

USOPC Age-Focused Competitive Analysis

Stergios Koutrouvelis - MS Analytics
Georgia Tech Spring 2022



Table of Contents

Introduction	2
Key Questions	2
The Data	2
Project Approach and Phases	3
Processing Framework	3
Exploratory Data Analysis	4
Feature Engineering	5
Machine Learning for Predictions	7
Exploratory Data Analysis	8
Analysis by Country of Origin and Event Type	8
Analysis by Athlete Country of Origin and Age	18
Analysis of Success by Country for Home vs Away Competition	28
Observations and Conclusions	29
Machine Learning for Predictions	30
Classification	30
Independent Event Data	31
Time Series Data	34
Regression	37
Independent Event Data	38
Time Series Data	38
Observations and Conclusions	39
Key Questions and Answers	42
Final Outcome and Recommendations for future work	44

Introduction

The scope of this analysis is for the United States Olympic and Paralympic Committee to gain a better understanding of the role that age plays in medal success for various Summer and Winter Olympic Sports. The original hypothesis is that peak age for winning medals varies by sport, event and gender. The goal is to evaluate this hypothesis, and if possible leverage this information (potentially in combination with other factors) to be able to more accurately predict which athletes will have a greater chance of winning medals in upcoming Olympic events. By leveraging the data provided by the USOPC, we will conduct a thorough analysis that can help us understand the health of Team USA's Olympic pipeline compared to those of the top medal-winning countries primarily through the lens of age.

Key Questions

The questions that this analysis will try to answer are the following:

- Can a country's percentage of athletes competing at an Olympic Games while in their peak age ranges in their respective sports predict medal success?
- Which athletes who competed in Tokyo 2021 (Beijing 2022) will be hitting their peak age in their respective sports in Paris 2024 (Milan 2026)?
- Can the number of young, "pre-peak" athletes that competed in Tokyo predict medal success for their countries in Paris 2024 and Milan 2026?
- Have many promising athletes who did not compete in Tokyo (and Beijing) but will reach peak age in Paris (and Milan) do Team USA and other countries have in their Pipelines?

At the end of this report, we will summarize our findings around the above questions as well as suggestions for extending the presented work for further improvements and adaptations.

The Data

The USOPC has provided an excellent and consistent set of datasets. Each dataset corresponds to one specific sport and covers multiple events of this sport. For example, the dataset for Snowboarding (which will be the main focus of this analysis) contains multiple different events such as Snowboardcross, Big Air, Slopestyle etc. Every row in the dataset represents one athlete's information and their performance in the specific event. Some of the key columns of these datasets are the following:

- Athlete Name and personal ID
- Athlete's nation
- Age of athlete at the time of the event
- Type of competition (e.g. Olympic Games, World Championship etc.)
- Sport/event or discipline
- Athlete Placement/Rank and Medal
- Result

One key characteristic of the datasets is the consistency of the schema. This enables the development of a framework that is scalable and can be used across sports in an automated way. For this reason, we opted for the development of an API that can ingest, transform and process the data and also plot consistent visualizations and develop machine learning models.

Project Approach and Phases

The approach of this project has emphasized not only on developing a standalone analysis but also on reproducibility and standardization by taking advantage of the data consistency as described previously.

For the purpose of this project we are going to be working with the **Snowboarding** dataset and more specifically with the **Snowboardcross** events, as this type of event has the largest number of data points and total athlete participation. The approach can be replicated and be applied and scaled to any dataset and event provided by the USOPC.

The progression can be split into the phases outlined below.

Processing Framework

During this phase, after some initial dive into the data, an API has been developed to make processing easier by developing functions and queries that can be used repeatedly for feature creation, visualizations and machine learning model development. In the initial stage, after going through the data, it was determined that the following columns contain useful information:

- **Class:** This is the competition class. This column can take the values of Elite, Juniors or Youth and is determined primarily by the age limitations of the particular event.
- **Competition Date:** As the name suggests, this is the date the competition has taken place.
- **Competition City:** This is the city where the competition has taken place. It is worth mentioning here that this does not include the country's name but only the city.
- **Competition Country:** This column is being generated from the framework. Based on the competition city column, API calls are being made through the Geopy library to get the competition country. The values are stored in a JSON file to speed up processing as depending on the size of the data, multiple calls can be very time consuming and computationally expensive.
- **Event/Athlete Gender:** As the name suggests this is the gender of the event that also matches the gender of the athlete.
- **Event:** This is the name of the event within the specific sport. For example, processing Snowboarding that is the focus of this analysis, individual events can be Snowboardcross, Slopestyle, Big Air etc.
- **Medal:** This column can take the values G, S, B or None indicating whether the athlete won a medal and if so what medal that was. In later stages of the processing a new

column will be created called won_medal which will be binary and will be used as one of our response variables.

- **Country:** This column represents the athlete's country of origin.
- **Home Competition:** This is one of the engineered columns based on the event country and the athlete's country of origin. If they match, the column value is true, indicating that the athlete is competing at an event held at home. We strongly believe that competing in the home country is an influencing factor in predicting success and we are going to validate this hypothesis based on both our exploratory analysis and the machine learning models.
- **Age:** This is the age of the athlete at the time of the event.
- **Rank:** The rank of the athlete at the event.

Exploratory Data Analysis

After the initial data cleanup, we performed exploratory data analysis across some major axes, in order to understand and identify patterns in the data. The main axes the analysis focused on are the following:

- Country
- Sport
- Event
- Age
- Gender
- Competition Location
- Medal Winning

Combining the factors above, the outline of the exploratory data analysis and graphs produced were the following:

- Exploratory Data Analysis by Athlete Country of Origin and Event to determine the following:
 - **Medal Rate:** Ratio of medals won over total athlete participation for a given country and given event. This factor provides some insight into how successful a country is in a particular type of event overall.
 - **Participation Count:** This parameter represents the total number of athletes from a particular country who have participated in the event under consideration. As countries vary in population and athletes sent to type of event, this will be a key factor in observing how many athletes participate in an event from each country and whether the sample is sufficiently large in order to extract conclusions.
 - **Participation Rate:** This parameter represents the distribution of a country's athletes across different events. It is computed as the ratio of a country's athletes participating in an individual event over the total number of a country's athletes.

This is insightful in terms of showcasing whether different countries participate in certain events more heavily than others.

- Exploratory Data Analysis by Athlete Country of Origin and Age to determine the following:
 - **Medal Rate:** This is the ratio of medals won across different ages. It is computed as the winning athletes for each age group over the total country's athletes in the particular age group. This is independent of the type of the event and focuses only on country of origin and age.
 - **Participation Count:** This represents the total number of athletes in the dataset broken down by country of origin and age group.
 - **Participation Rate:** This variable represents the proportion of a country's athlete that belong to different age groups irrespective of sport. It is computed as the number of a country's athletes in a given age group over the total number of a country's athletes. This is an indicator of whether there is a pattern of specific countries sending athletes of certain age groups to competitions.
- Exploratory Data Analysis for home competition success rate by country. In this analysis, we will be focusing on how different countries perform at home vs away. The graph will show the success rate for competitions held at home and away by computing the successful events over the total events. This will be insightful in determining how competing at home increases chances of success.

After developing the graphs we are going to discuss the observations and results.

Feature Engineering

Data preparation and feature engineering is needed for training and testing machine learning models. While evaluating our approach, two variations of training data have been considered.

In the initial phase, we opted for the simplest version of the data and treated each datapoint as an independent event. The feature set that has been utilized contained the following variables:

- **Gender:** Binary variable indicating the gender of the event and athlete. This is 1 if gender is male and 0 if gender is female
- **Country:** The athlete's country of origin after being processed with a one-hot encoding approach, where each column represents a country and is binary.
- **Is home competition:** Binary variable indicating whether the competition is being held at an athlete's home country or abroad.
- **Age:** The most critical factor for the purpose of this analysis indicating the age of the athlete at the time of the event.
- **Rank:** The rank of the athlete in the corresponding event (to be used as regression response variable).

- **Won_medal:** Binary variable indicating whether the athlete won a medal in the corresponding event (to be used as classification response variable)

In the second approach, we generated time series features, to leverage the power of past performance as an indicator of future performance. In order to achieve that, we compiled a dataset using the following procedure. For every event and athlete, we compiled the following information:

- Person ID
- Event date
- Event and sport name
- Home competition variable
- Athlete Rank in the event
- Won medal binary variable

For the attributes above and using the athlete-event data points defined above, we collected the three events the athlete participated for the same sport preceding the event in reference. This results in the following compilation:

- 4th preceding athlete event for given sport (feature prefix “ft_”)
- 3rd preceding athlete event for given sport (feature prefix “t_”)
- 2nd preceding athlete event for given sport (feature prefix “s_”)
- Event in reference (1st) athlete event for given sport. This is used as our prediction event (feature prefix “f_”).

So the time-series dataset looks as follows:

- **Person_ID:** The athlete identifier.
- **Country:** Athlete's country of origin, converted to one-hot-encoding variable based on a subset of medal winning countries.
- **Event:** event name (snowboardcross, halfpipe etc.).
- **Gender:** Binary variable indicating event and athlete gender.
- **4th latest competition date:** Date of the fourth to latest competition.
- **4th latest competition rank:** Rank of athlete in the respective competition.
- **4th latest competition athlete age:** Age of athlete during the respective competition.
- **4th latest competition won_medal:** Binary variable indicating whether the athlete won a medal in the respective competition.
- **3rd latest competition date:** Date of the third to latest competition.
- **3rd latest competition rank:** Rank of athlete in the respective competition.
- ...
- **2nd latest competition won_medal:** Binary variable indicating whether the athlete won a medal in the respective competition.
- **Latest competition date:** Date of the latest competition.

- **Latest competition rank:** Rank of athlete in the latest competition (to be used as regression response variable).
- **Latest competition athlete age:** Age of athlete in the latest competition.
- **Latest competition won_medal:** Binary variable indicating whether the athlete won a medal in the latest competition (to be used as classification response variable).
- **Latest competition is_home_competition:** Binary variable indicating whether the competition was held in the home country.

The features will be used in both approaches for classification and regression with different response variables as described in the next section.

Machine Learning for Predictions

In this phase, we are going to develop several machine learning models for predicting the outcome of a specific athlete participating in an event, using both the independent event data and the time series data and comparing the results.

The prediction can be treated as a classification or a regression problem depending on the response variable.

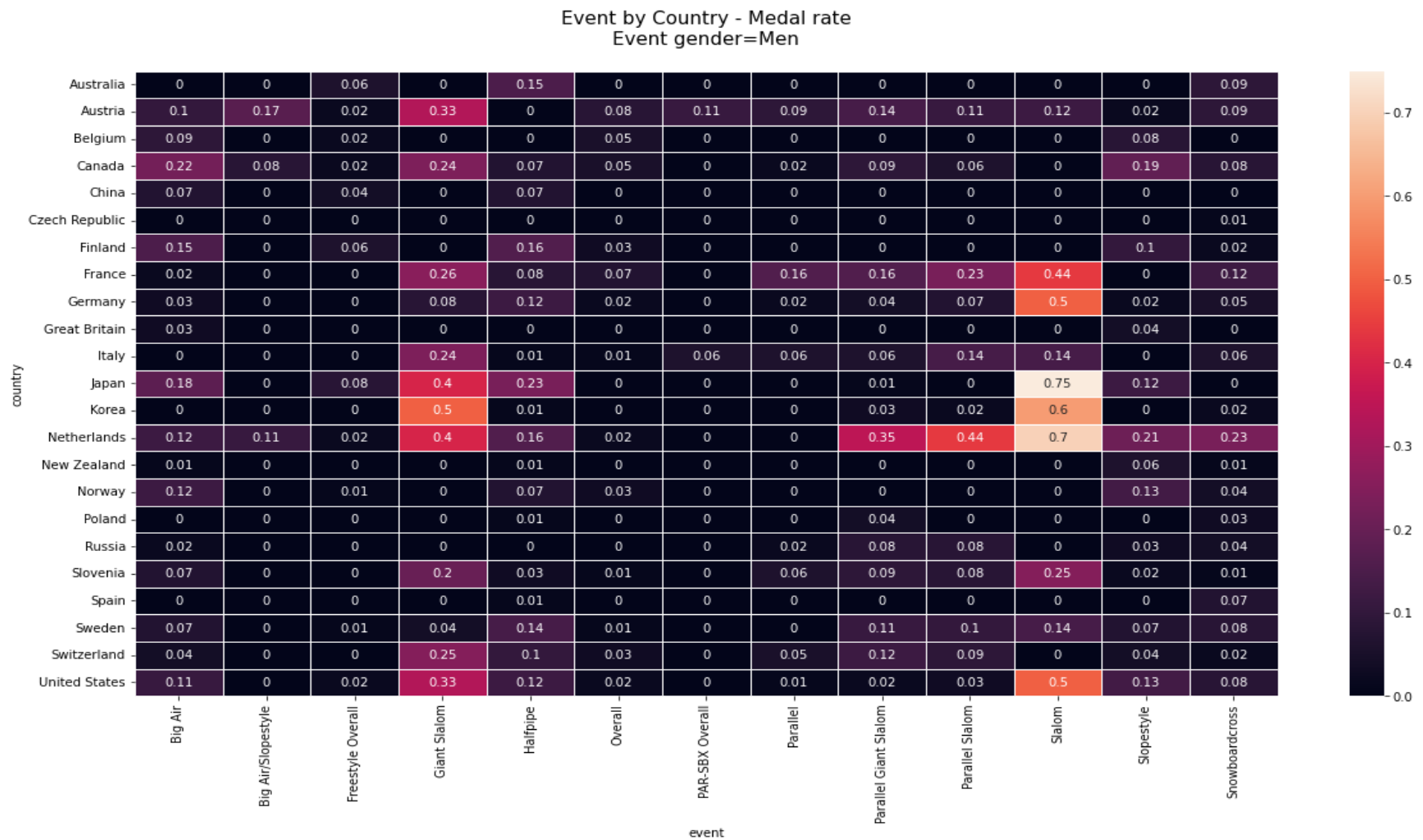
The first approach will focus on the won_medal variable as the response variable and treat the problem as classification. As a reminder, won_medal is 1 if the athlete won a medal in a given event or 0 if they did not. As expected this introduces heavy imbalance in the dataset, since for any given event, only three athletes win a medal and the majority does not. Through our analysis it has been determined that this ratio is around 1:12 for the majority of the events. The processing framework addresses this issue by passing class weights as hyperparameters when training the models. The alternative approach is to use the rank variable as our response variable and treat the problem as regression. Rank is the final ranking of the athlete in the event.

The features that will be used for prediction, and especially age, are highly nonlinear. For this reason we opted for ensemble methods that provide the best of-the-shelf accuracy and can handle nonlinear relationships and trends very well. The two algorithms that we will use are Random Forest and Gradient Boosting Decision Tree. In the hyperparameter tuning phase, we are going to employ cross-validation with 5 folds in order to determine the hyperparameters in an automated way and the best model's results will be reported in each case..

