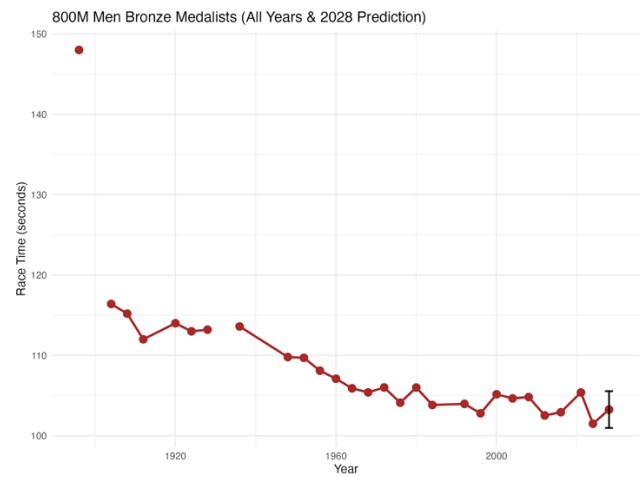


Fig. D13 Men's 800m bronze medalist race times vs year, including 2028 predictions and 95% CI's

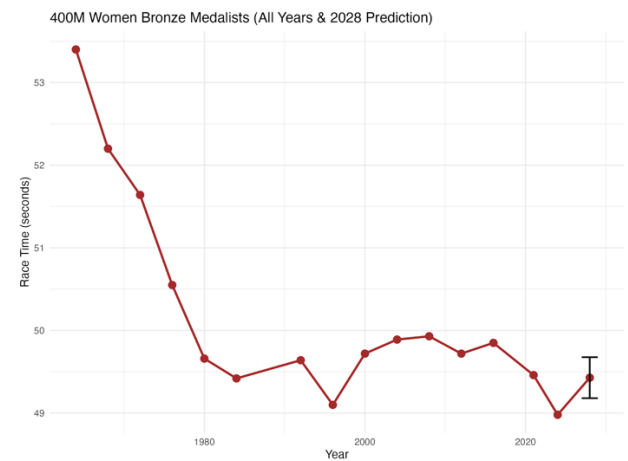


Fig. D11 Women's 400m bronze medalist race times vs year, including 2028 predictions and 95% CI's

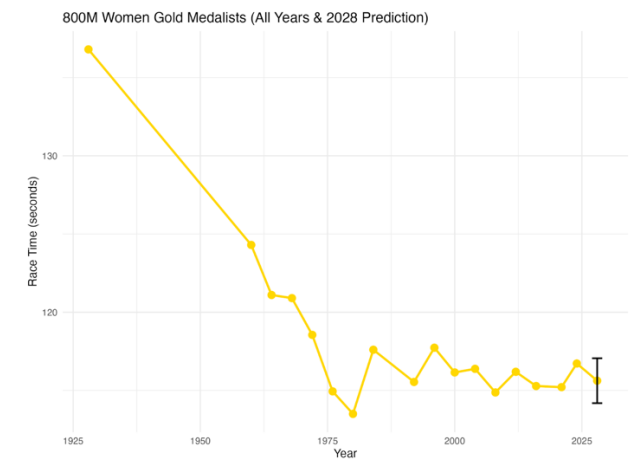


Fig. D14 Women's 800m gold medalist race times vs year, including 2028 predictions and 95% CI's

