

## Data Science Exam

Precious Nhamo

### Abstract

This document comprises the responses to **Questions 1 to 5** of the 2025 Data Science examination, along with separate analyses (e.g., a PowerPoint presentation) completed in accordance with the exam instructions and organised within the designated folder.

Data Science Exam

Question1 Baby Names

We report how we determined our sample size, all data exclusions (if any), all manipulations, and all measures in the study.

Question2 Music Taste

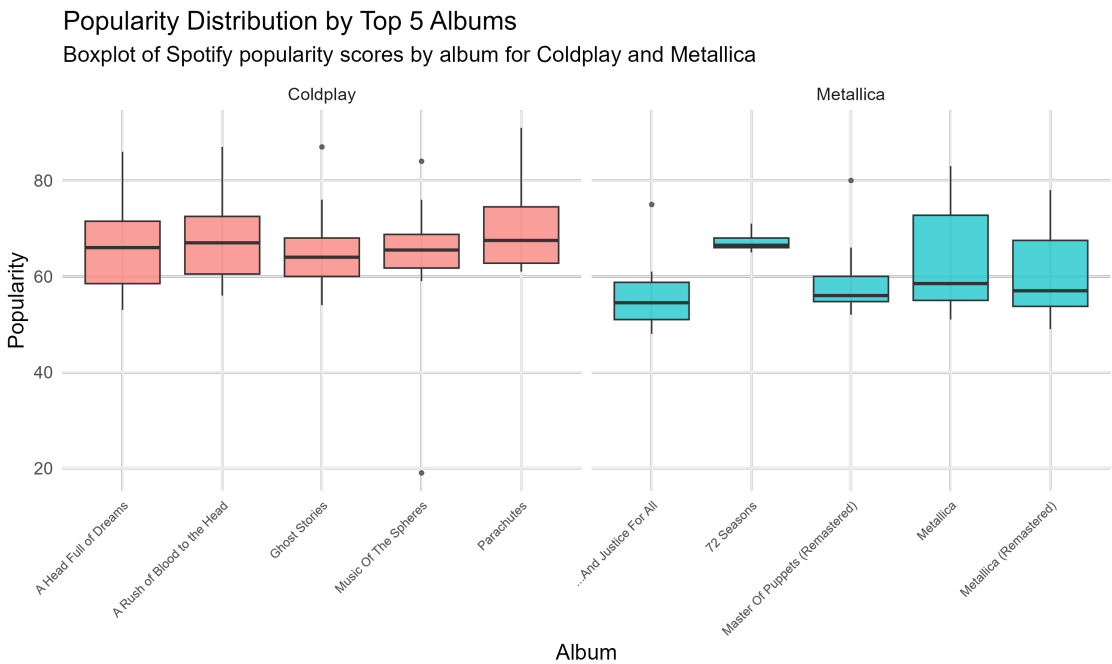


Figure 1

Longevity and Musical Progression of Coldplay and Metallica

The data reveals distinct trajectories for Coldplay and Metallica in terms of popularity, musical evolution, and industry adaptation. Coldplay demonstrated early dominance, charting five songs in their first decade compared to Metallica’s one (Figure 2). Their popularity scores on Spotify also show broader appeal, with a higher median and narrower interquartile range than Metallica’s (Figure 1).