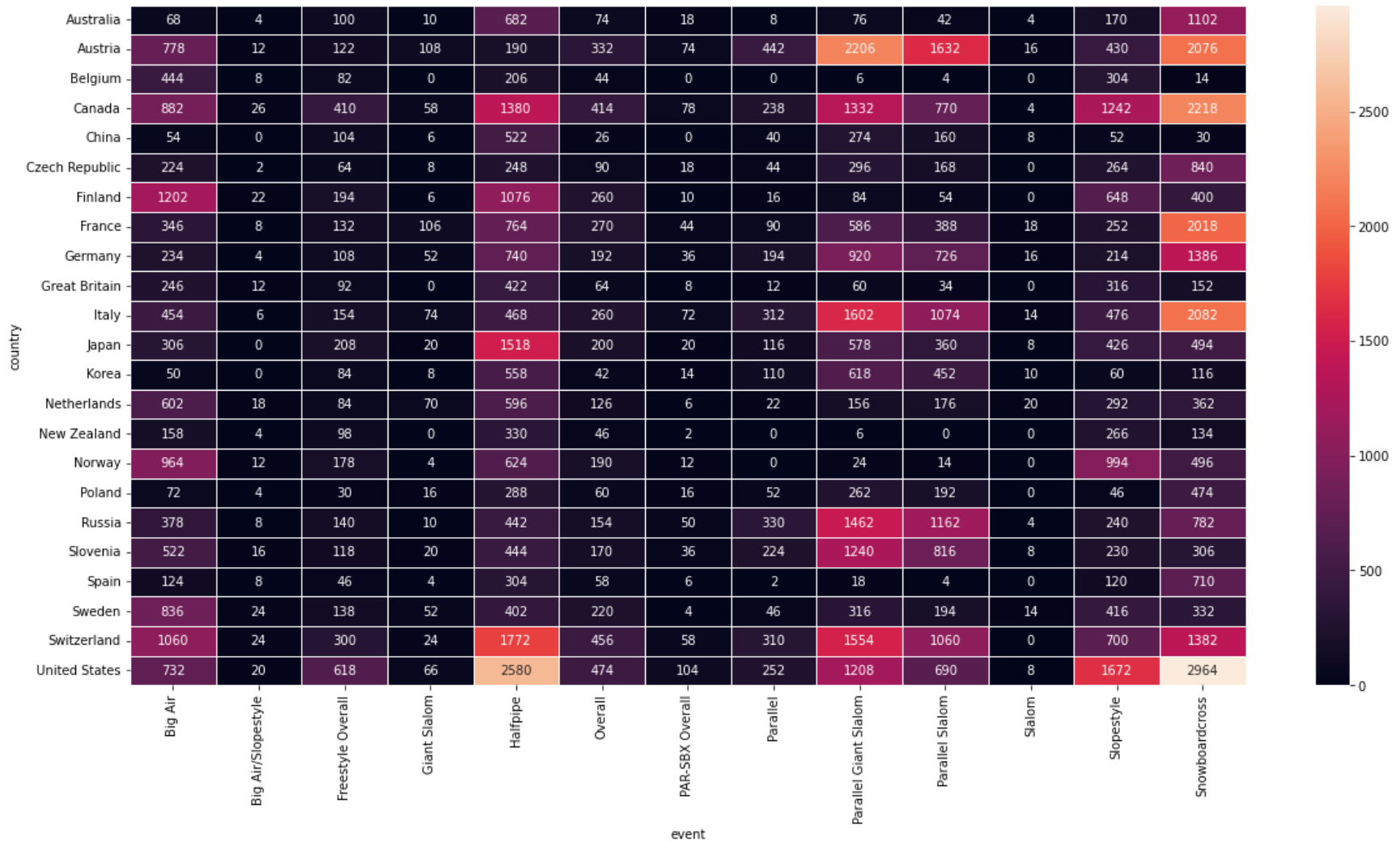
For the classification problem the metrics that we will use in evaluating the performance are precision, recall and f1-score. The regression algorithms will be evaluated on mean absolute error and mean squared error

Exploratory Data Analysis

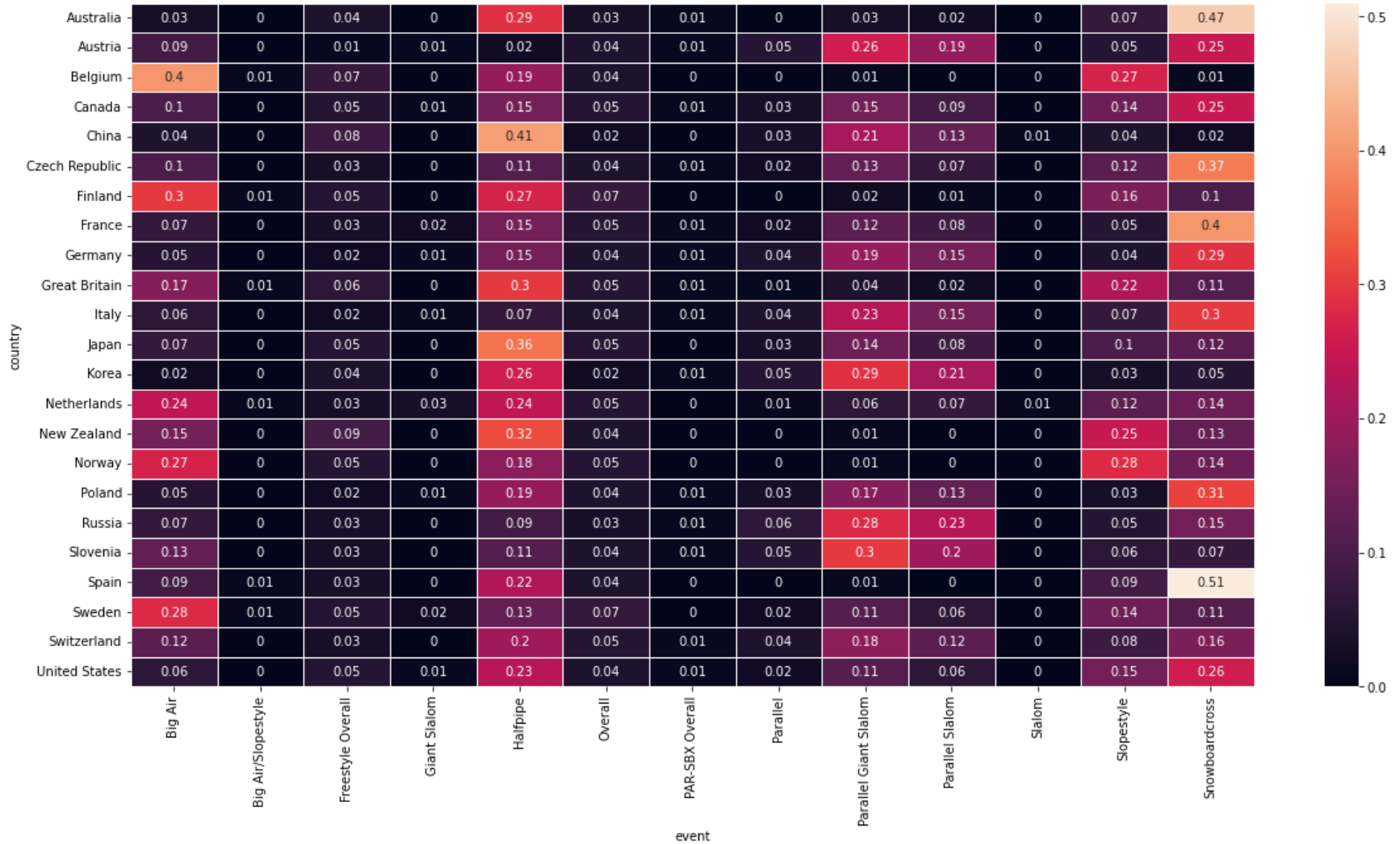
Analysis by Country of Origin and Event Type



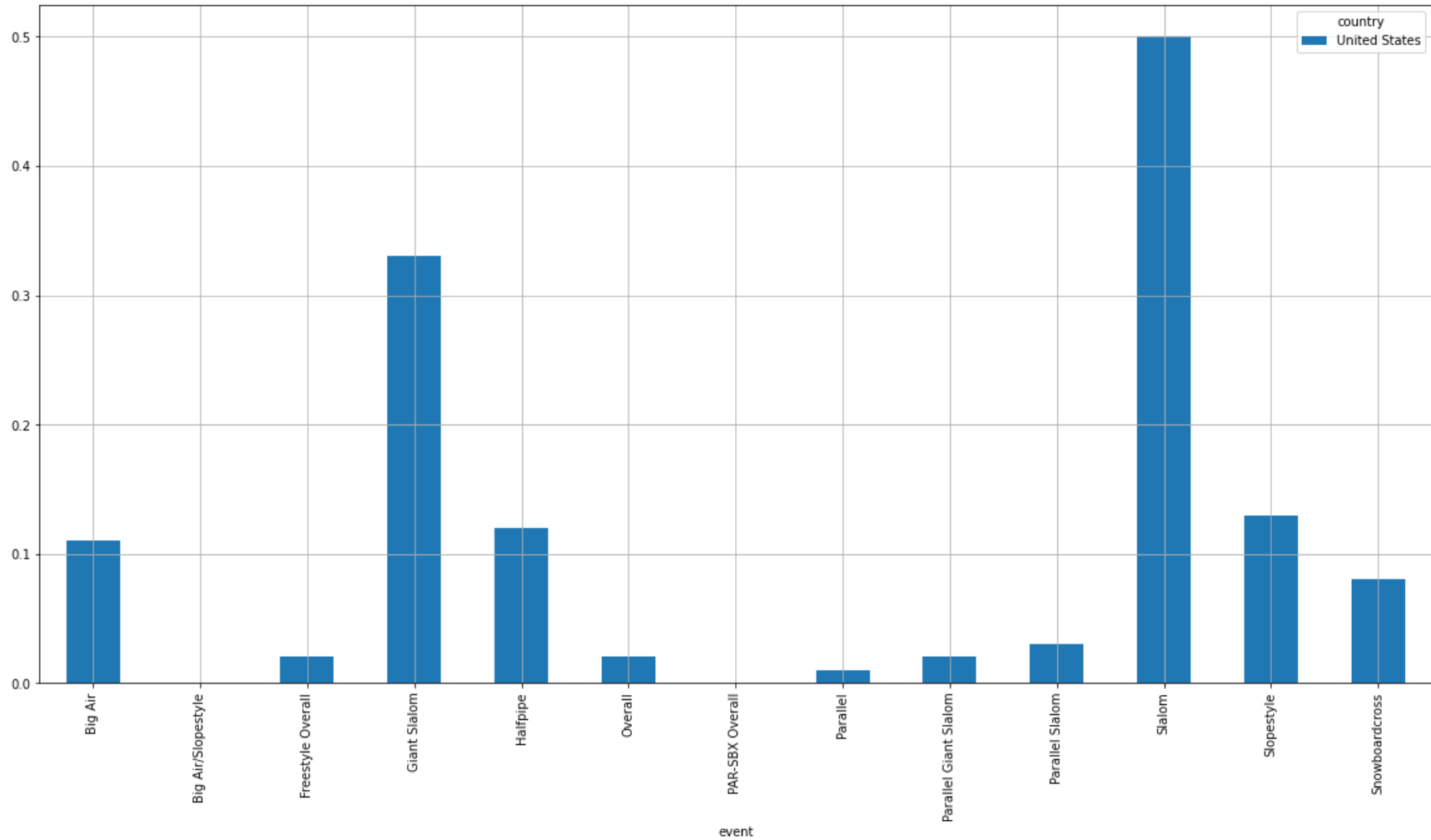
Event by Country - Participation count
Event gender=Men



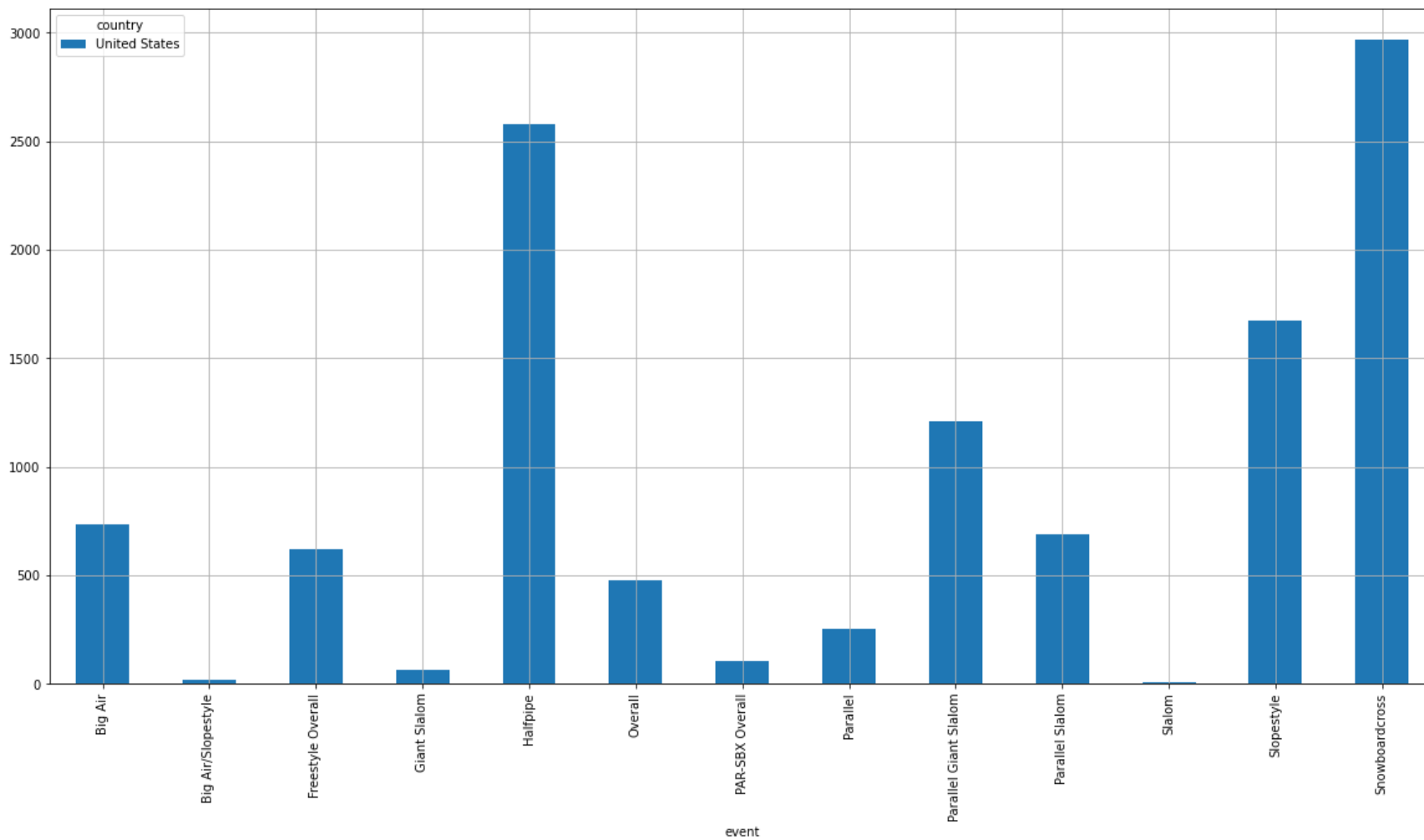
Event by Country - Participation rate
Event gender=Men



