# Assessment 2

UTS
UNIVERSITY OF TECHNOLOGY SYDNEY

32146 - Data Visualization and Visual Analytics

## Australian Open Analysis

Somesh Shanbhag

Student ID:
25525837

Somesh
Shanbhag

25525837

AO
australian open

# Table of Contents

# 1. Introduction

The Australian Open tennis finals dataset is a rich historical record that spans more than a century capturing the evolution of men's and women's tennis from 1905 to 2024. Each row in the dataset represents the final match of the tournament in a given year and includes detailed fields such as the year, gender, champion and runner -up names along with their nationalities and seed numbers. Match scores, match durations in minutes and various game win/loss metrics along with win ratios and first-to-win ratios. This report recognizes the importance of understanding the evolution of tennis by analyzing historical data trends. This dataset not only provides insights into the performance of individual champions over time but also gives an overview of how the characteristics of the match outcomes, nationality shifts and even gender representation have evolved with sport's global development. In this report, we will discuss the format and structure of the dataset, explain the data transformations and calculations and propose a suite of visual analytics graph using only tree maps, parallel coordinates, geographic maps and scatter charts. Our visualization will be created in tableau and excel with detailed instructions, annotations and recommendations to ensure clarity, effective storytelling and data accuracy.

# 2. Dataset Summary

## 2.1 Data Structure and Format

The dataset is structured as a tabular file with each row representing one finals match and each column capturing a specific attribute of that match. The attributes include:

- Year: Indicates the year the match was played.
- Gender: Specifies whether the match was in the women's or Men's tournament.
- Champion: the name of the winning player
- Champion nationality & Champion country: provide information on the nationality and country of origin of the champion.
- Score: The match score detailing how the champion won the match.
- Champion Seed: The seed number of the champion indicating ther ranking based on pre-tournament analysis
- Mins: Duration of the match in minutes.
- Set Results: Fields like "1st-won", "1st-loss", "2nd-won", "2nd-loss", etc. detail performance in individual sets.
- Runner-up Information: This includes the running-up's name, nationality, country and seed.
- Games won, Games lost, Win Ratio and Set Win Ratios: These columns capture further quantifiable performance metrics calculated for each match.

- Upset Measurement and Calculations: Match competitiveness measures how close the difference between games won and games lost is. If it's big, then the champion dominated the game but if it is small then it was an extremely competitive match. Match Upset measures if the finals were won by a champion that was seeded lower than the Runner Up and Upset Level Difference measures are a conditional on Match upset being 1 and is Used to measure the difference between the seeded levels of the winner and the runner up. Unseeded players were given the level 33 as there are only 32 seeded players in the Australian Open. This was done so that Tableau is able to do the calculations with a numerical value.

## 2.2 Data Characteristics

- Temporal Dimensions: The data spans from the earliest records starting in 1905 up to the most recent matches in 2025. This wide temporal range allows us to examine long-term trends such as changes in match duration, scrolling patterns and competitive balance.
- Gender Segregation: Separate rows exist for men's and women's finals which allows for gender based comparisons of game dynamics.
- Performance metrics: The dataset includes both categorical (names, nationalities, gender) and numerical variables ( score details , games win/loss, win ratios). These mixed data types support a variety of statistical and comparative analysis.
- Scoring data: The score field is a rich textual field that may require further parsing to extract set-level results and to calculate additional metrics such as average games per set
- Seed Information: This data can be correlated with match outcomes to explore possible relationships between pre-tournament expectations (seed value) and ultimate match results.

## 2.3 Observed Trends and Outliers

In preliminary exploration of the dataset, there are several trends and outliers that have emerged:

- Dominance Patterns: Some Champions like Novak Djokovic and Serena Williams appear repeatedly across different years indicating periods of dominance.
- National Representation: there is a noticeable shift from predominately local and regional players in the early years with many Australian players to a more international pool in the modern Era.
- Gender Representation: From the data its clear that in the early there were no female competitions taking place. This changed from 1922 onwards and has

expanded from just Australian female participants to international female participants.

- Outliers in Win Ratios: Some matches feature extraordinary win ratios or set wins percentages (E.g. near perfect set wins) hat merit further investigation. These outliers correspond to either particularly dominant performance or in some cases matches with atypical conditions such as retirement.

## 3. Data Transformation and Calculation

Before creating visualizations, it is essential to prepare and transform the data to enrich the analysis. Some of the key transformations and calculations include:

### 3.1 Parsing Score Data

The raw score entries such as 6-3, 7-6(7-4) etc. was already parsed initially to extract

- Set-By-Set Scores: This involves splitting the score string using commas and further breaking it into each set and then into individua game wins and losses
- Determination of Match types: Some matches were best of three sets while others were best-of-five. This differentiation is crucial in computing average performances across sets.
- Unabbreviated Country names: Since Tableau cannot associate countries with their abbreviations a separate column was created with the champions nationality clearly explained by the country.

### 3.2 Calculation of Additional Metrics

Several new metrics derived from the available data to support meaningful comparisons:

- Total Games per Match: Calculated by summing all games won and lost by the winning player.
- Win Ratio: Win Ratios were calculated based on the total games played per match which tell us how efficient and dominant the winning player was in the match. It also tells us more about their adaptability, match stamina and competitive edge.
- Champion versus Runner up Differential: The difference in games won between champions and runner up provides insight into match competitiveness.

### 3.3 Data Aggregation by Nationality and Gender

- Aggregated Win Ratios: Average win ratios were computed for each nation.
- Frequency of appearance: The number of finals appearances by nationality and gender was calculated to identify dominant periods and countries.

- Trend lines by Era: Data was segmented into time periods to observe how the competitive landscape evolved over time.

## 3.4 Handling Missing and Anomalous Values

The dataset contains some cells with "#Div/0!" error for win ratios or missing set outcomes. These missing values were replaced with Nulls by Tableau to allow for the creation of continuous, meaningful charts without disruption. Extreme values were marked so that they could be investigated further in the visualization phase. Year dates for 1977 have been handled as Null as due to a scheduling change the tournament was held twice that year. Most of the records are missing match duration data due to which no analysis can be done using this field and this field has been ignored.