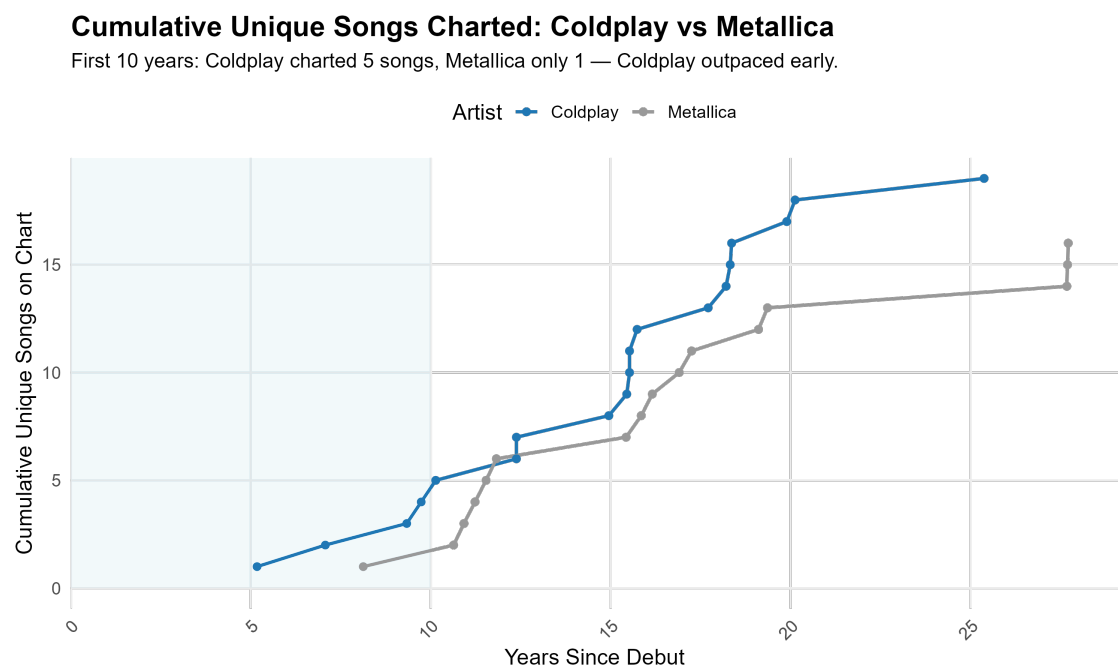


Figure 2

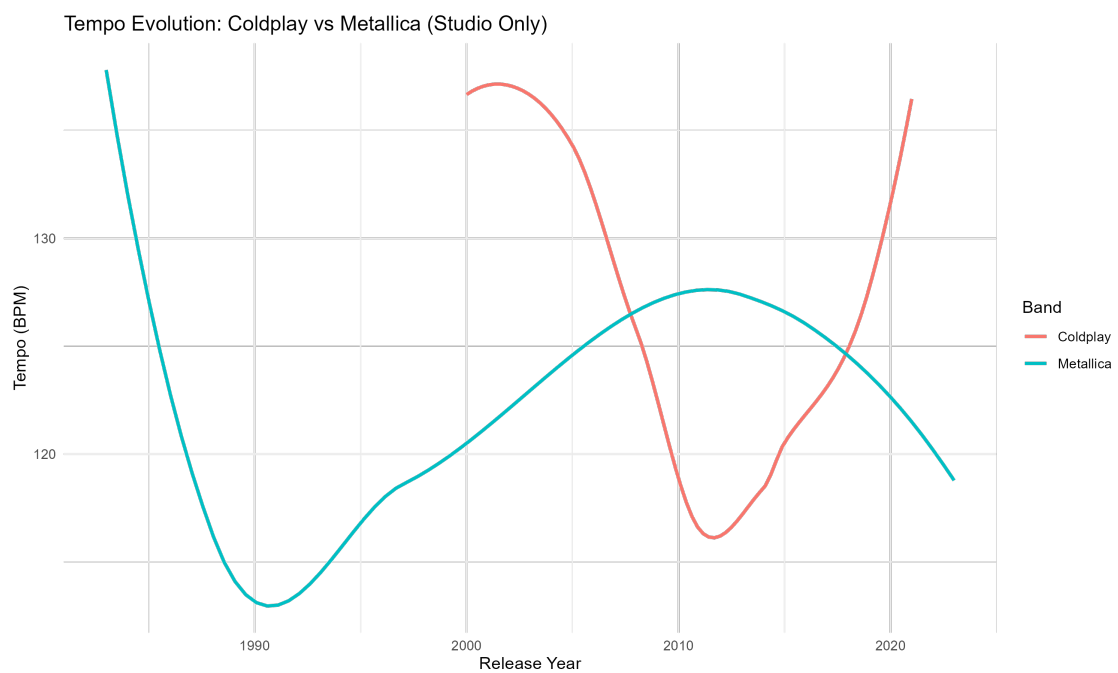


Figure 3

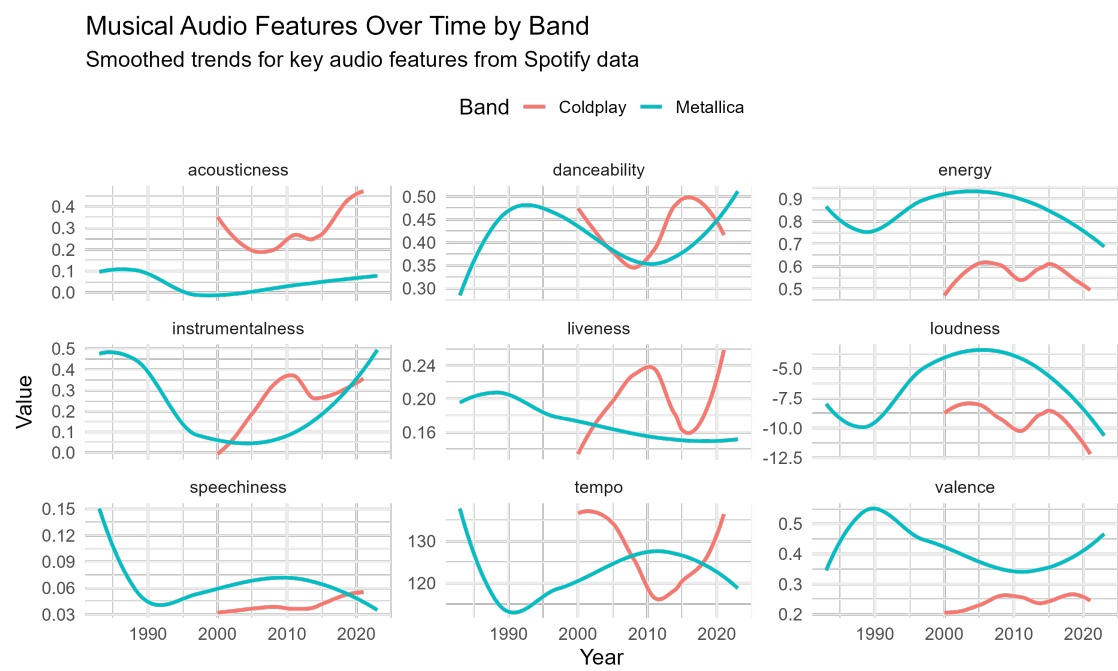


Figure 4

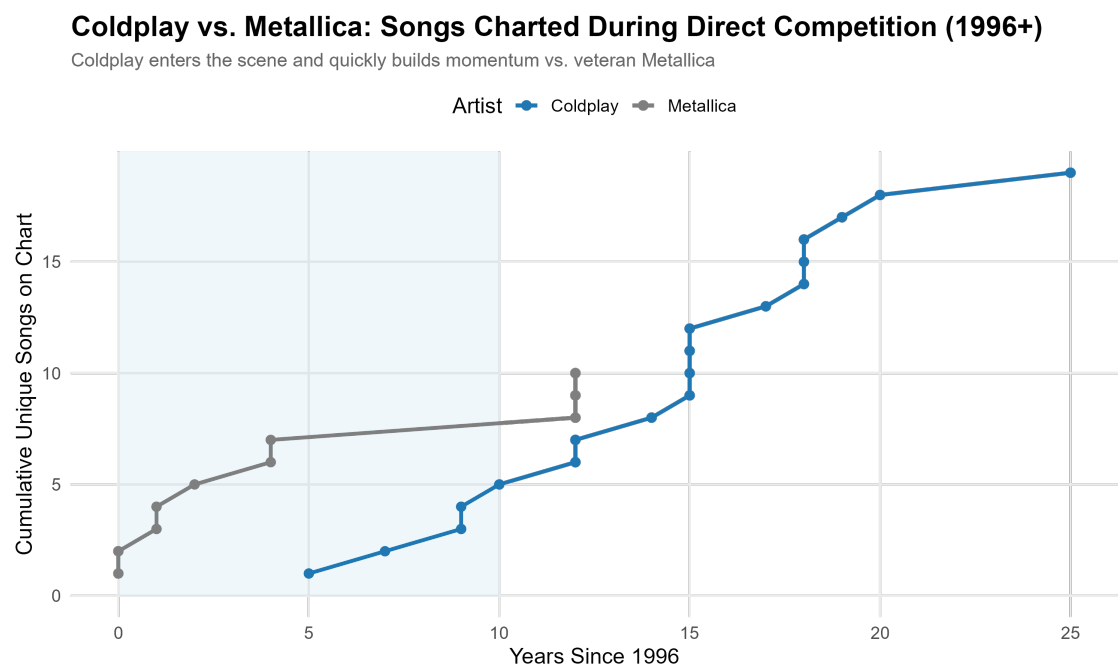


Figure 5

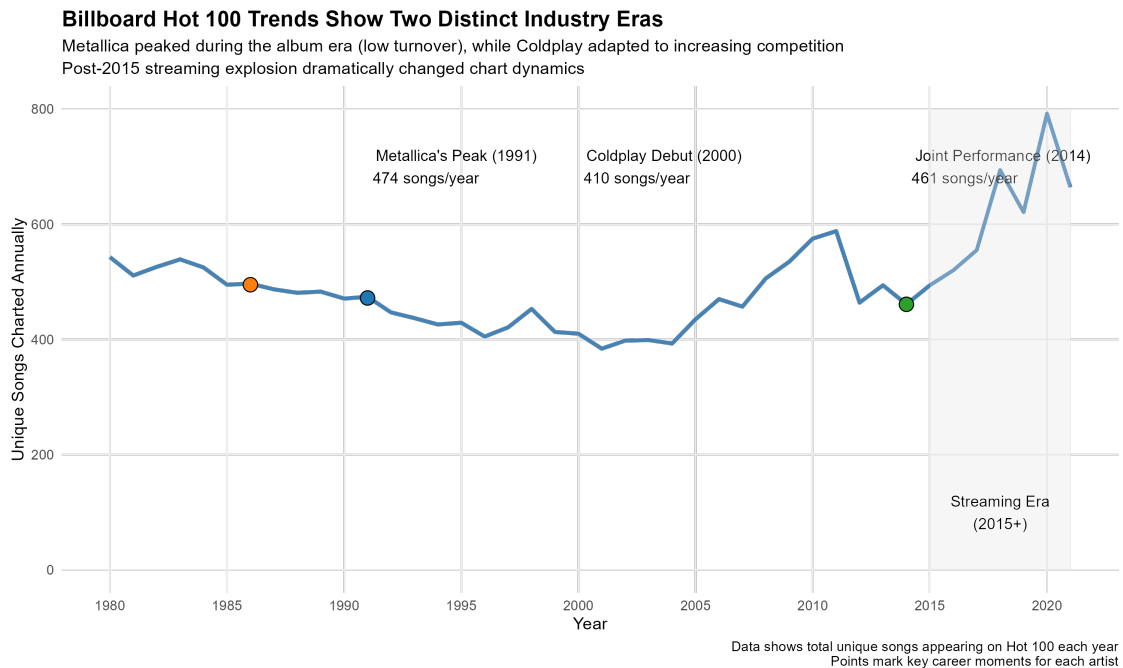


Figure 6

Musically, Coldplay's tempo has remained stable (Figure 3), while their audio features, such as danceability and valence, trended positively over time (Figure 4). Metallica, conversely, maintained higher instrumentalness and energy, reflecting their heavier style.

Billboard data highlights their adaptation to industry shifts: Metallica peaked during the album era (1991), while Coldplay thrived post-2000, leveraging streaming's rise (Figure 6). During direct competition (1996+), Coldplay's momentum outpaced Metallica's (Figure 5).

Coldplay's consistent, accessible sound contrasts with Metallica's enduring heavy metal identity. Both bands exemplify longevity but reflect divergent strategies in navigating musical trends

Question 3 : Netflix Content Strategy Analysis

In light of Netflix’s recent subscriber attrition and share price volatility, a strategic review was conducted to inform potential market entry for a new streaming venture. This review draws on IMDb ratings and global production data to assess what drives success in streaming content.