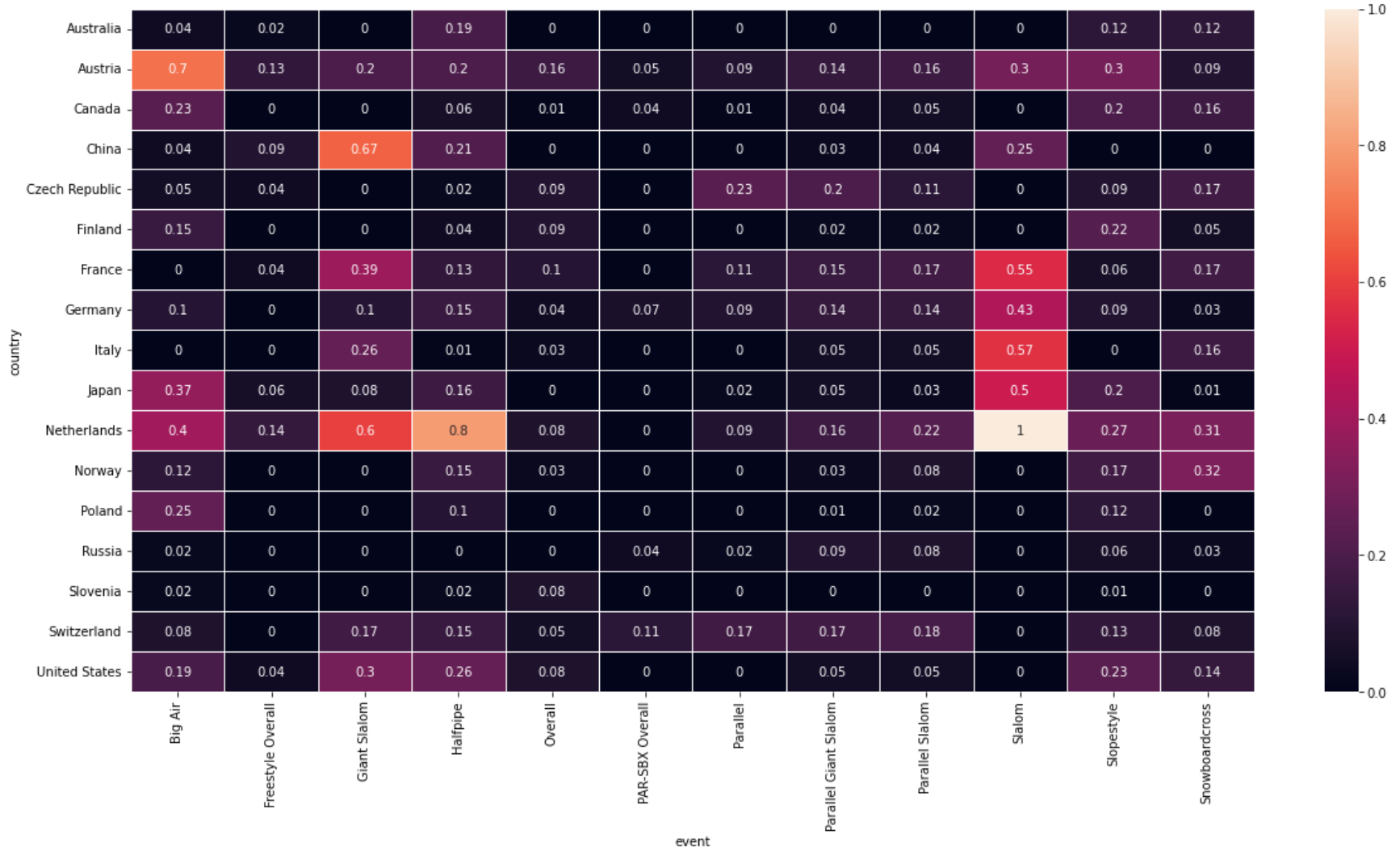
Event by Country - Medal rate
Event gender=Men - Country=United States



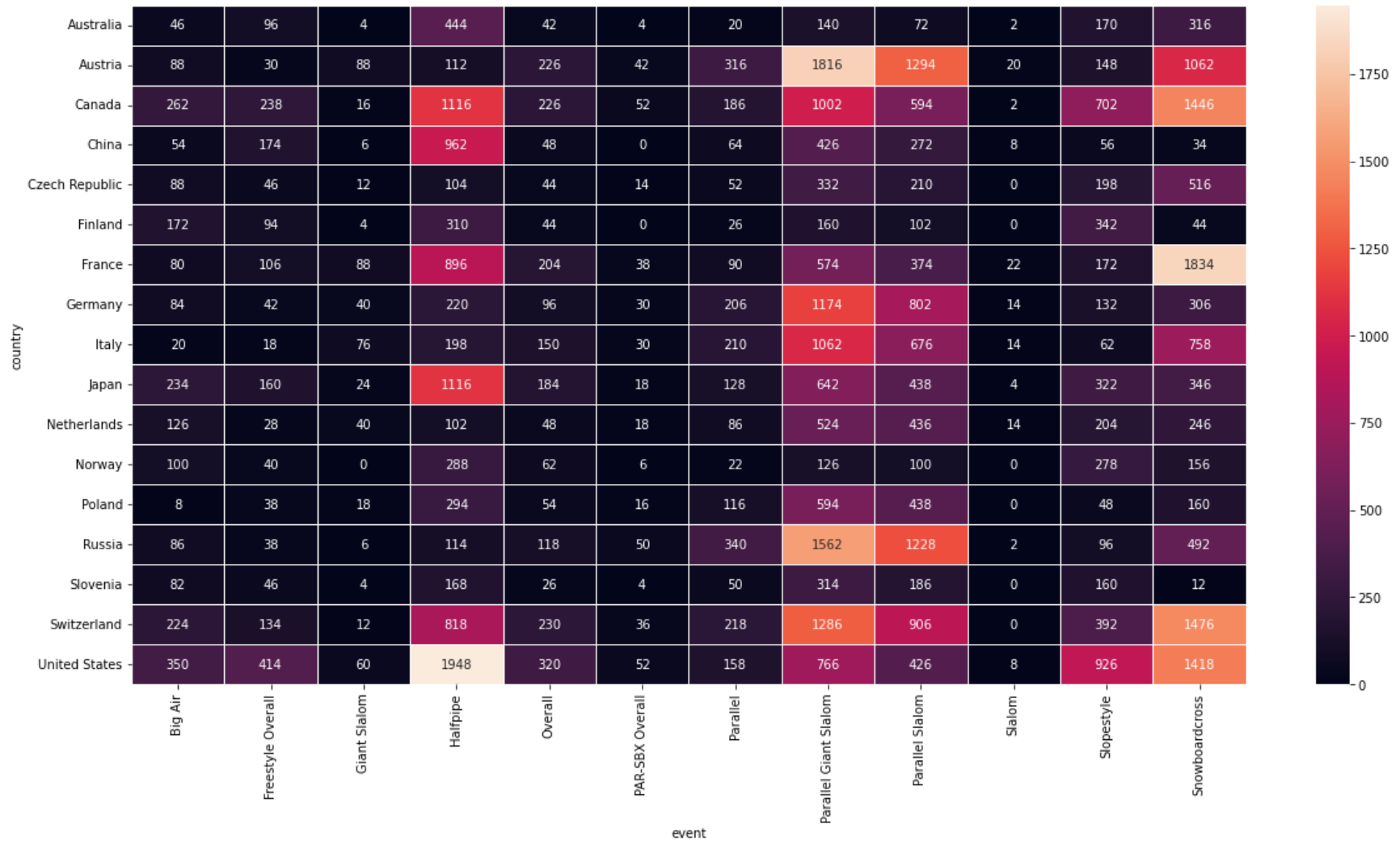
Event by Country - Participation count
Event gender=Men - Country=United States



Event by Country - Medal rate
Event gender=Women



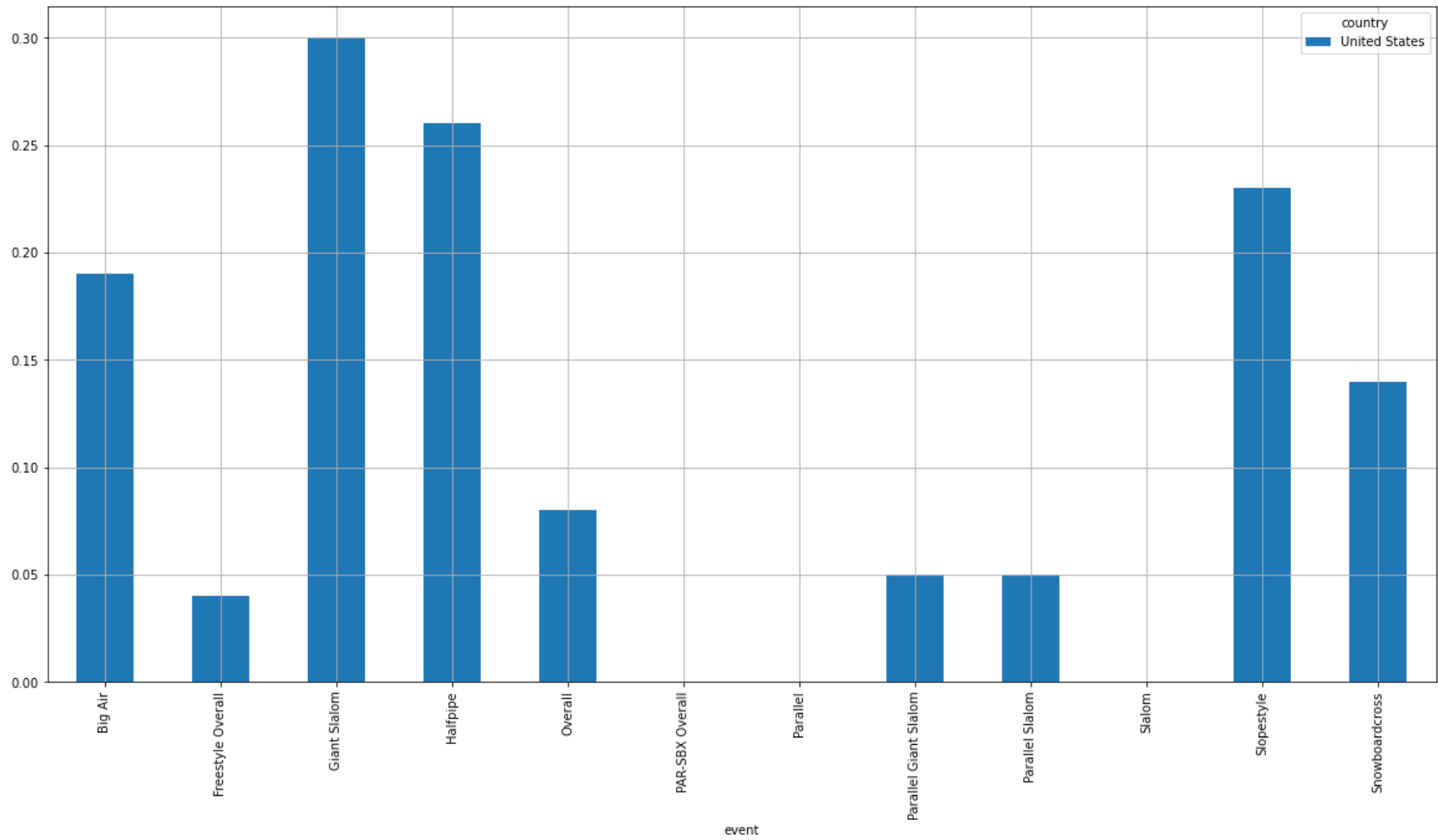
Event by Country - Participation count
Event gender=Women



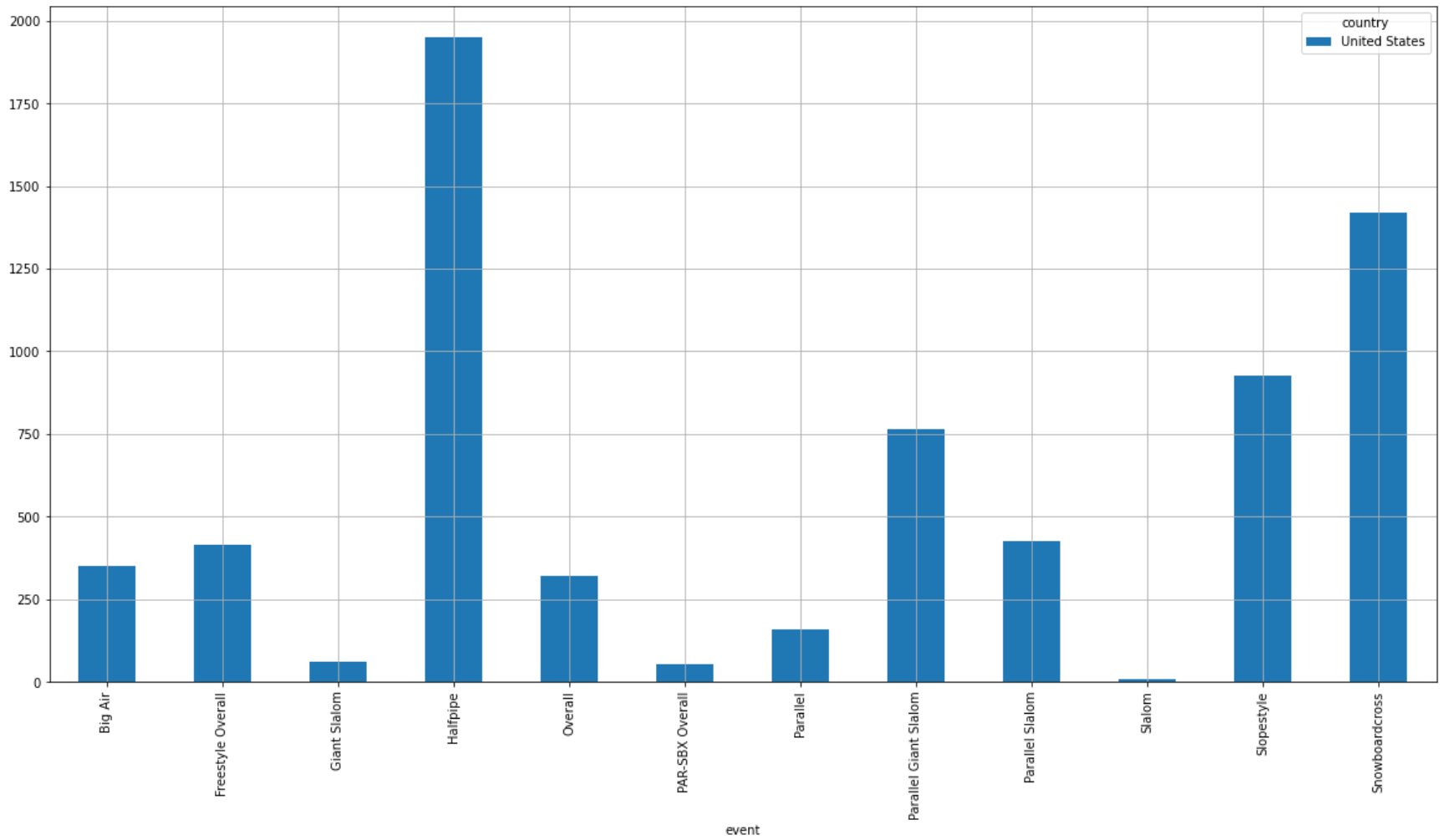
Event by Country - Participation rate
Event gender=Women