## 4. Data Analysis and Visualization

In this section, the analysis is presented through four required visualizations: a tree map, parallel coordinate lots, geographic maps and scatter charts. Detailed explanations are provided on how these visualizations were implemented in the tableau along with an explanation of the insights they reveal. Each visualization helps answer specific analytical questions on player performance, trends by nationality, gender-based progressions and changes over time.

## 4.1 Treemap Visualisation

### 4.1.1 Purpose

A Treemap is a space efficient hierarchical visualization that displays proportions in a nested format. In this report, it has been used to visualize the frequency of tournament wins by different countries and Top players wins (which is measures as 5 or more tournament wins).
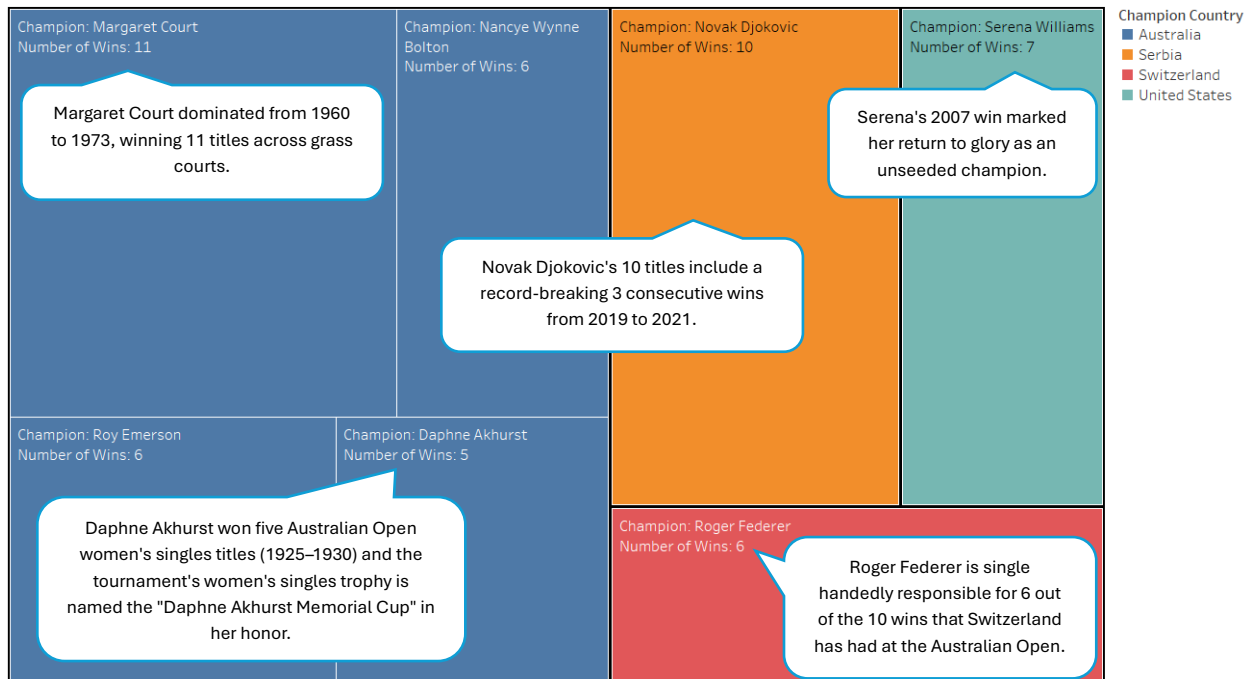
Top Players by Country

Champion: Margaret Court
Number of Wins: 11

Margaret Court dominated from 1960 to 1973, winning 11 titles across grass courts.

Champion: Nancye Wynne Bolton
Number of Wins: 6

Champion: Novak Djokovic
Number of Wins: 10

Champion: Serena Williams
Number of Wins: 7

Serena's 2007 win marked her return to glory as an unseeded champion.

Novak Djokovic's 10 titles include a record-breaking 3 consecutive wins from 2019 to 2021.

Champion: Roy Emerson
Number of Wins: 6

Champion: Daphne Akhurst
Number of Wins: 5

Daphne Akhurst won five Australian Open women's singles titles (1925–1930) and the tournament's women's singles trophy is named the "Daphne Akhurst Memorial Cup" in her honor.

Champion: Roger Federer
Number of Wins: 6

Roger Federer is single handedly responsible for 6 out of the 10 wins that Switzerland has had at the Australian Open.

Champion Country
- Australia
- Serbia
- Switzerland
- United States

*Figure 1: Top Player Wins with 5 or more Australian Open Wins*

Top Runner-Ups By Country

Runner Up: Esna Boyd
Number of Losses: 6
Their Average Seed: 2.000
Lost to Average Seed: 2.000

Runner Up: Jan Lehane
Number of Losses: 4
Their Average Seed: 3.250
Lost to Average Seed: 2.500

Runner Up: Andy Murray
Number of Losses: 5
Their Average Seed: 4.200
Lost to Average Seed: 1.400

Runner Up: Rafael Nadal
Number of Losses: 4
Their Average Seed: 3.500
Lost to Average Seed: 6.750

Esna Boyd, a trailblazer in Australian tennis, clinched one Australian Open title (1927) but also endured six runner-up finishes, showcasing her resilience and consistent presence in the finals.

Andy Murray, a five-time Australian Open finalist, has never won the title, often losing to higher seeds like Djokovic and Federer in fiercely contested matches.

Runner Up: Thelma Coyne Long
Number of Losses: 4
Their Average Seed: 2.750
Lost to Average Seed: 1.500

Despite injuries, hard court challenges, and rivals like Djokovic, Nadal persevered to claim the Australian Open title twice, in 2009 and 2022.

Runner Up: John Bromwich
Number of Losses: 5
Their Average Seed: 2.200
Lost to Average Seed: 2.400

John Bromwich, a seven-time Australian Open finalist, claimed two titles (1939, 1946) but faced defeat in five finals, underscoring his remarkable consistency in the tournament's history.