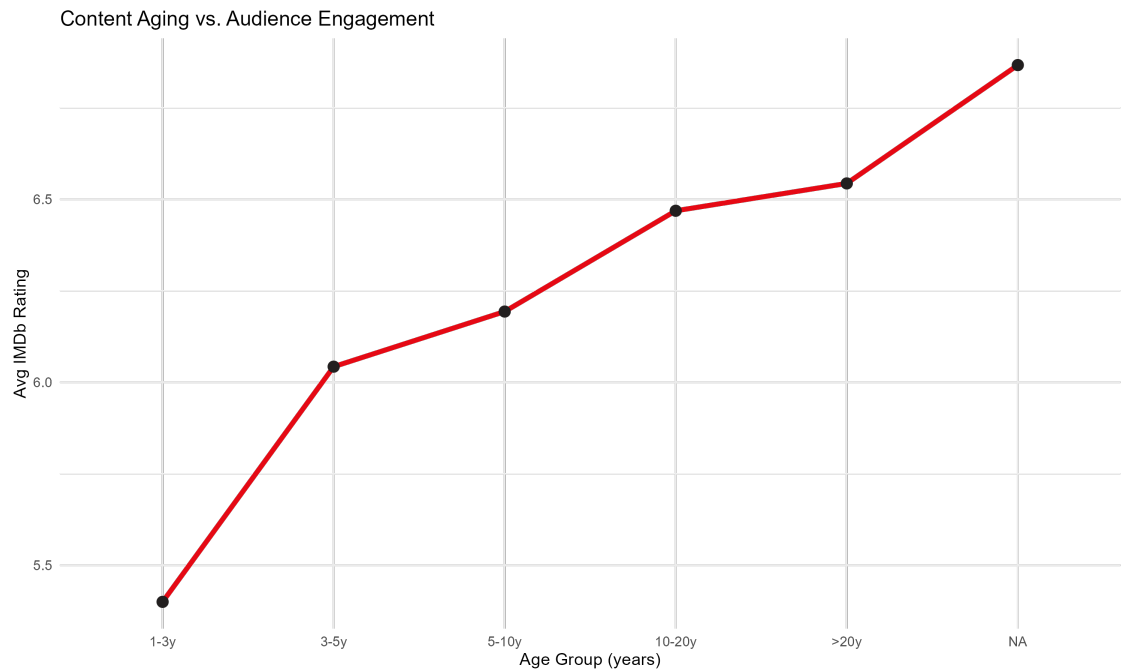


Figure 7

Key Findings

**1. Quality vs. Volume Trade-Off.** Japanese titles, while constituting only 1.6% of Netflix’s catalogue, achieve superior average IMDb ratings (mean = 6.7, SD =  $\pm 0.3$ ). By contrast, the United States and India collectively contribute over 50% of Netflix’s content but yield lower average ratings (mean = 6.1–6.3) (Figures 8, 15) **Implication:** High-volume strategies may dilute perceived quality.

**2. Genre Performance Differentials.** Content genres vary markedly in audience reception. Documentaries (mean = 7.6) and dramas (mean = 7.0) lead on quality metrics,