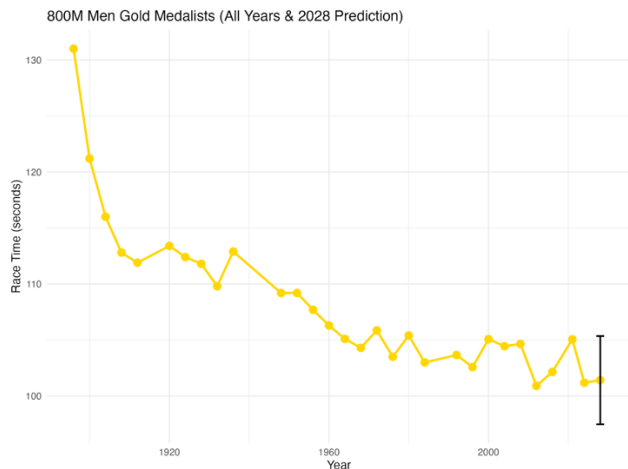


Fig. D12 Men's 800m gold medalist race times vs year, including 2028 predictions and 95% CI's

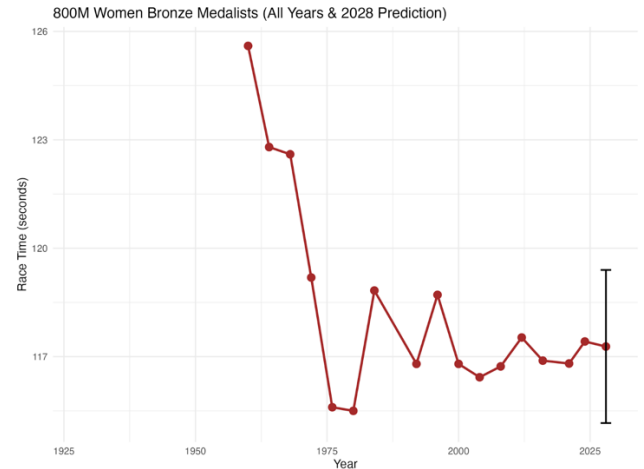


Fig. D15 Women's 800m bronze medalist race times vs year, including 2028 predictions and 95% CI's



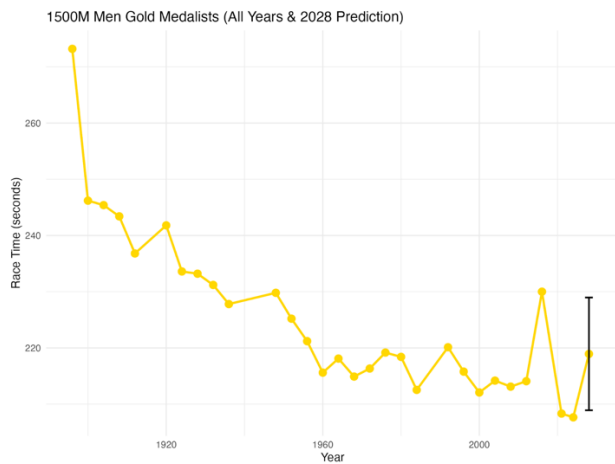


Fig. D16 Men's 1500m gold medalist race times vs year, including 2028 predictions and 95% CI's

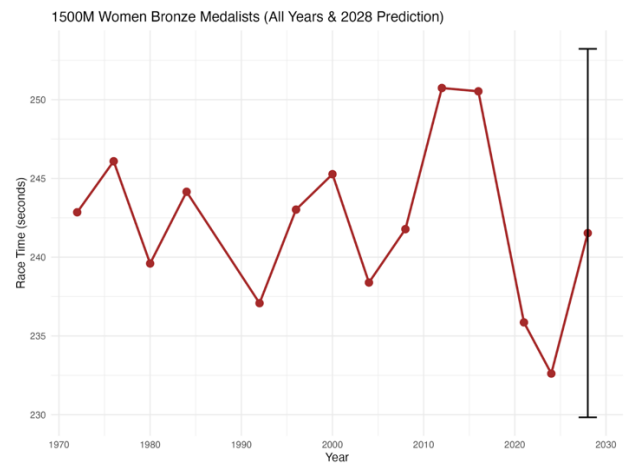


Fig. D19 Women's 1500m bronze medalist race times vs year, including 2028 predictions and 95% CI's