Runner Up: Chris Evert
Number of Losses: 4
Their Average Seed: 1.500
Lost to Average Seed: 2.000

Runner-up Country
- Australia
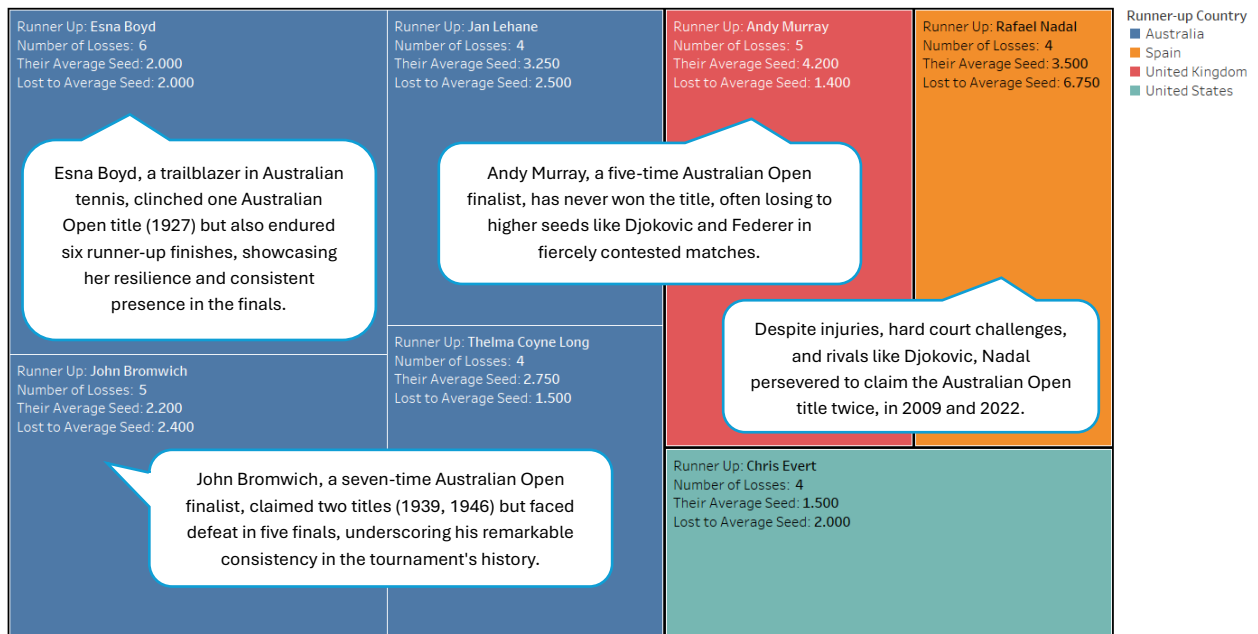- Spain
- United Kingdom
- United States