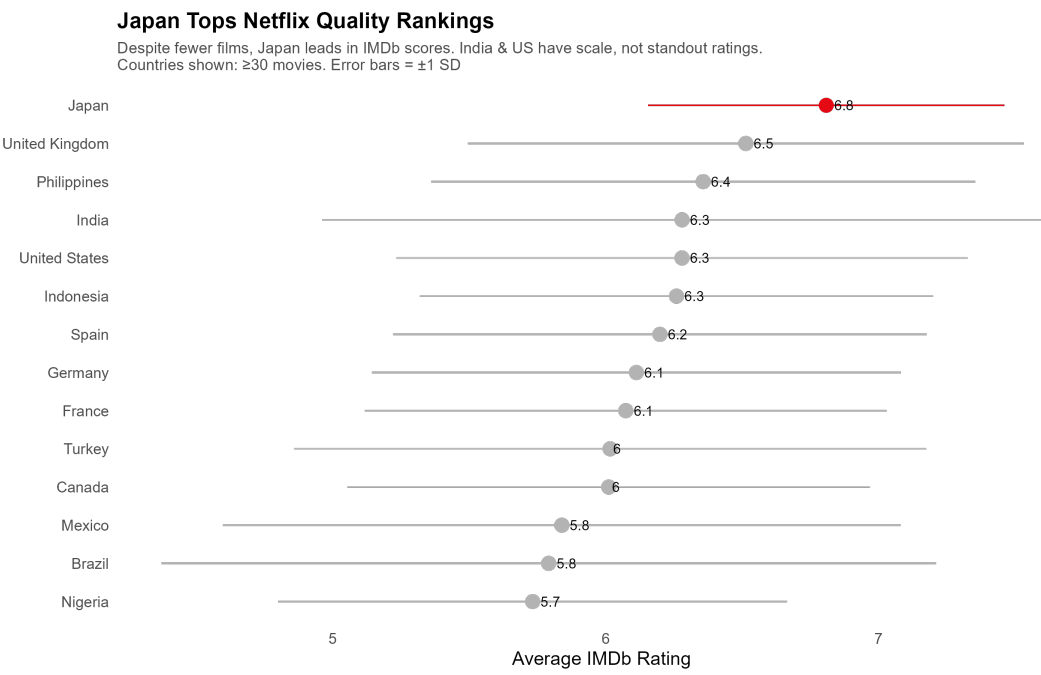


Figure 8

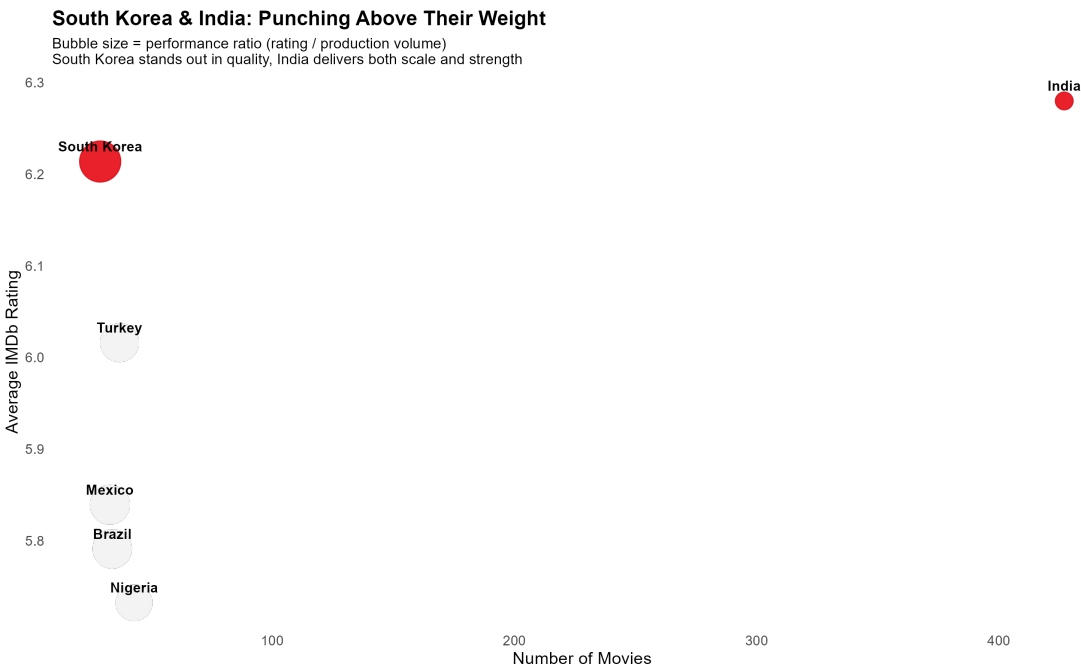


Figure 9