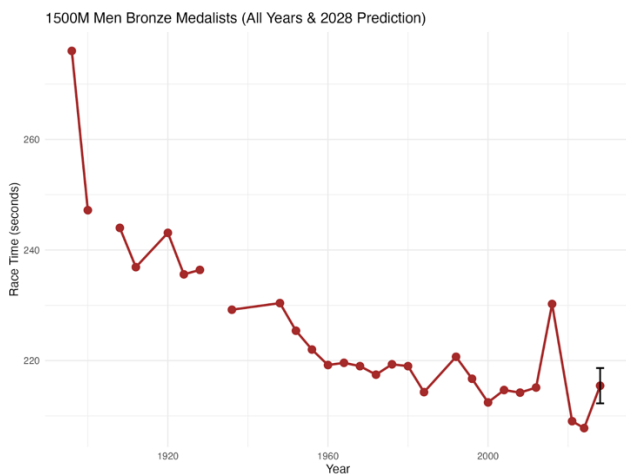


Fig. D17 Men's 1500m bronze medalist race times vs year, including 2028 predictions and 95% CI's

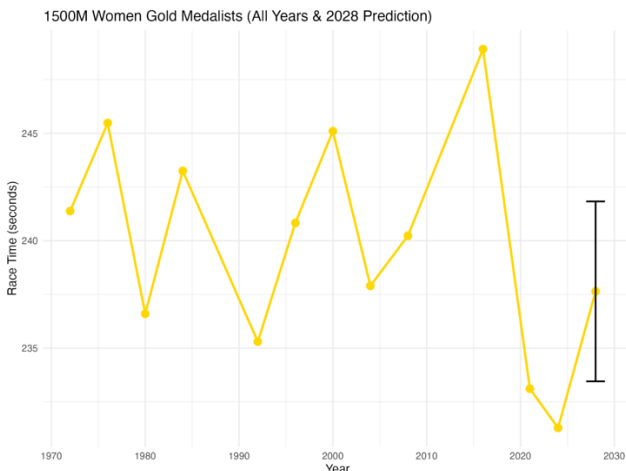


Fig. D18 Women's 1500m gold medalist race times vs year, including 2028 predictions and 95% CI's

## APPENDIX E ADDITIONAL FIGURES

Call:  
lm(formula = log(Result) ~ Year + male + Medal + Event, data = Train\_olympics)

Residuals:

	Min	1Q	Median	3Q	Max
	-0.046457	-0.014899	-0.002468	0.010242	0.218538

Coefficients: (1 not defined because of singularities)

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	7.254e+00	8.520e-02	85.145	< 2e-16 ***
Year	-1.244e-03	4.275e-05	-29.100	< 2e-16 ***
male	-2.471e+00	6.808e-03	-362.933	< 2e-16 ***
MedalG	-1.255e-02	2.832e-03	-4.431	1.26e-05 ***
Event100M Women	-2.369e+00	6.686e-03	-354.290	< 2e-16 ***
Event1500M Men	3.087e+00	6.050e-03	510.209	< 2e-16 ***
Event1500M Women	7.197e-01	7.992e-03	90.056	< 2e-16 ***
Event200M Men	6.974e-01	6.027e-03	115.697	< 2e-16 ***
Event200M Women	-1.660e+00	6.919e-03	-239.864	< 2e-16 ***
Event400M Men	1.495e+00	5.847e-03	255.627	< 2e-16 ***
Event400M Women	-8.582e-01	7.261e-03	-118.202	< 2e-16 ***
Event800M Men	2.354e+00	5.803e-03	405.731	< 2e-16 ***
Event800M Women	NA	NA	NA	NA

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.02666 on 346 degrees of freedom  
(8 observations deleted due to missingness)  
Multiple R-squared: 0.9994, Adjusted R-squared: 0.9994  
F-statistic: 5.356e+04 on 11 and 346 DF, p-value: < 2.2e-16

Fig. E1 Model 1 summary