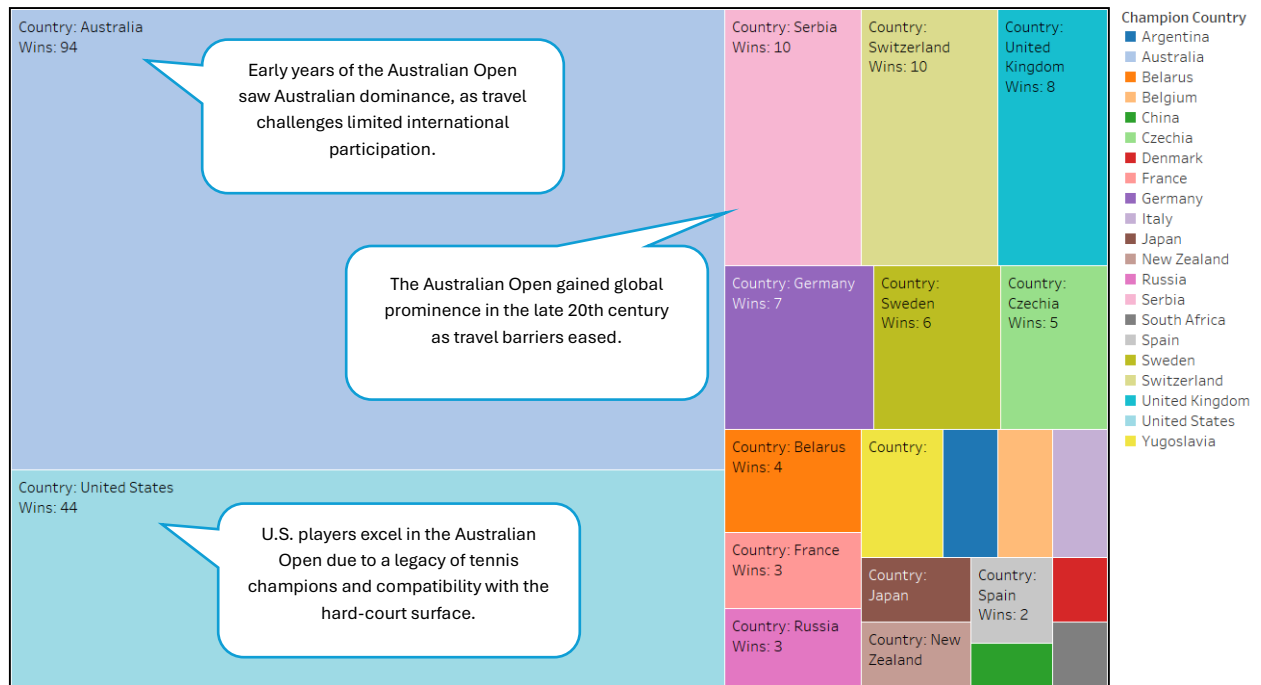*Figure 2: Top Runner Up by Country*

## Wins By Country



*Figure 3: Australian Open Wins by Country*
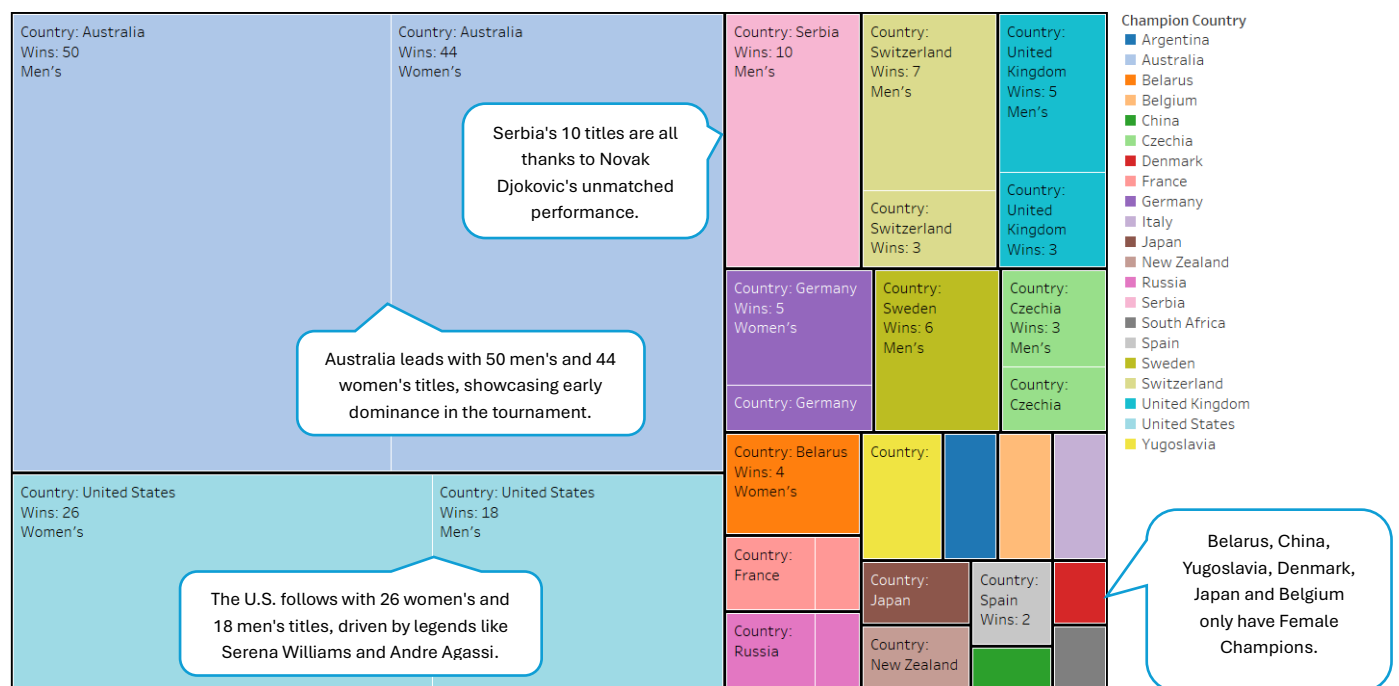
## Wins By Gender per Country



*Figure 4: Australian Open Win by Gender per Country*

### 4.1.2 Chart Description

For Top 5 player wins, A filter was added to only display those field that had an Australian open count higher than 5. Then champion and count of Australian open field were added to color and size respectively. For Top Runner-Ups, all champion stats were replaced with runner up stats. For Australian Open wins by Country, Champion country field is dragged to rows shelf and Number of records to the size mark. Champion country is also dragged to color to differentiate by color. The color palette was adjusted to ensure clear visibility. Labels were added clearly showing the country and its win count. The fourth graph was just created by adding gender to Detail to see the contribution of total wins by gender for each country. Relevant information has been annotated to give share more information regarding certain results as shown by the treemap. Annotations were added to give some context to the data.

### 4.1.3 Insights and Trends

The tree map highlights the concentration of champion appearances by country. For example, Australia overwhelmingly dominated the early 20th century which is evidenced by large blocks representing frequent wins. In contrast, more recent years show smaller, more fragmented blocks signifying a diversification in champion nationalities as the tournament became more global. Later decades saw a mix of champions from USA, Europe and other regions reducing dominance by any single nation. While winning this tournament once itself is an incredible achievement, The Best top Player with the highest record for Australian Open wins is held by Margaret Court at 11 wins in the female Category and Novak Djokovic with 10 wins in the male category. There are also many others such as Roy Emerson( 6 titles), Roger Federer (6 titles), Serena Williams (7 titles), Daphne Akhurst (5 titles) who have showed exceptional skill and talent being able to win this competition multiple times. From the Runner-Ups treemap, it's clear that Esna Boyd, Andy Murray and John Bromwich put up an amazing fight in each of their finals but lost many finals without winning. Out of these Andy Murray is the only player to be a 5-time finalist that has never won the Australian Open. By visualizing based on gender, the treemap can also show how female champions from countries like the USA and Belarus had risen to prominence in recent years. From the treemap, it's clear that some countries Belarus, China, Yugoslavia, Denmark, Japan and Belgium have excellent female tennis players who have won the Australian Open as opposed to their male counterparts.