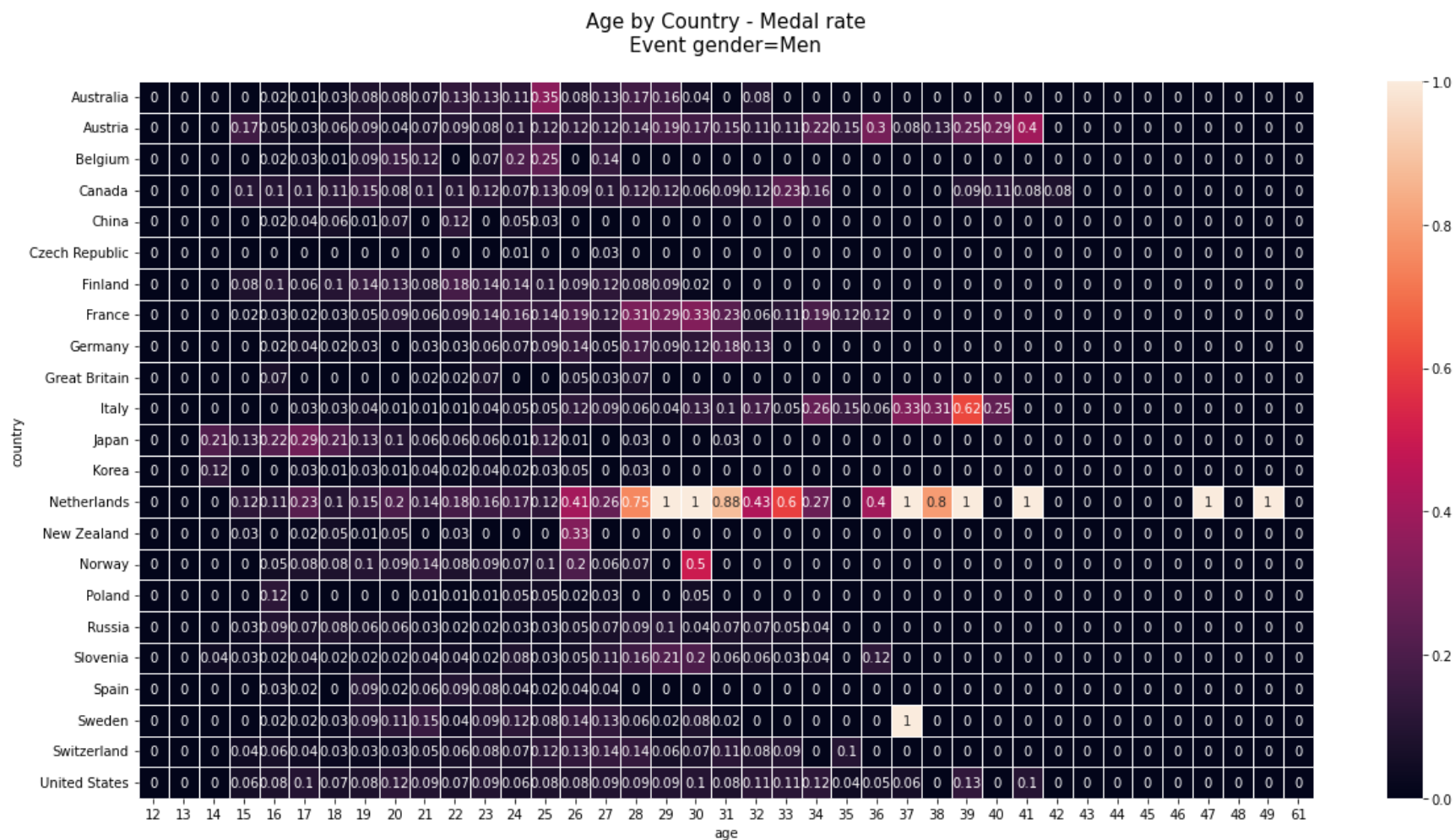
Event by Country - Medal rate
Event gender=Women - Country=United States



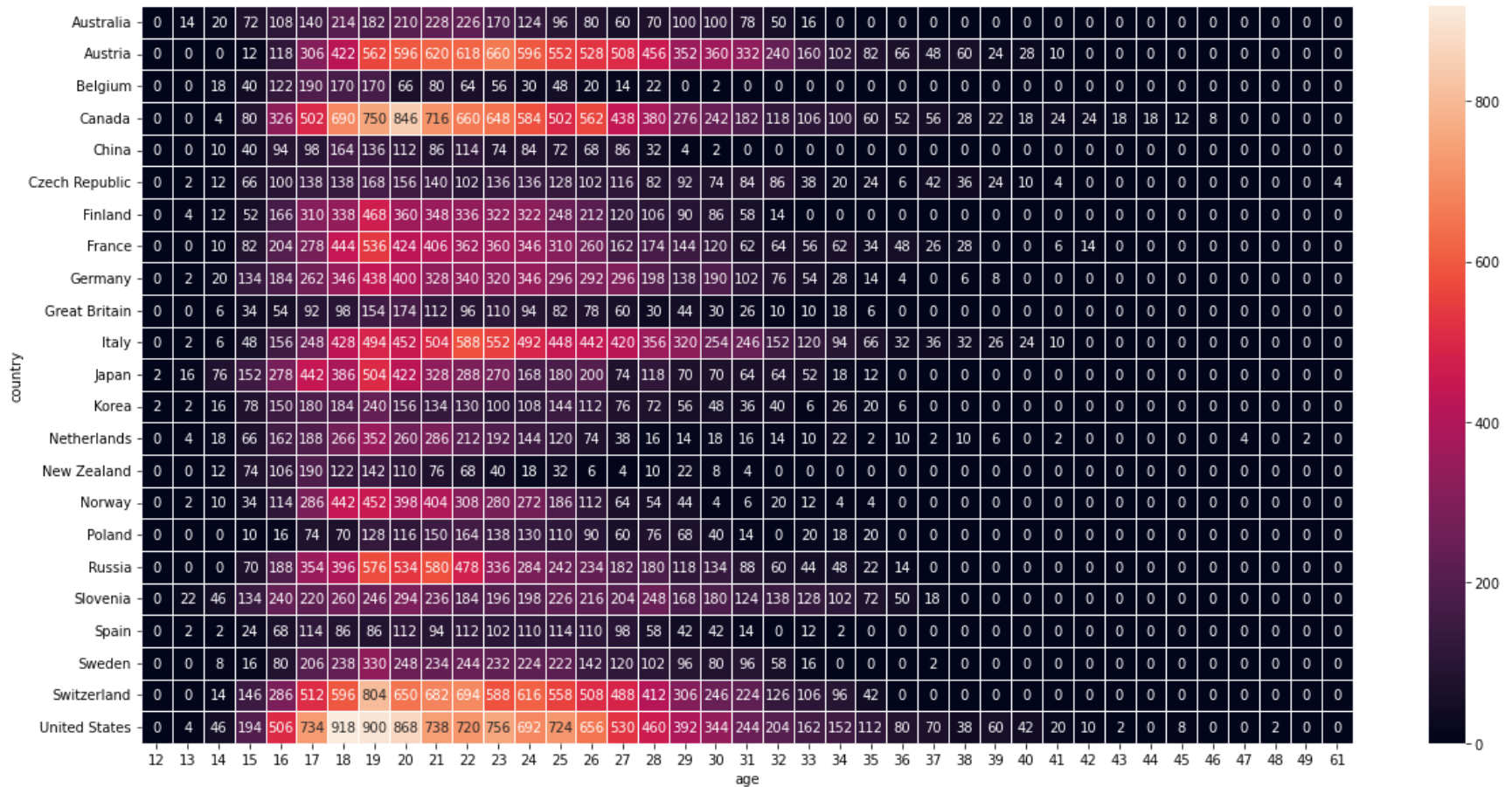
Event by Country - Participation count
Event gender=Women - Country=United States



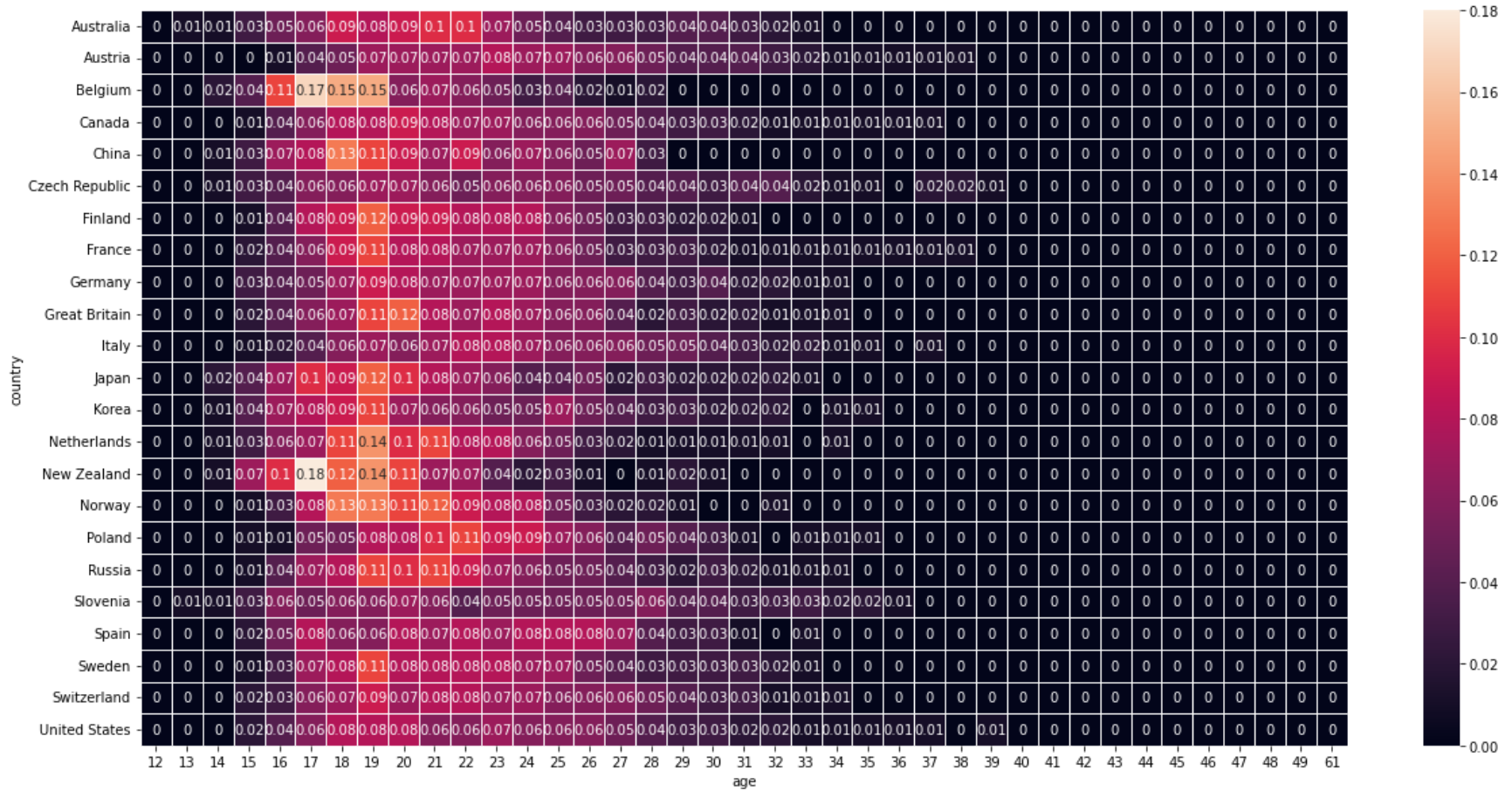
Analysis by Athlete Country of Origin and Age



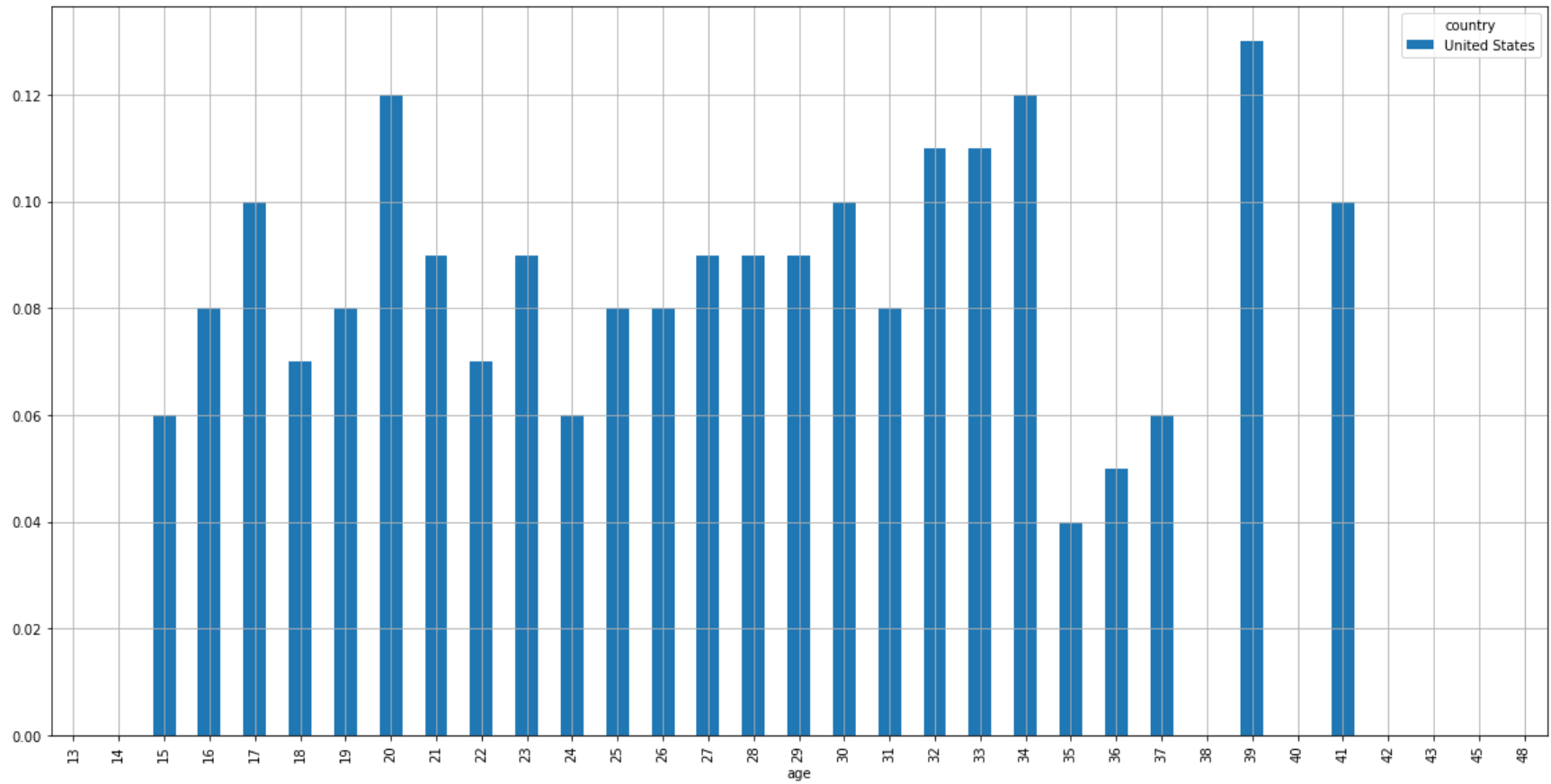
Age by Country - Participation count
Event gender=Men