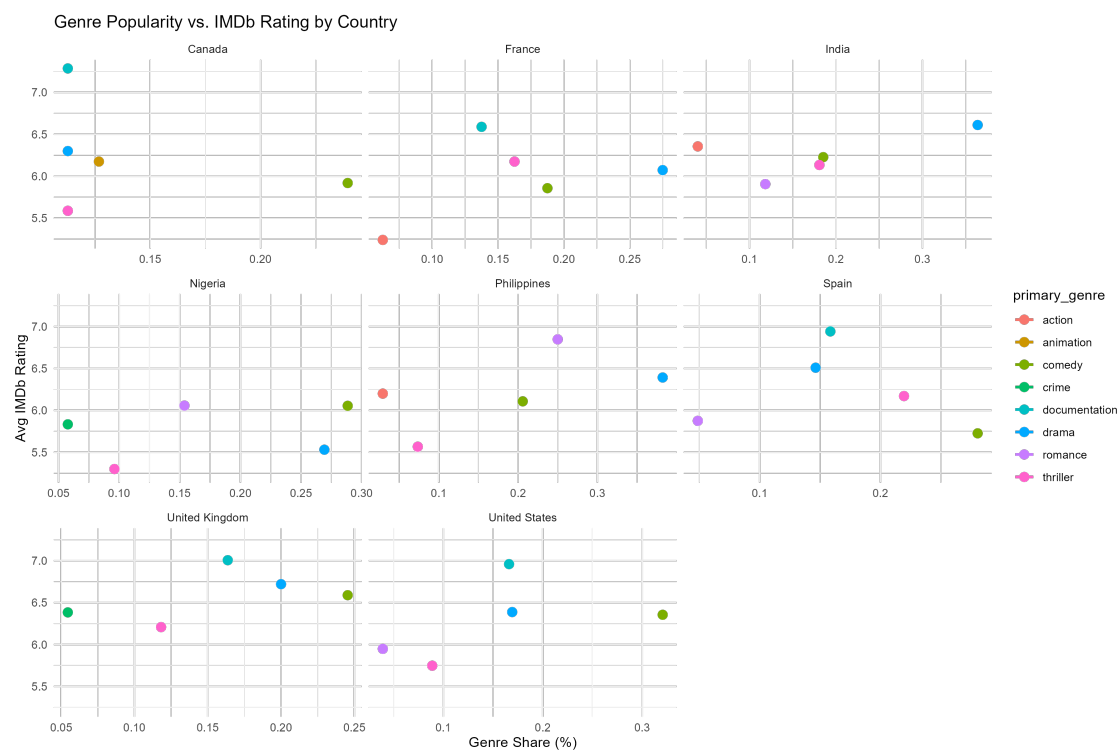


Figure 10

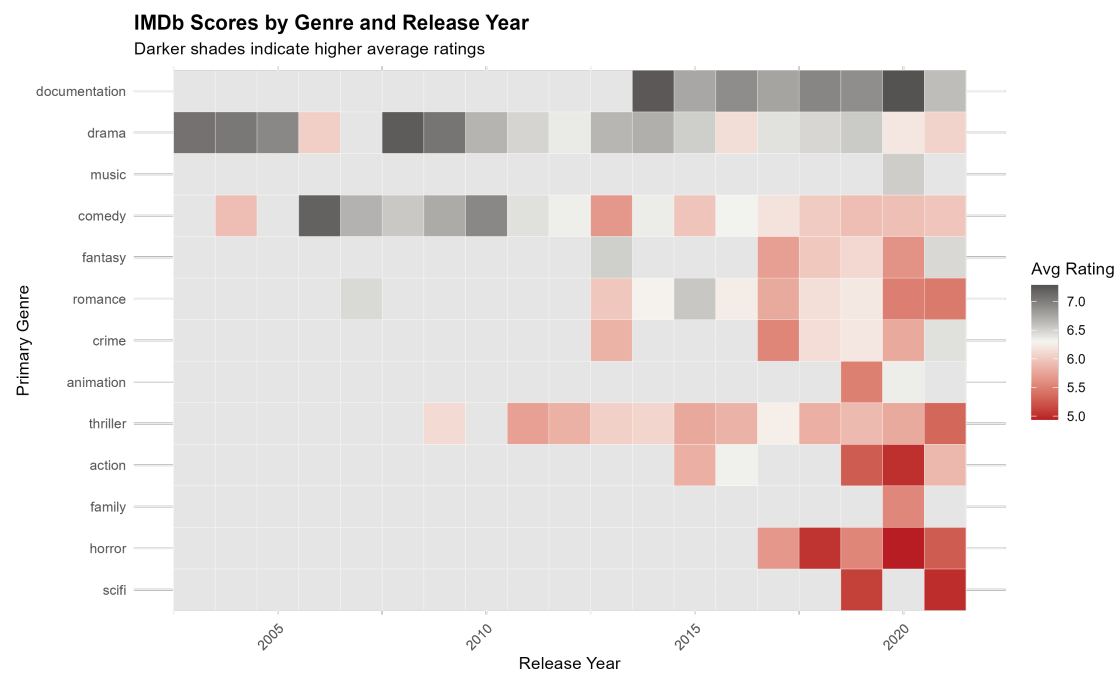


Figure 11

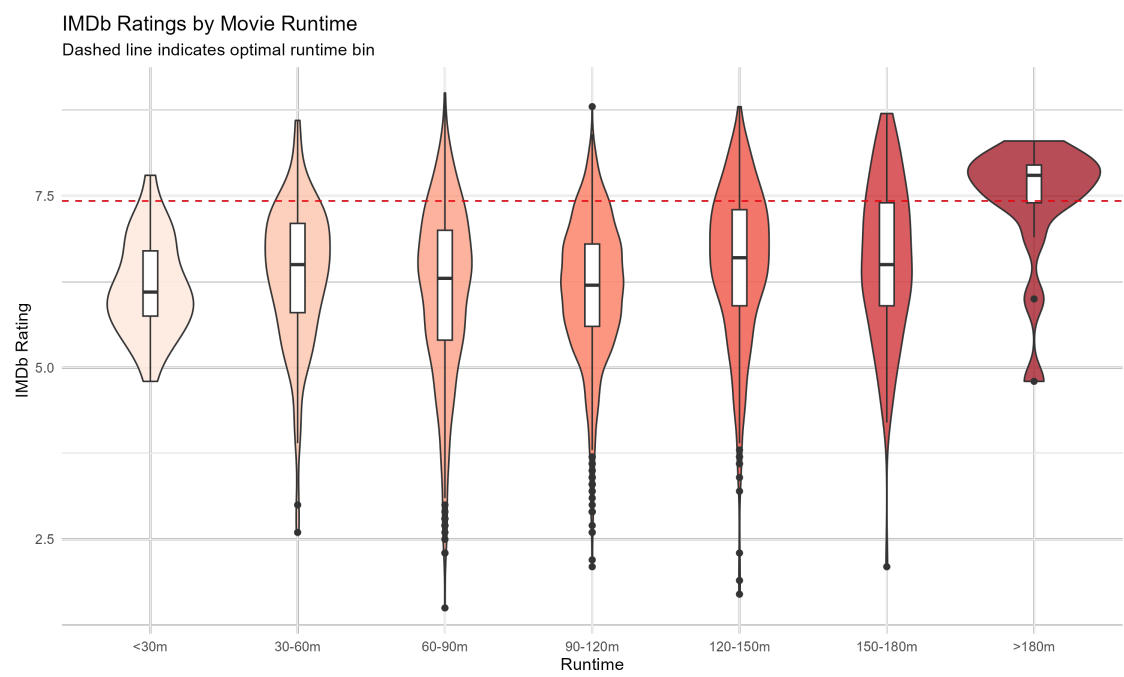