## 4.2 Parallel Coordinate Plot

### 4.2.1 Purpose

Parallel coordinate plots are ideal for multivariate analysis. In our scenario, they are used to visualize a pattern across multiple performance metrics and to uncover any correlation across time periods or overall match outcomes.
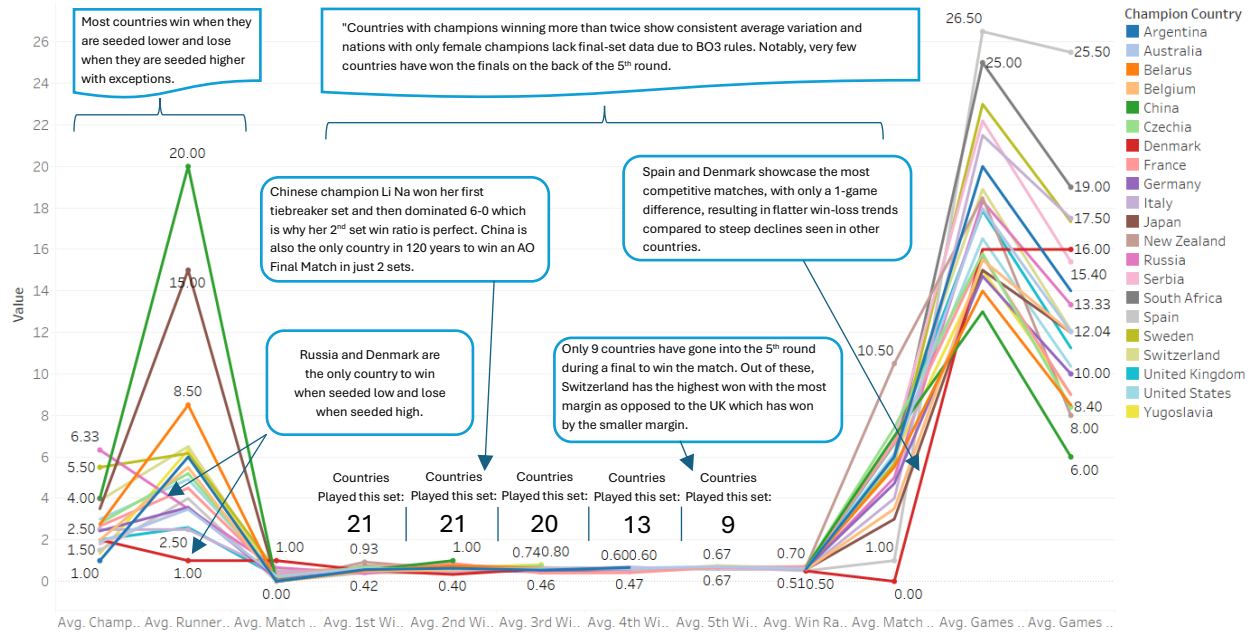
## Performance by Country



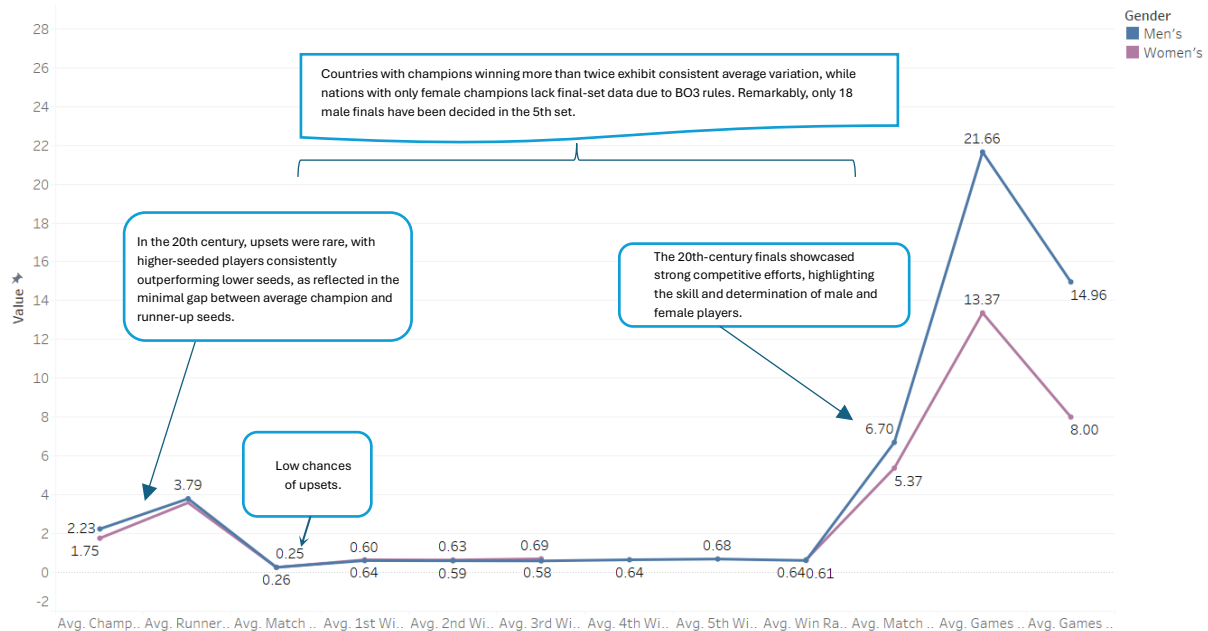*Figure 5: Performance by Country*

## Performance by Gender(20th Century)



*Figure 6: Performance Statistics by Gender for matches played in the 20th Century*
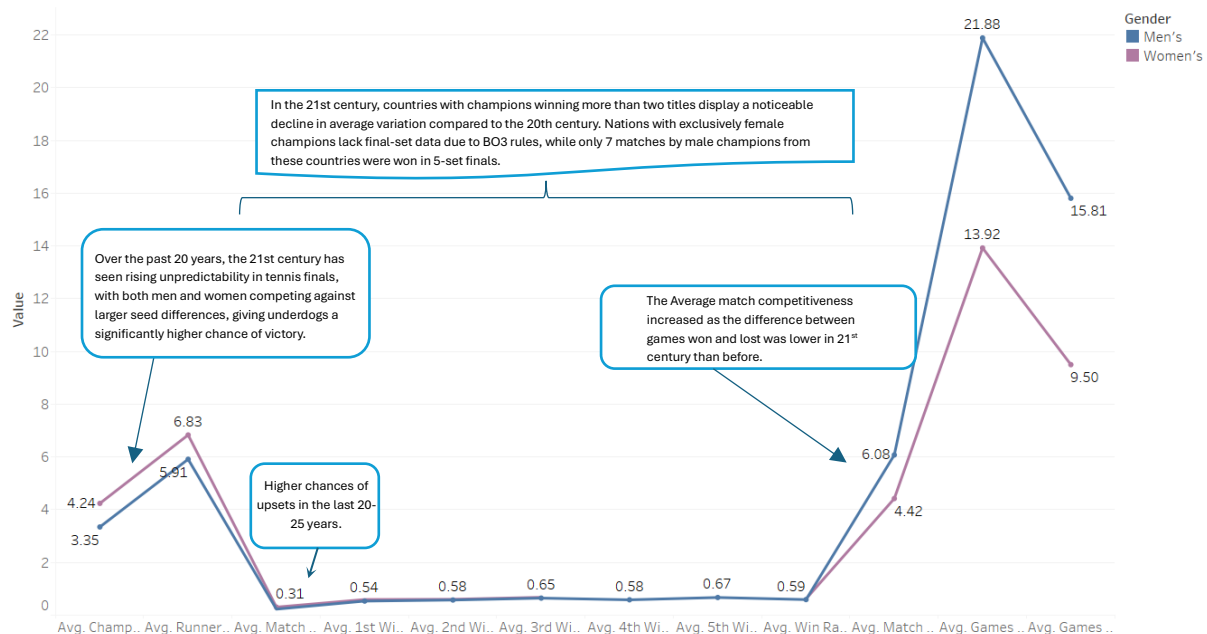
Performance by gender(21st Century)