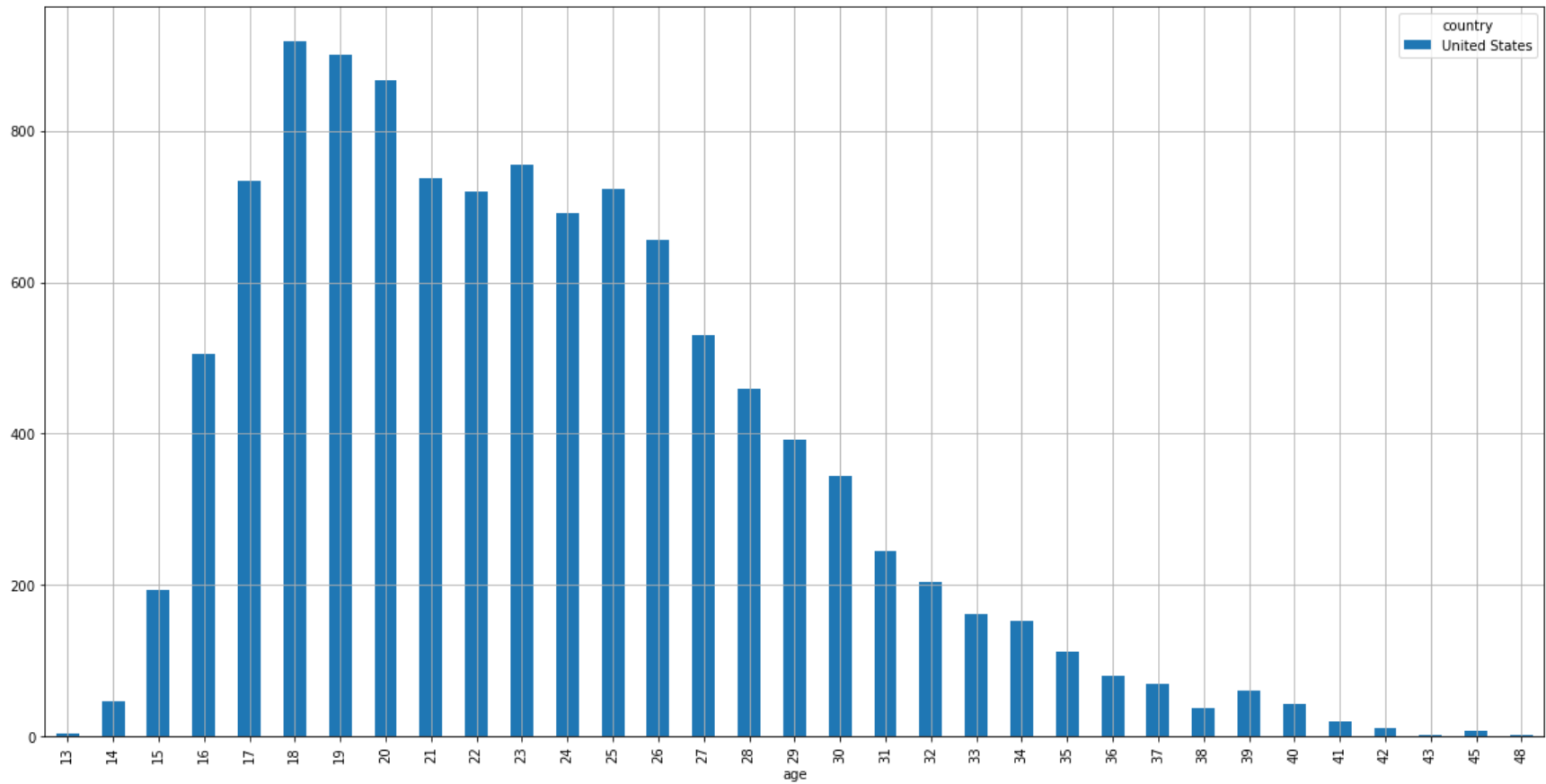
Age by Country - Participation rate
Event gender=Men



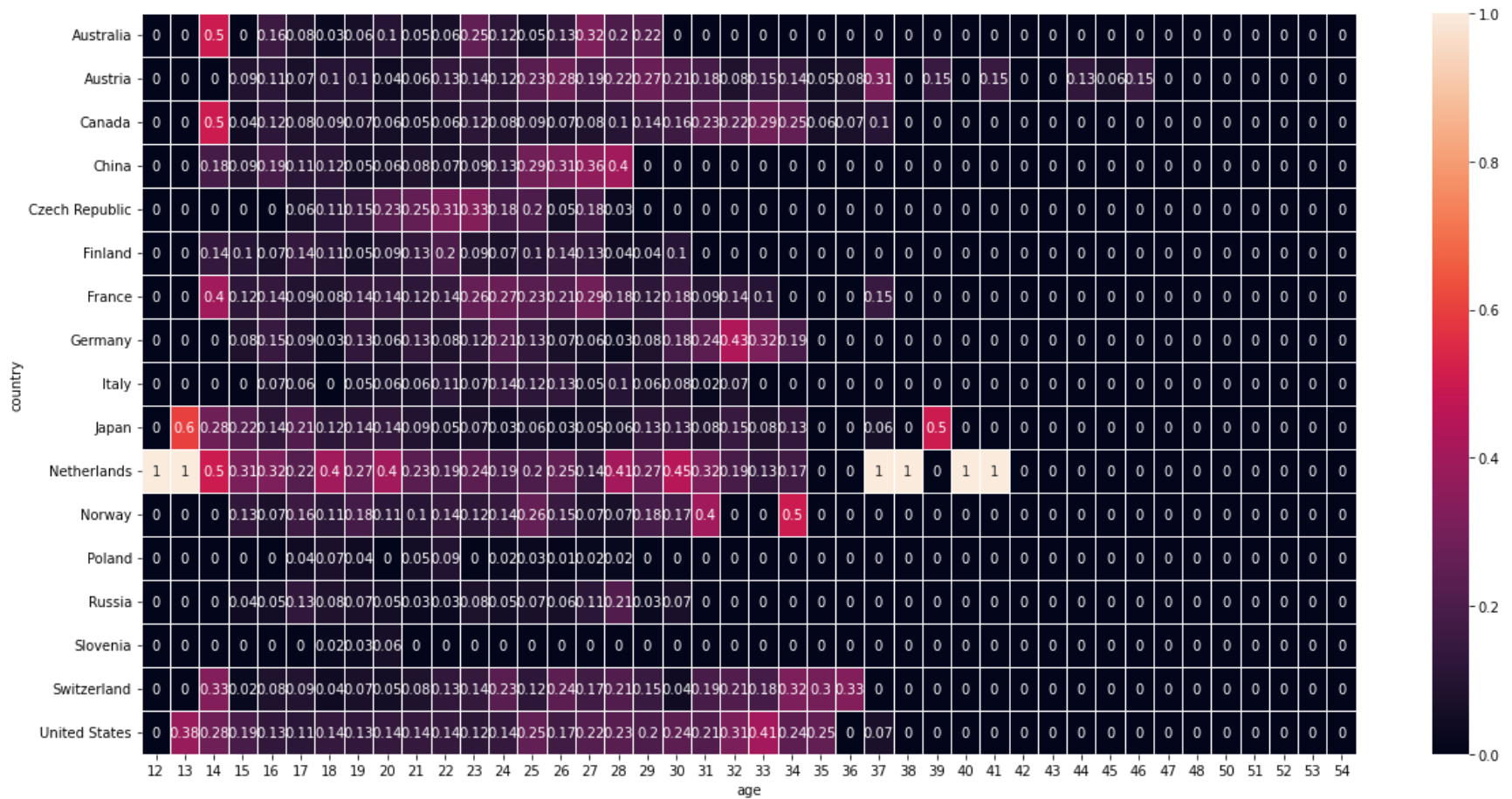
Age by Country - Medal rate
Event gender=Men - Country=United States



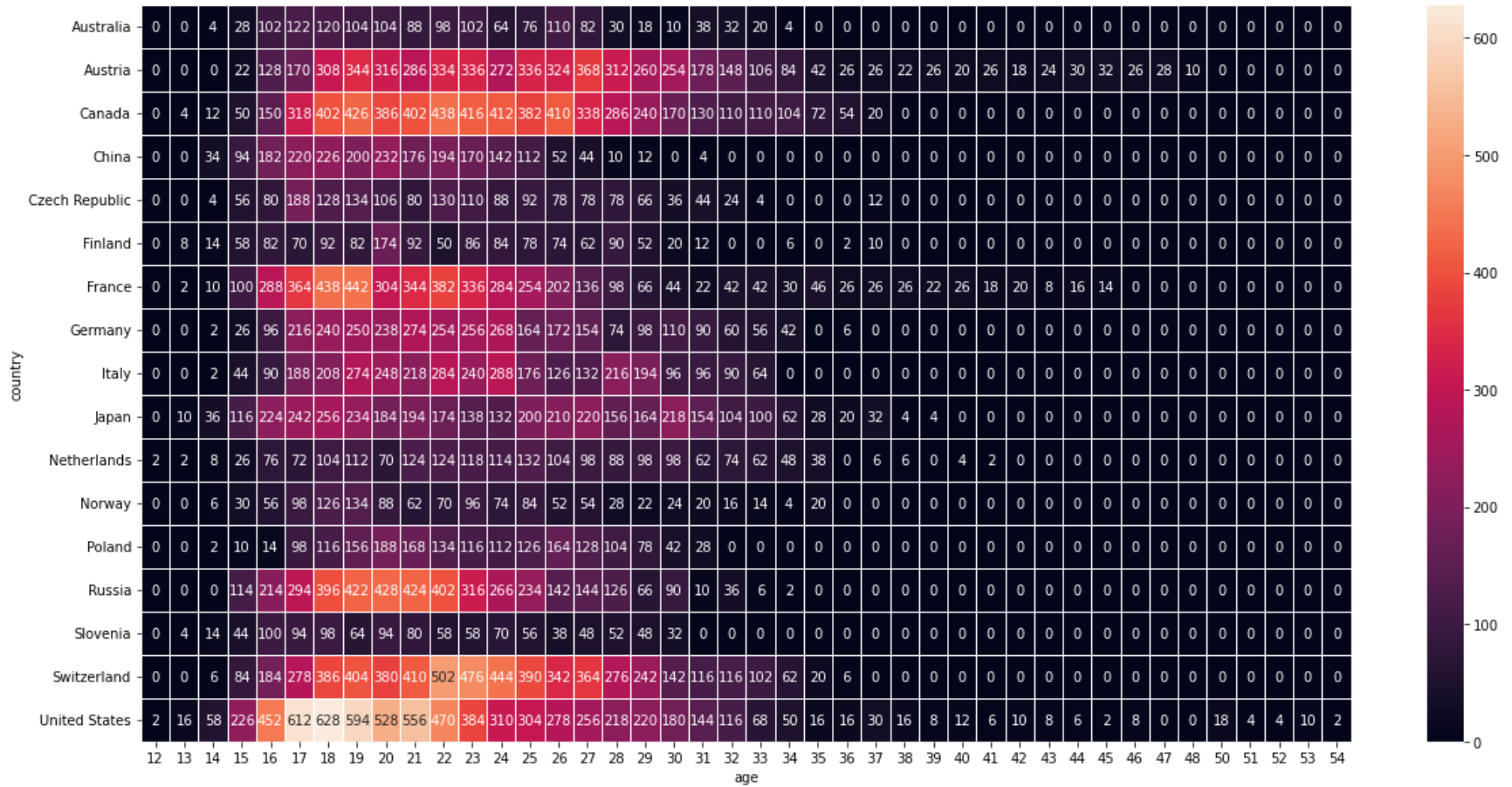
Age by Country - Participation count
Event gender=Men - Country=United States



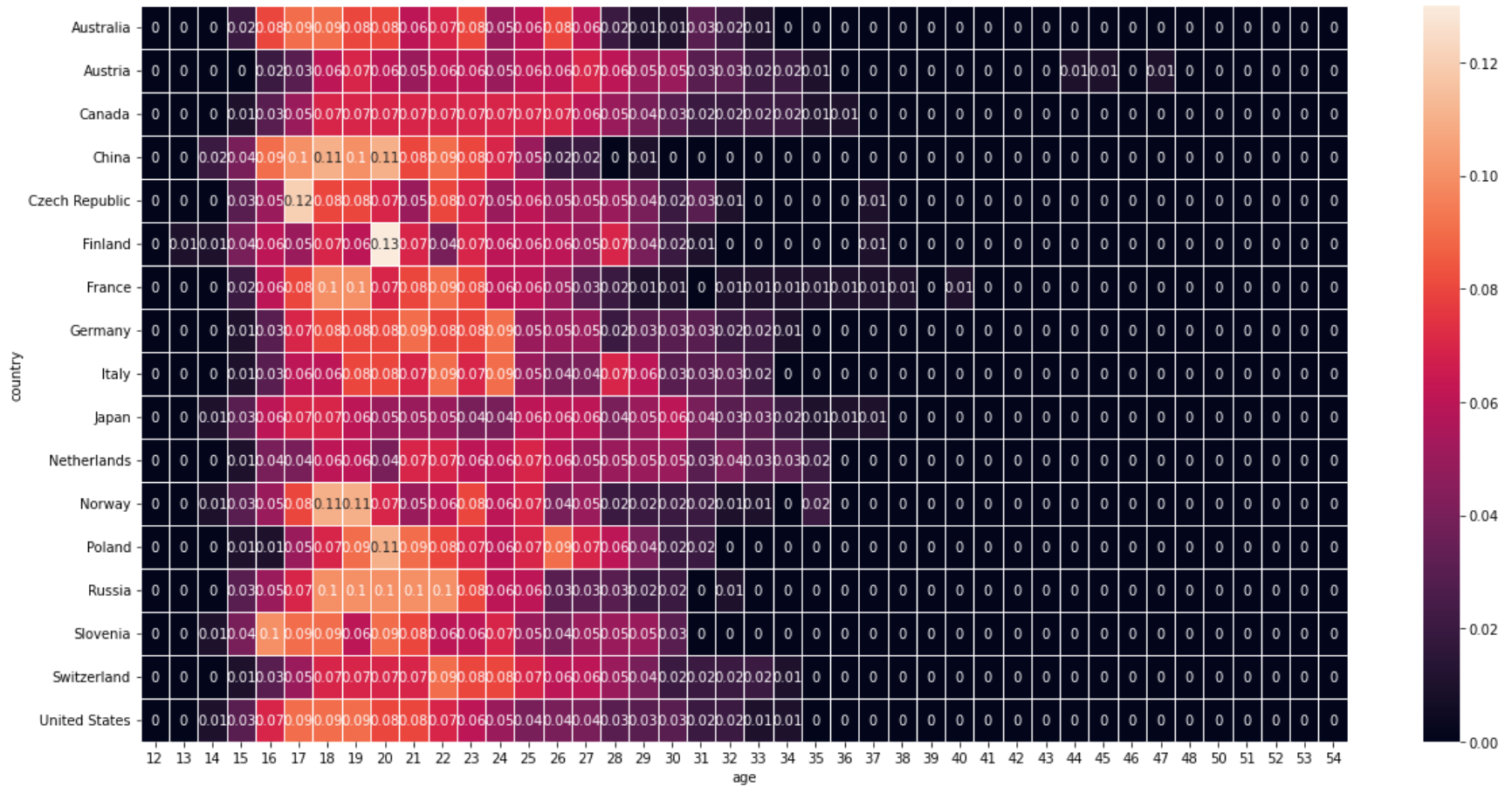
Age by Country - Medal rate
Event gender=Women