In the 21st century, countries with champions winning more than two titles display a noticeable decline in average variation compared to the 20th century. Nations with exclusively female champions lack final-set data due to BO3 rules, while only 7 matches by male champions from these countries were won in 5-set finals.

Over the past 20 years, the 21st century has seen rising unpredictability in tennis finals, with both men and women competing against larger seed differences, giving underdogs a significantly higher chance of victory.

The Average match competitiveness increased as the difference between games won and lost was lower in 21st century than before.

Higher chances of upsets in the last 20-25 years.

*Figure 7: Performance by Gender for matches played in the 21st Century*

### 4.2.2 Chart Description

First all, measure name and measure values were dragged onto the Tableau canvas. Then all the measure were filtered and Average games won, average games lost, Average win Ratio, Average win ratio for each set, Average match competitiveness, Average Upset level difference were added to comparatively measure. The number of countries that took part in playing every set was counted and calculated and annotated to see if there were any outliers or patterns. The champion Nationality was used on the color shelf to differentiate each country's line. In the next graph, Champion country was replaced by Gender in the color mark and Year was added to the filter and two graphs were created one for the 20th Century and one that highlights the statistics of players in the 21st century. Any strange or unexpected variations in the graph have been labeled and annotated.

### 4.2.3 Insights and Trends

The parallel coordinates plot shows that top-seeded champions (low seed numbers) consistently achieve high stable wins rations and set win percentages while a few underdogs like Russia ( 3 wins, 2 Upsets) and Denmark stand out with greater variability and "upset" victories. Early round win rates remain uniformly strong (~ 60-70% and overall match competitiveness reflects a balance of dominant and tightly contested finals. Gender differences mirror format: men ( BO5) log more games won and lost exhibit higher late set win spikes  while women ( BO3) show slightly lower game volumes but similar upset

frequencies. From the 20<sup>th</sup> to the 21<sup>st</sup> century , upset counts have risen equally for both genders underling an era of deeper fields and heightened competitiveness across the draw.

## 4.3 Geographic Map Visualisation

### 4.3.1 Purpose

The geographic map visualization enables us to capture the spatial distribution of champions by linking their country data with geographic coordinates. This visualization is used to illustrate the shift from local to global dominance over time, compare the geographic distribution of champions and runners-up and highlight which countries have become tennis powerhouses and how this pattern has changed.

Upset Wins by Country (20th Century)



*Figure 8: Match Upsets By Country in the 20th Century*

Upset Wins by Country (21st Century)



Figure 9: Match Upsets by Country in the 21st Century

### 4.3.2 Chart Description

First champion country dimension and count of Australian Open was added to the detail mark of the map. Sum of match upset was added to the color mark. All metrics were aggregated to country level. For the next variation, we can add a Year Filter to compare upset variations in different time periods. These chart now tells us the number of upsets wins a country has enjoyed with respect to its total wins where their champion beat a higher seed player to win the Australian Open and how much it has between 20[th] and 21[st] century.

### 4.3.3 Insights and Trends

The geographic plots shows Australia with the highest number of Australian open wins (94) and match upsets (25). This reinforces the historical legacy of Australian tennis which was driven by local participation earlier on when international travel was limited. The data reflects that as the tournament opened up powerhouse nations like the United States (44 wins 14 upsets) and United Kingdom(28 wins , 3 upsets) emerged. This signals a shift from a regional to a global event where top tennis nations built strong competitive programs. Noticeable differences in the number of match upsets suggest variability in how often the expected outcomes were subverted. While Australia has higher upset count (25) the relative lower number for nations like Russia who won 3 Australian open tournaments out of which 2 were upsets where their player beat higher seeded players to win the title. This clearly shows excellent training programs set up by these countries to give new talent a

fighting chance on the international stage. Upon analyzing the upset percentage, countries like Denmark, Russia, Italy and Sweden emerge as top contenders. Historically, over 50% of the tournaments won by champions from these nations have been upset victories where lower seeded players triumph over higher seeded opponents in the finals. This pattern indicates that athletes from these countries not only have the capability to challenge expectations but also excel under pressure, often turning presumes underdogs into tournament winners.

In the early 20$^{th}$ century, due to travel restriction, most wins were secured by Australian champions which gave Australia a commanding win lead over the other countries. However, US and UK were close behind and have had an excellent run as well. Majority of the upsets in the early century came from these countries as well. However In the 21 century, the picture is very different. Australia has only one title win to date while many countries such as Russia and Japan have excellent 2 or more wins and upsets. However, US continues to maintain its dominance in Australian tennis with an exceptional margin clearly showcasing the talent and training that US tennis players have. Overall, no one country hold dominance in recet times and the unpridictaility of the game has increased with top seeds loosing finals to lower seed players through pure skill and endurance.

## 4.4 Scatter Plot Visualization

### 4.4.1 Purpose

Scatter charts are used to compare two numerical metrics simultaneously making them ideal for performance comparison. In this analysis, they serve to visualize the relationship between seed ranking and total games won, as well as identify any outliers that indicate unusual performance or trends and provide clear point-by-point comparison that can be further segmented by gender or nationality.
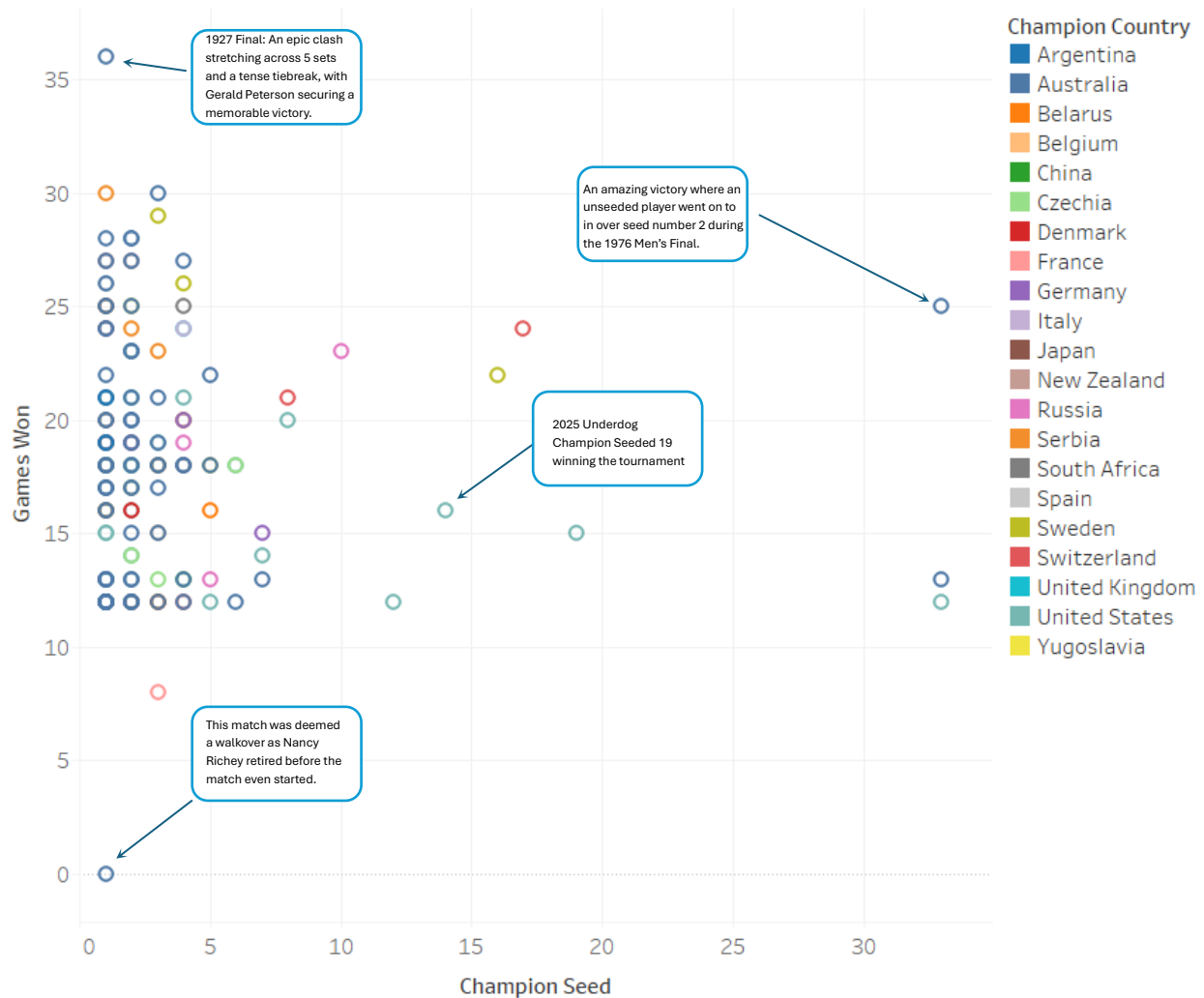
Figure 10: Scatter plot showcasing the relationship between champion seed and Games won.

## 4.4.2 Chart Preparation

This scatter plot was constructed in Tableau to explore the relationship between a champion's pre-tournament seed (X-axis) and their total games won in the final match(Y-axis) To build it champion seed was added to columns Games won was added to rows Champion country was added to color.

## 4.4.3 Insights and Trends

From the graph a few things are clear. Top seeds dominate the high games zone. The densest cluster appears at Champion Seed = 1 with games- won values mostly between 20-30 which confirms that #1 seeds not only win finals most often but also tend to rack up the highest game totals in those matches which if longer are indicative of men's finals as

these match follow BO5 rules. A handful of champions seeded greater than 10 each have won 12-24 games. An example is 2025 women's champion who sits at (19,15) who is an impressive underdog that manage to win the tournament despite her lower seed rankings. These outliers highlight exceptional performance where low seeded players overcame the odds.one point at (1,0) reflects a year where a walkover or retirement occurred in the final before any games were played. Seed 2-6 becomes a secondary cluster where game won values are usually between 22-28. This suggests that while #1 seeds lead in sheer volume, seeds 2-6 also deliver consistently strong performances in the finals. The single highest point at roughly (1,36) corresponds to a champion Gerald Patterson who played an exceptionally long muti set final with a tiebreaker reaching 18-16 but finally going on to clinch the title from Seed no. 3 John Hawkes in 1927.

## 5. Advantages and Disadvantages of Virtualization

Below are four advantages and disadvantages for each chart type.

### 5.1 Treemap

*Advantages*

1. Space Efficiency: Treemaps can display a large amount of hierarchical data ina compact area.
2. Immediate visual impact: They allow users to quickly identify dominant categories based on the size of the blocks.
3. Color Coding: They can be easily enhanced using colors to represent additional variables, such as a gender or win ratio.
4. Relative comparison: Users can immediately see the relative contribution of each category which is idea for understanding national performance in the dataset.