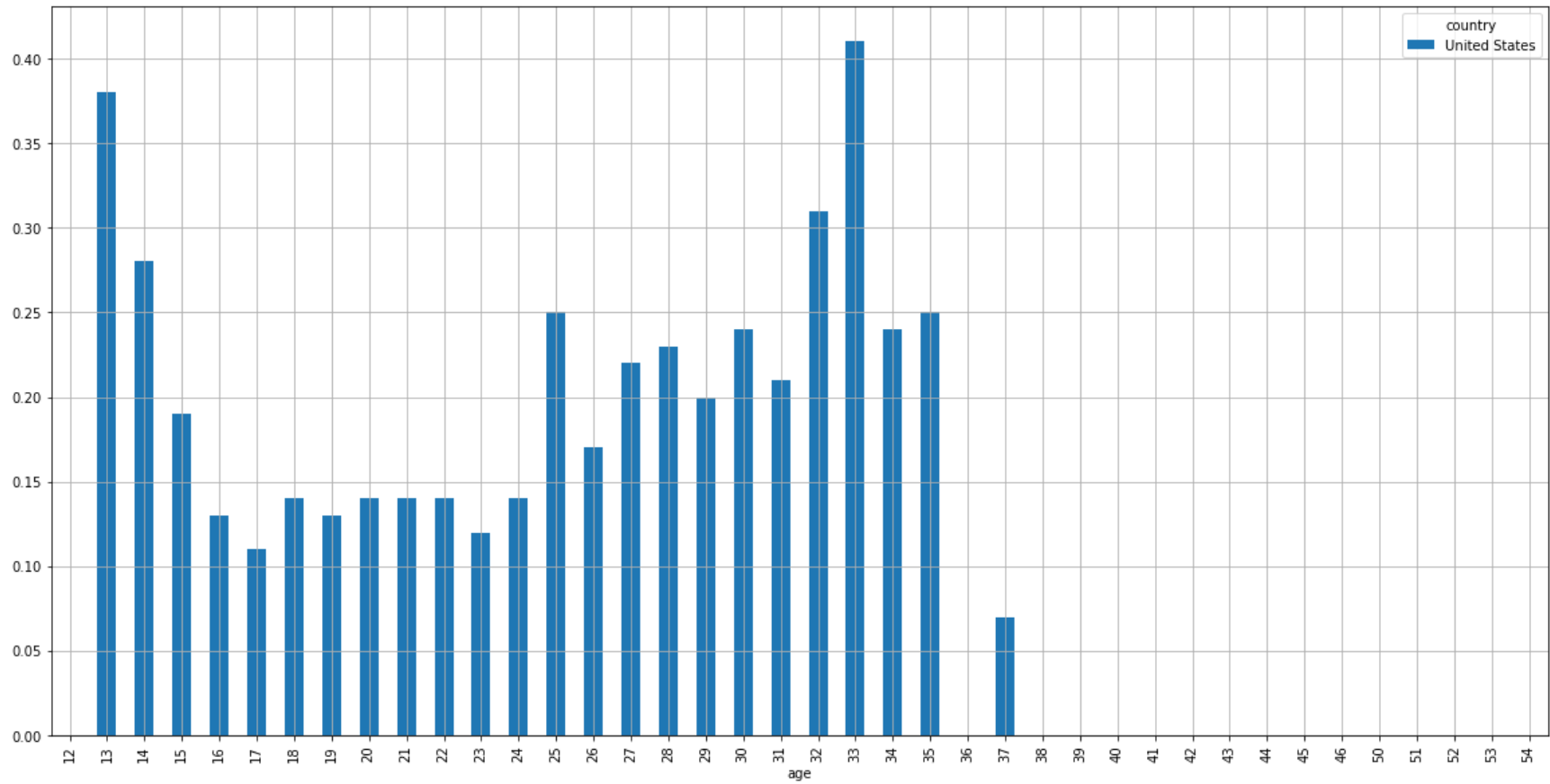
Age by Country - Participation count
Event gender=Women



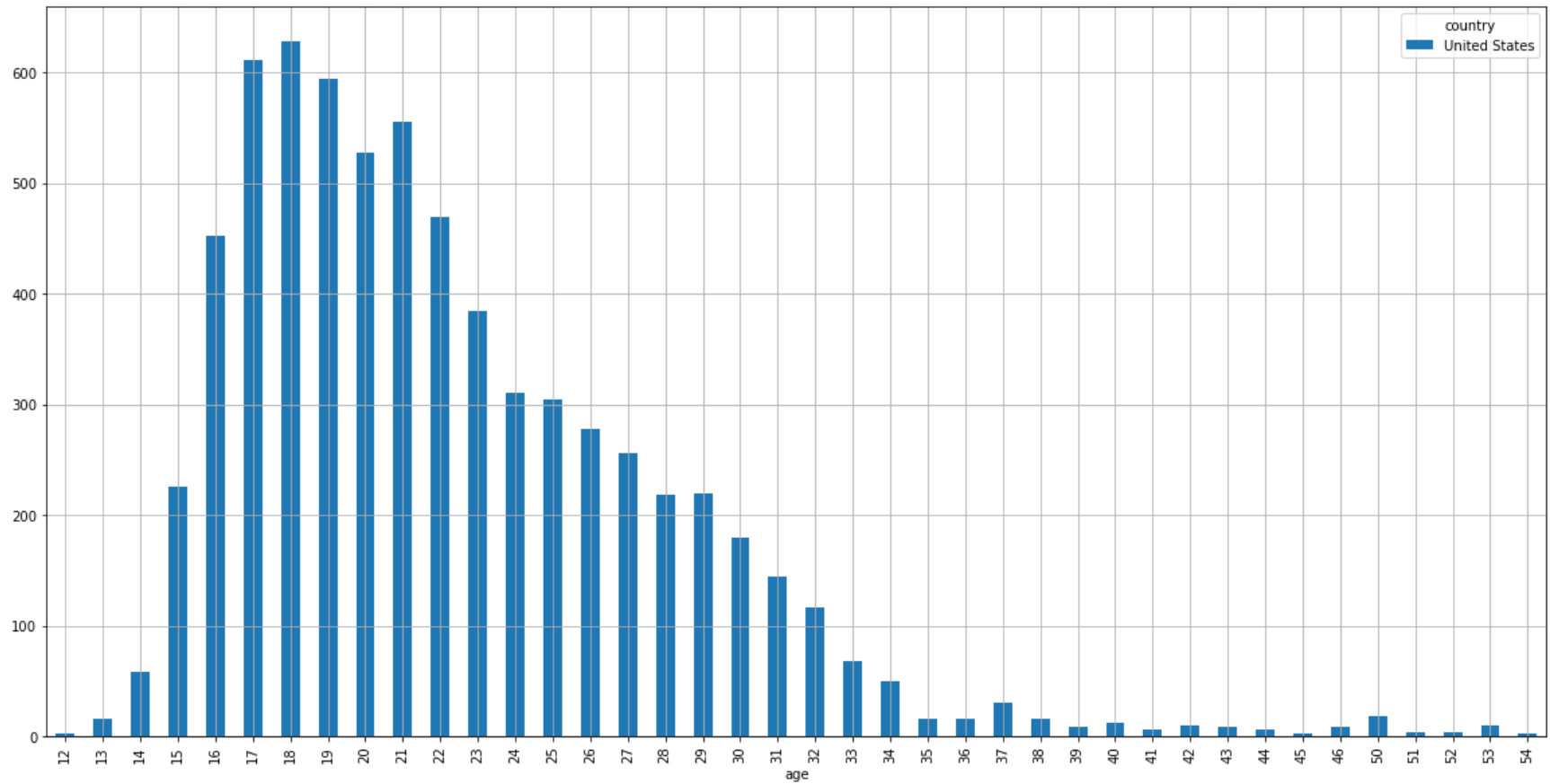
Age by Country - Participation rate
Event gender=Women



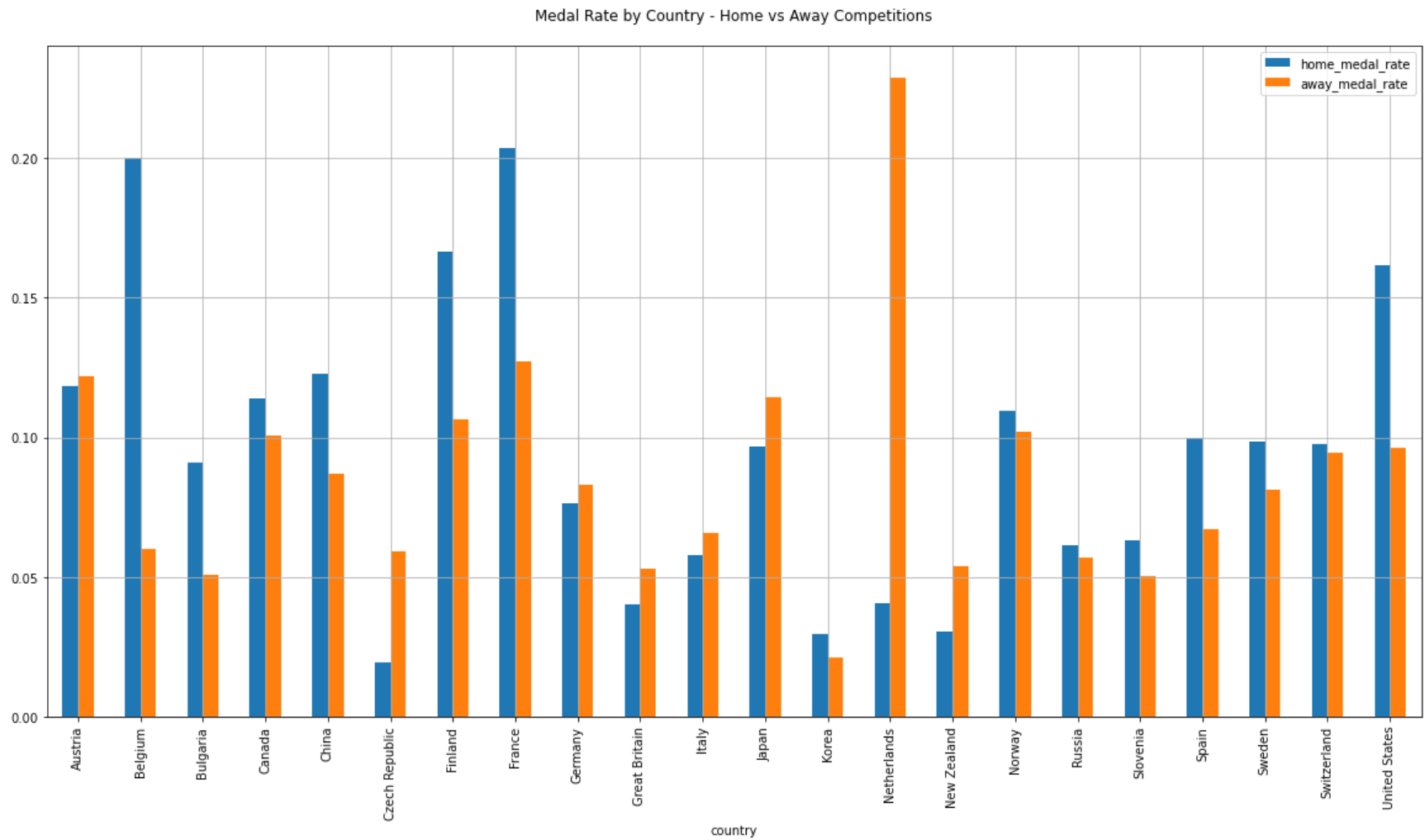
Age by Country - Medal rate
Event gender=Women - Country=United States



Age by Country - Participation count
Event gender=Women - Country=United States



Analysis of Success by Country for Home vs Away Competition



Observations and Conclusions

By observing the Country by Sport charts we can see the following patterns:

- There are wide variations in terms of event focus by country of origin as illustrated by the participation count by sport graph.
- The most popular sports by participation rate are Snowboardcross, Slopestyle, Parallel Giant Slalom, Halfpipe, Parallel Slalom and Big Air.
- Big Air has a significantly higher participation rate in Men and Women.
- Some countries show higher medal rates than others.
- The Netherlands appears to be the most successful using medal rate as the metric.
- The variation in country medal rate indicates the significance of country of origin as a factor that predicts success.
- Based on this analysis, countries and sporting events will be factored in when making predictions.

The patterns that arise from the Country by Age group statistics are the following:

- Peak ages, with highest medal rates show considerable variance with no explicit peak ages overall. This is an indicator that there is no clear peak age for any sport and event.
- Very high medal winning rates are observed by the Netherlands in higher ages.
- High medal rates appear to be more frequent in age ranges with lower participation counts, indicating potentially smaller competition.
- The highest participation rate by country follows a similar distribution across countries, but is centered around ages 19-20. This means most participants for this sport are around that age.
- The participation rate for women has a wider variance compared to the one for men and ages are a little bit more evenly distributed.

Finally, going over the home competition success plot, the following inferences can be made:

- Most countries are more successful when their athletes are competing at home versus when they are competing at a competition held in a foreign country.
- Netherlands exhibits an unusually high away winning rate but is generally an exception. This is likely due to the limited number of competitions held at home ground and the relatively small number of athletes the country has compared to other nations.
- For the United States in particular home competitions are much more successful.
- This is an indicator that a feature capturing whether an upcoming competition is being held at home versus away can provide valuable information when predicting the probability of winning a medal.

This concludes the exploratory data analysis portion of the report. In the next section we are going to discuss the Machine Learning methods used and the results of their predictive power.

Machine Learning for Predictions

As discussed earlier, for the machine learning phase we are going to treat the events in our dataset in two different ways. In the first approach we are going to utilize every event independently and in the second one we are going to leverage previous athlete performance by combining the events and forming a time series dataset where every row will represent the athlete's outcome in one event and the features of the three events that precede it.

In both cases we are going to predict the rank and the won_medal variable by using regression and classification respectively. The train/test split will be 80% and 20% respectively and we will use 5-fold Cross-Validation for hyperparameter tuning.