*Disadvantages:*

1. Complexity with many categories: if there are too many categories then the treemap can become cluttered and harder to interpret.
2. Limited Detail on Individual Data Points: While t is good for an overview, treemaps may not provide granular details about each record.
3. Difficulty in comparing Sizes precisely: Quantitative comparisons between blocks can be challenging without annotations.
4. Potential Misinterpretation: if the color scheme or sizing is not chosen carefully, to may mislead viewers regarding the important of some categories.

## 5.2 Parallel Coordinates

*Advantages:*

1. Multivariate data analysis: they allow simultaneous visual comparison of multiple performance metrics.
2. Pattern and correlation identification: Trends and relationships across several dimensions become apparent.
3. Flexibility: Users can filter and highlight individual lines(Players) to study variations in performance.
4. Effective Outlier Detection: Outliers can be quickly detected as they deviate starkly from the overall patterns.

*Disadvantages:*

1. Overplotting: With a large dataset, lines can overlap significantly making interpretation difficult.
2. Steep Learning curve: They may be less intuitive for audiences not familiar with high dimensional visualization.
3. Clustering Appearance: Without careful design and filtering , the plot can appear messy and obscure critical insights.
4. Difficulty in Quantitative Analysis: Extracting precise numerical values from the plot is challenging compared to other chart types.

## 5.3 Geographic Map

*Advantages:*

1. Spatial Understanding: Provides a clear visualization of geographic distributions and regional trends.
2. Contextual Relevance: Helps tie performance data to specific regions or countries, which is useful for global comparisons.
3. Interactivity: Geographic maps in tableau can be interactive, allowing users to drill down ito specific regions.
4. Instant recognition: The map format is familiar to most users and is immediately engaging.

*Disadvantages:*

1. Over-Simplification of Data: Geographical visualizations might hide the underlying detail if its too aggregated.
2. Reliance on Accurate Geocoding: Inaccurate or incomplete location data can lead to mis representations.
3. Limited Quantitative Detail: While Effective for spatial Relationships, it is less useful for detailed numerical Analysis.

4. Scalability Issues: with too many datapoints or overlapping regions, the map can become cluttered and less informative.

## 5.4 Scatter Chart

*Advantages:*

1. Clear Correlations: They effectively display the relationship between two continuous variables.
2. Identification of outliers: Outliers stand out clearly from clusters of data points.
3. Simplicity: they are easy to interpret and understand for both experts and non-experts audiences.
4. Flexibility in Segmentation: Data points can be color coded, sized or shaped to represent additional dimensions( e.g. Gender, nationality)

*Disadvantages:*

1. Limited to two variables: Primarily shows relationships between just two variables at a time which may require multiple charts for comprehensive view.
2. Potential overplotting: with large datasets, overlapping points can obscure the true relationships.
3. Scaling issues: The Relative scale between variables may distort perception if not properly normalized.
4. Requires Supplemental Information: Often needs additional annotations or trend lines to fully communicate insights.

## 6. Conclusion

This report set out to illustrate more thana century of Australian Open finals through four carefully chosen visualization techniques which are tree maps , parallel coordinates, geographic maps and scatter charts, each implemented in tableau with complimentary excel support. By methodically preparing the raw data, parsing textual scores , calculating new metrics and aggregation by nationality and gender a Robust foundation was bult for analysis that speaks both to tennis aficionados and data enthusiasts alike.

Key Takeaways:

1. A shifting Global Landscape: Early Australian dominance gave way to a truly international contest. Tree maps vividly showed the transition from a blue dominated Australia on the early 20th century to a kaleidoscope of champions from the USA , Serbia, Spain and beyond. Gender inclusion grew from a female field introduced in 19222 to today's fully global women's draw with powerhouses like the USA and Belarus emerging over recent decades.

2. Performance Profiles and Outliers: Parallel Coordinates revealed that regardless of era or gender, champions tend to maintain high win rations through every set. Yet underdog stories stand out such as Russia and Denmark each boast multiple "upset" titles despite lower seeds, underscoring the unpredictability that makes tennis so compelling. Champion seed consistently correlated with match dominance as the densest clusters of games won are at seed 1-6 while rare low seed triumphs ( e.g. 2025's 19 seeded women's championship) become fan favorite narratives of grit and sheer determination.

3. Geography of success and Surprise: Mapping champions and their match upsets exposed not only where greatness was born but where it most often disrupted expectations. Australia and the USA lead in volume, yet Sweden and the UK achieved titles with zero upsets such "quiet champions" who met but never defied expectations while smaller nations like Russia and Denmark punched above their weight in upset percentages.

4. Match Dynamics and endurance: The scatter plot of seed versus games won crystallized how physical format shapes competition: men's best of five finals yield higher game counts than women's best of three creating a strategic and endurance differential that parallel broader gender based narratives in tennis.

5. Visualisation Best practices: Annotations anchored surprising outliers such as walkovers, record long finals and Cinderella runs, transforming raw points into vivid stories. Interactivity via filters by era, gender, seed range invited users to explore what-if scenarios and dive deeper into specific matchups. Consistent color coding and clear axis labels ensuring accessibility allowing both data novices and seasoned analysts to glean insights at a glance.

In compiling this report, not just a century of champions was charted but rigorous data preparation, thoughtful transformation and well-chosen visualization helped enrich our understanding of Tennis' greatest triumphs and most thrilling upsets. As tennis continue to evolve with new champions, new nations and new narrative, this analytical framework will remain a powerful tool for fans and analyst to celebrate, question and predict the ever-unfolding drama of the Australin Open.

## 7. Bibliography

- *University of Technology Sydney - Sign In*. (2025). Uts.edu.au. https://canvas.uts.edu.au/courses/34202/assignments/208749

- Wikipedia Contributors. (2019, November 27). *Australian Open*. Wikipedia; Wikimedia Foundation. https://en.wikipedia.org/wiki/Australian_Open