Once again, for the purpose of this study we are using the Snowboarding dataset and the event of Snowboardcross. The same process can be applied to any sport and event.

Classification

For the classification problem we developed two machine learning models and trained them on both the independent event dataset and the time series dataset and used the won_medal variable as our response variable. Rank of the final event was dropped from the dataset as it directly corresponds to the medal. The two models are a Random Forest Classifier and a Gradient Boosting Classifier. The performance results from both models are summarized below.

Independent Event Data

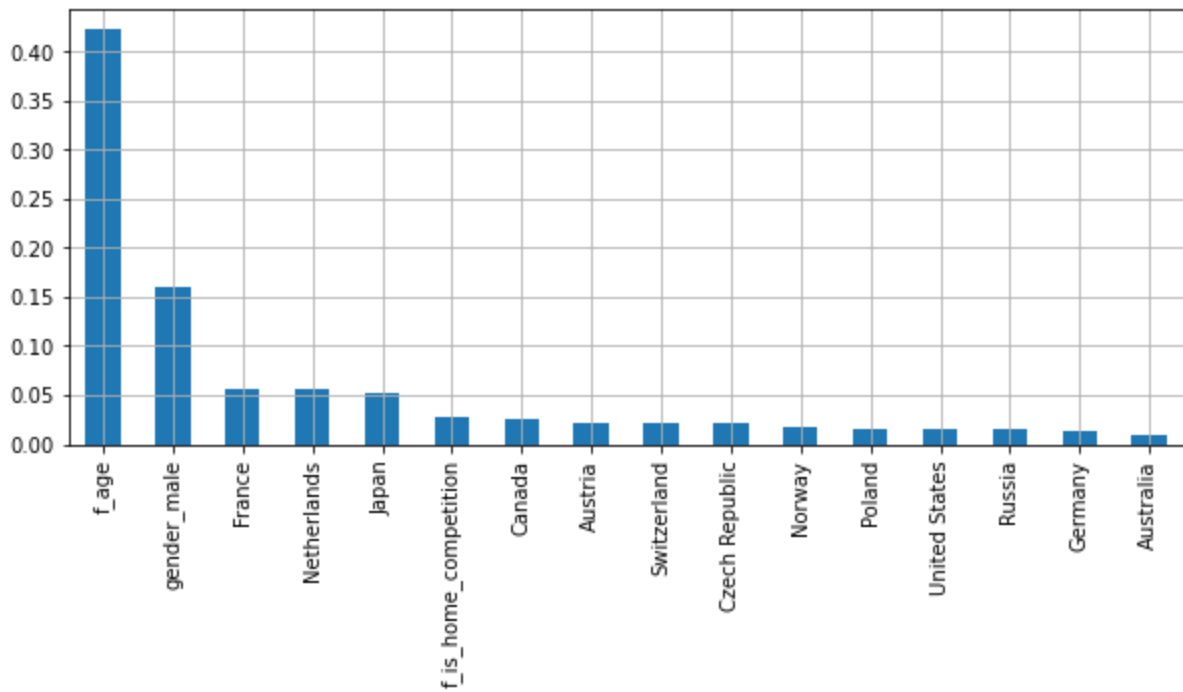
Random Forest Classifier - Independent Event Data					
Training Set Results	Medal Won	Precision	Recall	F1-Score	No. Datapoints
	False - 0	0.96	0.74	0.83	12689
	True - 1	0.18	0.66	0.29	1127
	Accuracy	-	-	0.73	13816
	Macro Average	0.57	0.70	0.56	13816
	Weighted Average	0.90	0.73	0.79	13816
Test Set Results	Medal Won	Precision	Recall	F1-Score	No. Datapoints
	False - 0	0.95	0.71	0.81	3152
	True - 1	0.16	0.57	0.25	302
	Accuracy	-	-	0.70	3454
	Macro Average	0.55	0.64	0.53	3454
	Weighted Average	0.88	0.70	0.76	3454

Training Set		Predicted N	Predicted P
	Actual Negative	9353	3336
	Actual Positive	379	748
Test Set		Predicted N	Predicted P
	Actual Negative	2253	899
	Actual Positive	130	172

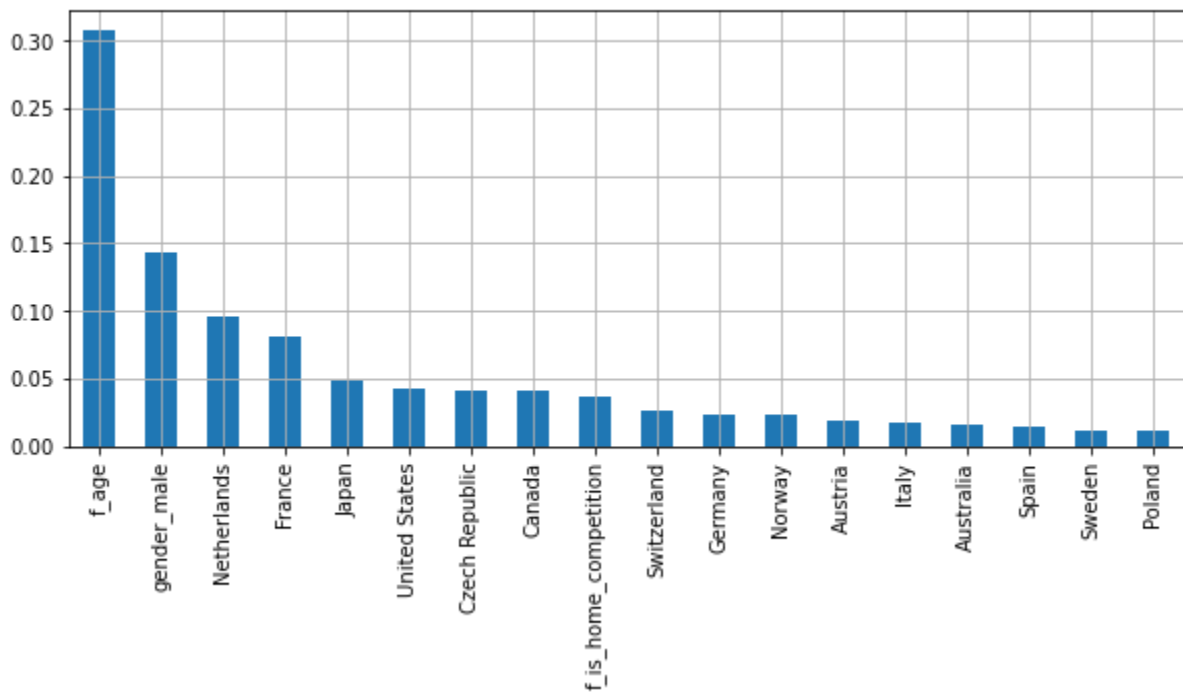
Gradient Boosting Classifier - Independent Event Data					
Training Set Results	Medal Won	Precision	Recall	F1-Score	No. Datapoints
	False - 0	0.96	0.71	0.82	12689
	True - 1	0.17	0.67	0.27	1127
	Accuracy	-	-	0.71	13816
	Macro Average	0.57	0.69	0.55	13816
	Weighted Average	0.90	0.71	0.77	13816
Test Set Results	Medal Won	Precision	Recall	F1-Score	No. Datapoints
	False - 0	0.95	0.70	0.80	3152
	True - 1	0.16	0.59	0.25	302
	Accuracy	-	-	0.69	3454
	Macro Average	0.55	0.64	0.53	3454
	Weighted Average	0.88	0.69	0.76	3454

Training Set		Predicted N	Predicted P
	Actual Negative	9063	3626
	Actual Positive	373	754
Test Set		Predicted N	Predicted P
	Actual Negative	2205	947
	Actual Positive	124	178

Random Forest Classifier - Independent Event Data - Feature Importances



Gradient Boosting Classifier - Independent Event Data - Feature Importances



Time Series Data

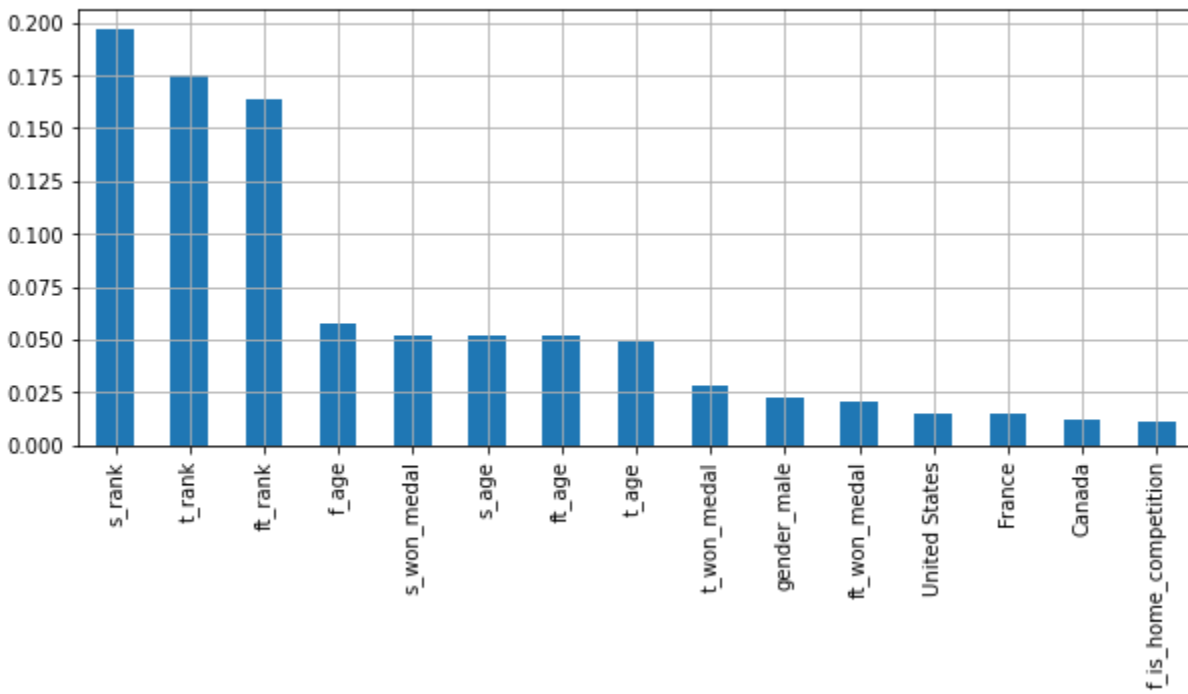
Random Forest Classifier - Time Series Data					
Training Set Results	Medal Won	Precision	Recall	F1-Score	No. Datapoints
	False - 0	0.99	0.97	0.98	11431
	True - 1	0.75	0.94	0.83	1097
	Accuracy	-	-	0.97	12528
	Macro Average	0.87	0.96	0.91	12528
	Weighted Average	0.97	0.97	0.97	12528
Test Set Results	Medal Won	Precision	Recall	F1-Score	No. Datapoints
	False - 0	0.95	0.95	0.95	2871
	True - 1	0.42	0.42	0.42	262
	Accuracy	-	-	0.90	3133
	Macro Average	0.68	0.68	0.68	3133
	Weighted Average	0.90	0.90	0.90	3133

Training Set		Predicted N	Predicted P
	Actual Negative	11085	346
	Actual Positive	64	1033
Test Set		Predicted N	Predicted P
	Actual Negative	2720	151
	Actual Positive	152	110

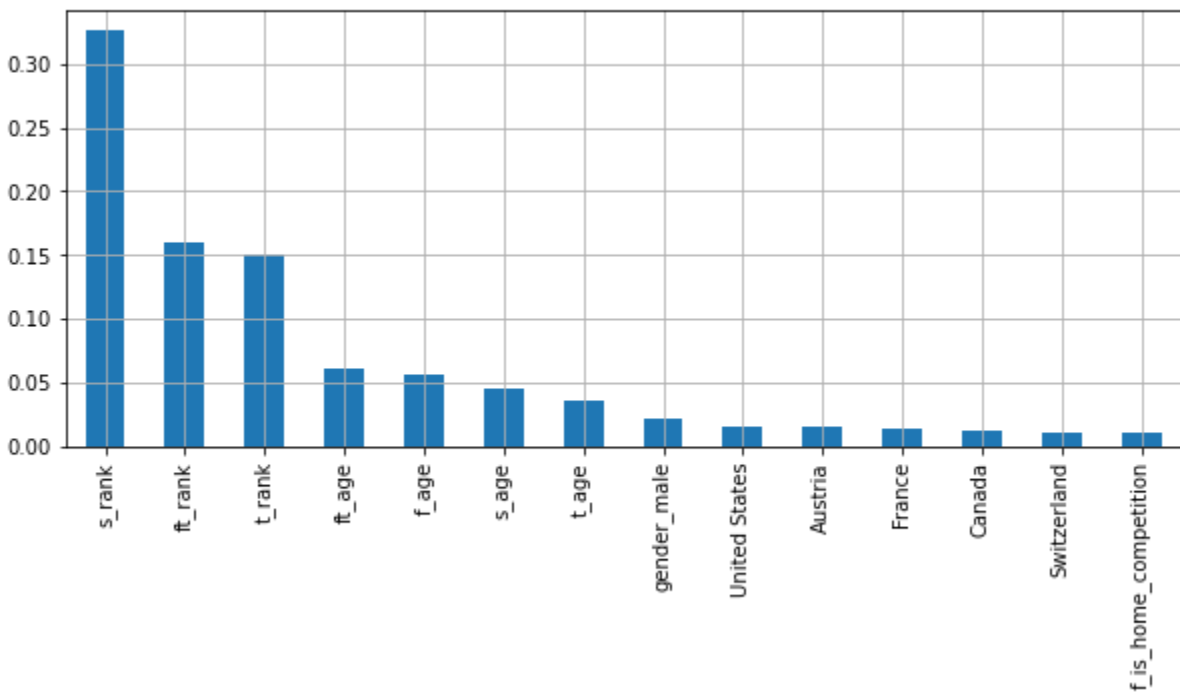
Gradient Boosting Classifier - Time Series Data					
Training Set Results	Medal Won	Precision	Recall	F1-Score	No. Datapoints
	False - 0	1.00	1.00	1.00	11431
	True - 1	0.97	1.00	0.98	1097
	Accuracy	-	-	1.00	12528
	Macro Average	0.98	1.00	0.99	12528
	Weighted Average	1.00	1.00	1.00	12528
Test Set Results	Medal Won	Precision	Recall	F1-Score	No. Datapoints
	False - 0	0.94	0.96	0.95	2871
	True - 1	0.39	0.27	0.32	262
	Accuracy	-	-	0.90	3133
	Macro Average	0.66	0.62	0.63	3133
	Weighted Average	0.89	0.90	0.90	3133

Training Set		Predicted N	Predicted P
	Actual Negative	11397	34
	Actual Positive	0	1097
Test Set		Predicted N	Predicted P
	Actual Negative	2758	113
	Actual Positive	191	71

Random Forest Classifier - Time Series Data - Feature Importances



Gradient Boosting Classifier - Time Series Data - Feature Importances



Regression

Similar to the classification approach, two machine learning models were developed for regression. The two models are a Random Forest Regressor and a Gradient Boosting Regressor. Both models were trained using the independent event data and the time series data in two separate phases and used the rank of the final event as the response variable. Of course, won_medal has been dropped from the dataset as it predicts the rank. The results from the two models are summarized below.

Independent Event Data

Random Forest Regressor		
	Mean Absolute Error	Mean Squared Error
Training Set Results	12.12	246.65
Testing Set Results	12.72	274.30

Gradient Boosting Regressor		
	Mean Absolute Error	Mean Squared Error
Training Set Results	11.81	237.68
Testing Set Results	12.74	278.71

Time Series Data

Random Forest Regressor		
	Mean Absolute Error	Mean Squared Error
Training Set Results	8.31	118.15
Testing Set Results	9.72	168.92

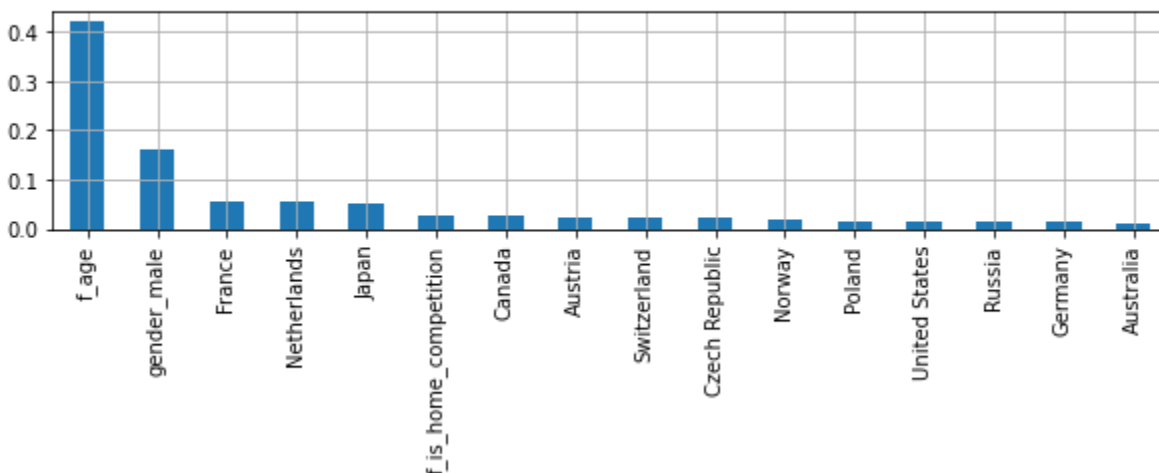
Gradient Boosting Regressor		
	Mean Absolute Error	Mean Squared Error
Training Set Results	9.53	141.57
Testing Set Results	10.78	188.68

Observations and Conclusions

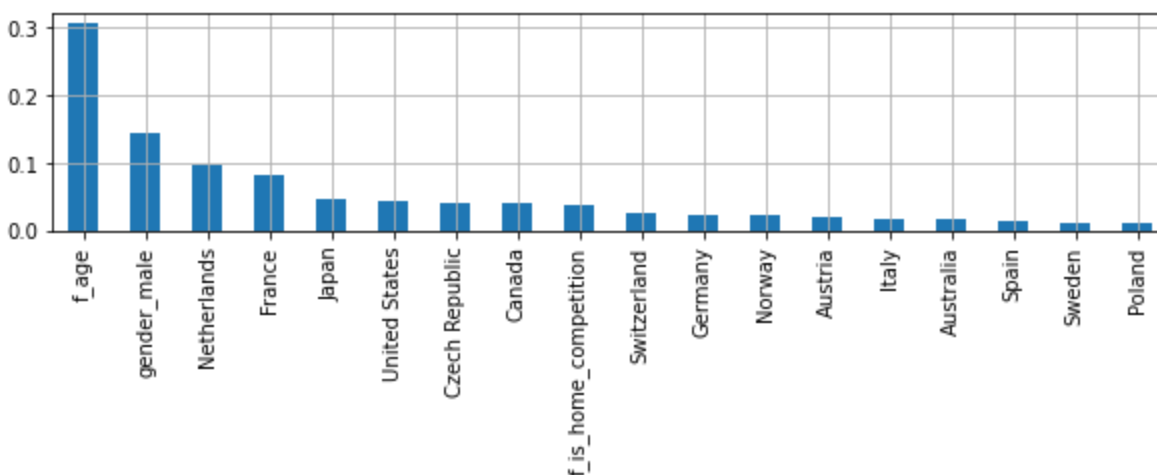
The analysis and model results showed that the Random Forest algorithm has outperformed the Gradient Boosting in both the regression and classification approaches. However the performance was significantly better when using the time series dataset as compared to the independent event dataset.

After training the model with the independent event data, both the regression and classification models showed relatively poor results. The best performing classifier was the Random Forest. The best precision score on the testing set after hyperparameter tuning was only 0.16 and the Mean Absolute Error from the regressor on the testing set was 12.72. The low performance in both the training and the testing set shows that the model is suffering from high bias. The feature importance graph illustrates that age is an important factor, followed by countries of origin that athletes are overall more successful. The high bias problem indicates that this dataset is overly simplistic and does not successfully capture the complexity of the dataset.

Random Forest Classifier - Independent Event Data - Feature Importances

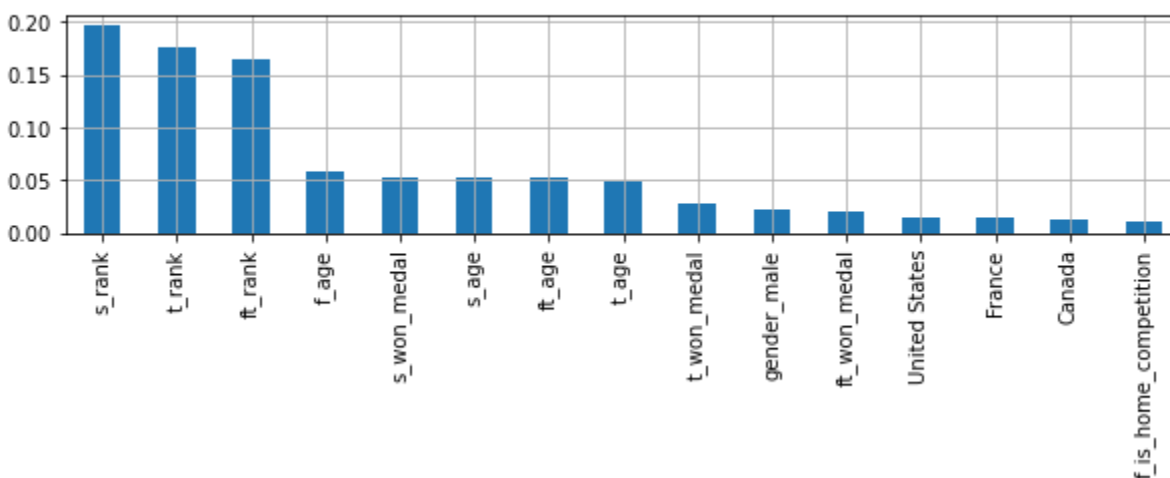


Gradient Boosting Classifier - Independent Event Data - Feature Importances

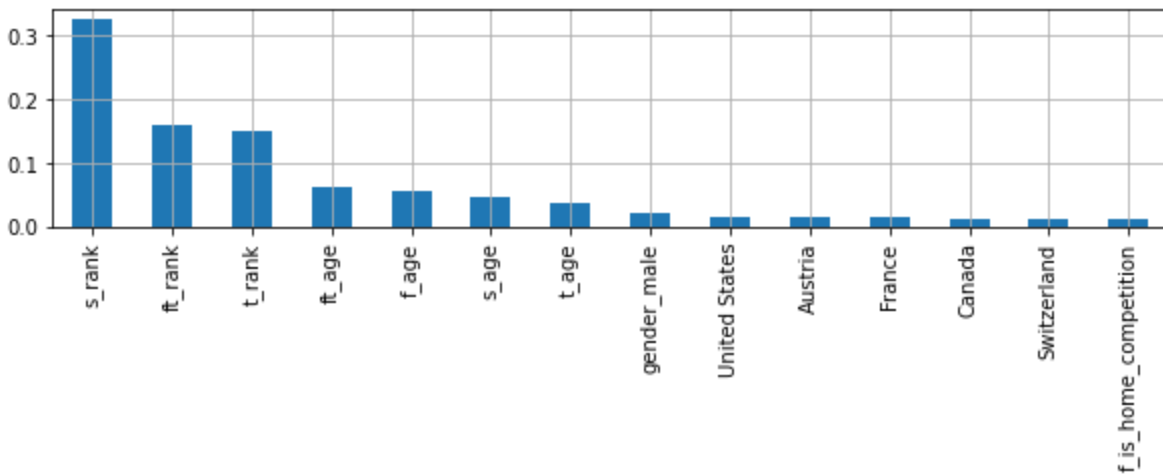


After training the model with the time series data, both the regression and classification models performed significantly better. The best performing classifier showed precision of 0.42 and recall of 0.42 in the testing set after hyperparameter tuning, which effectively translates to roughly a 42% chance of accurately predicting that an athlete will win a medal, whereas for non-medal winning athletes, both precision and recall is 0.95. The regressor results showed a Mean Absolute Error of 9.72 in accurately predicting the rank of the athlete. The feature importance analysis shows that performance and age in the preceding events have a very powerful effect in predicting the outcome, with country of origin for the most successful countries and the home competition variable coming after these.

Random Forest Classifier - Time Series Data - Feature Importances



Gradient Boosting Classifier - Time Series Data - Feature Importances



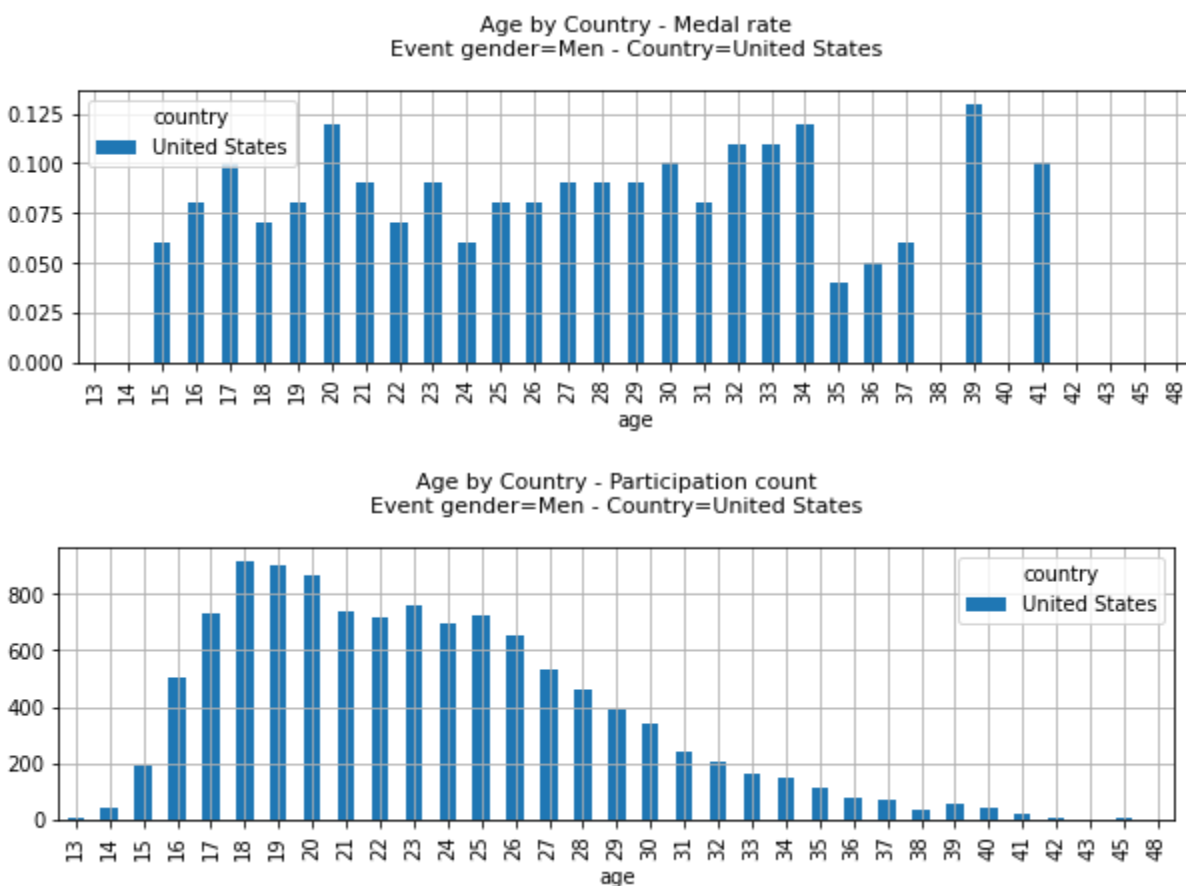
Overall, the Random Forest Classifier with the time series data, is the model of choice that provides the best predictive accuracy for the dataset of Snowboarding. The approach can be expanded in other sports as well. The predictions rely primarily on the performance and age of the athlete in the three events preceding the event we are trying to predict and the outcome is the probability of the athlete winning a medal. Alternatively, the Random Forest Regressor can be used in a similar way to predict the rank of the athlete for the event in question.

Key Questions and Answers

- Can a country's percentage of athletes competing at an Olympic Games while in their peak age ranges in their respective sports predict medal success?

Although medal success shows some higher rates in certain age groups by country, the difference is marginal and not significant enough to predict success. Thus, age cannot single-handedly be used to predict success and there is no clear cut peak age for medal success.

- Which athletes who competed in Tokyo 2021 (Beijing 2022) will be hitting their peak age in their respective sports in Paris 2024 (Milan 2026)?



Taking the figure above as a reference point for the United States with Snowboarding as the sport of reference, there is no specific peak age for athletes that are more successful. This is in line with previous conclusions of the analysis that even though there is a slightly higher success rate in some age ranges, there is no distinct “peak age” that athletes are more successful. Thus, with age alone as a factor it’s nearly impossible to accurately predict success in the upcoming olympics based on age and success in previous ones.

- Can the number of young, “pre-peak” athletes that competed in Tokyo predict medal success for their countries in Paris 2024 and Milan 2026?

Using the classifiers or regressors built during this project it is very feasible to predict success in the upcoming Olympics. However, prediction is not limited to Tokyo only (though this is an option), but other preceding events and athlete performance during those can be used as well.

- How many promising athletes who did not compete in Tokyo (and Beijing) but will reach peak age in Paris (and Milan) do Team USA and other countries have in their pipelines?

Similarly to the previous question, any events that the athlete previously participated in along with their age and performance during those events, can be used with the models developed to predict their probability of success in the upcoming Olympics.

Final Outcome and Recommendations for future work

As a final outcome of this study, the time series approach is a robust and sustainable methodology for predicting the success of an athlete not only in Olympic Game events but also in World Championship events. The only requirement of the present state of the model is the age and rank of an athlete in three preceding events, as well as the country of origin of the athlete, the gender and the location of the competition in question to determine whether it takes place in the athlete's home country. This information can successfully predict medal winners with 42% and non-winners with 95% precision and recall.

Further enhancements can be achieved by accounting for biometric data of the athlete, as intuitively this information can provide extra predictive power.