Figure 12

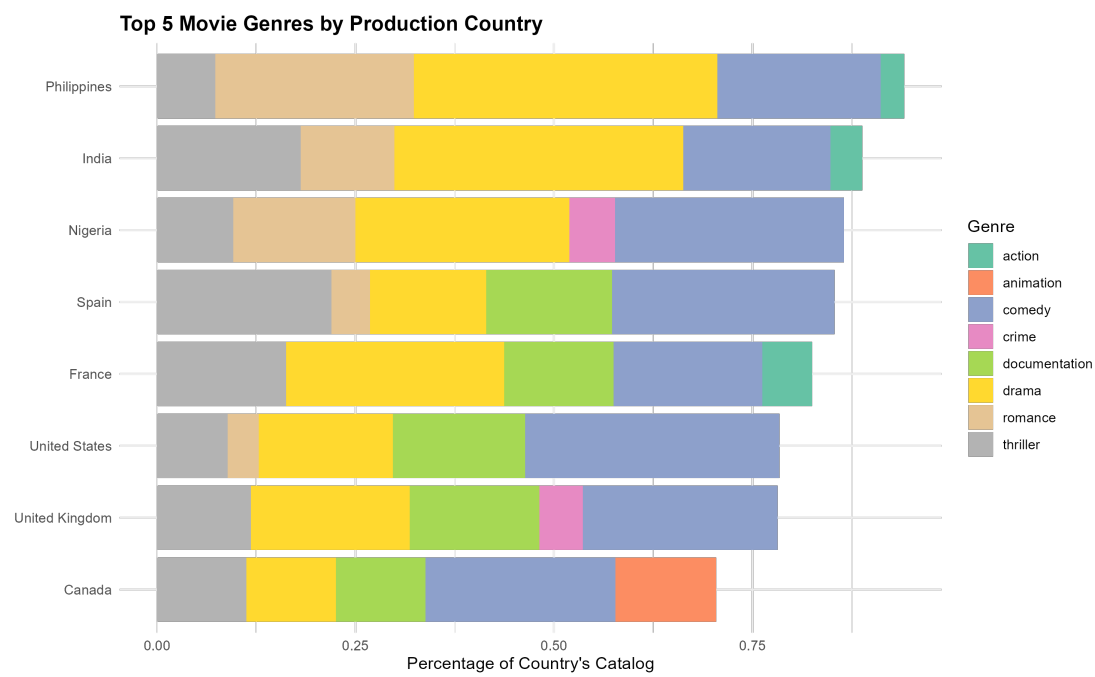


Figure 13

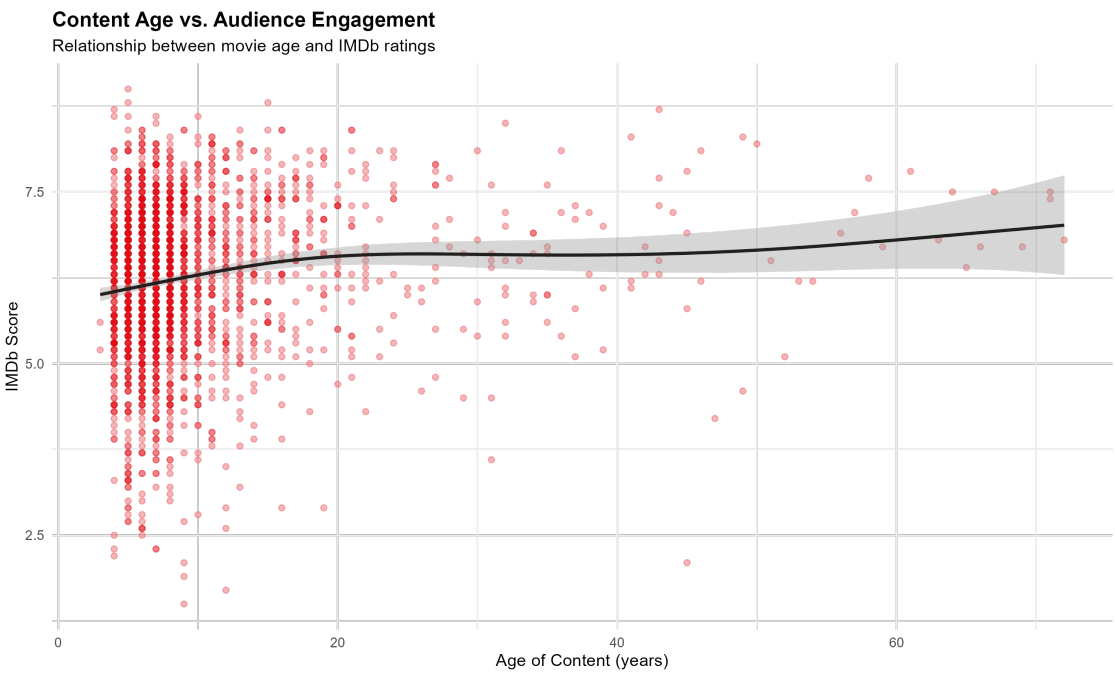


Figure 14

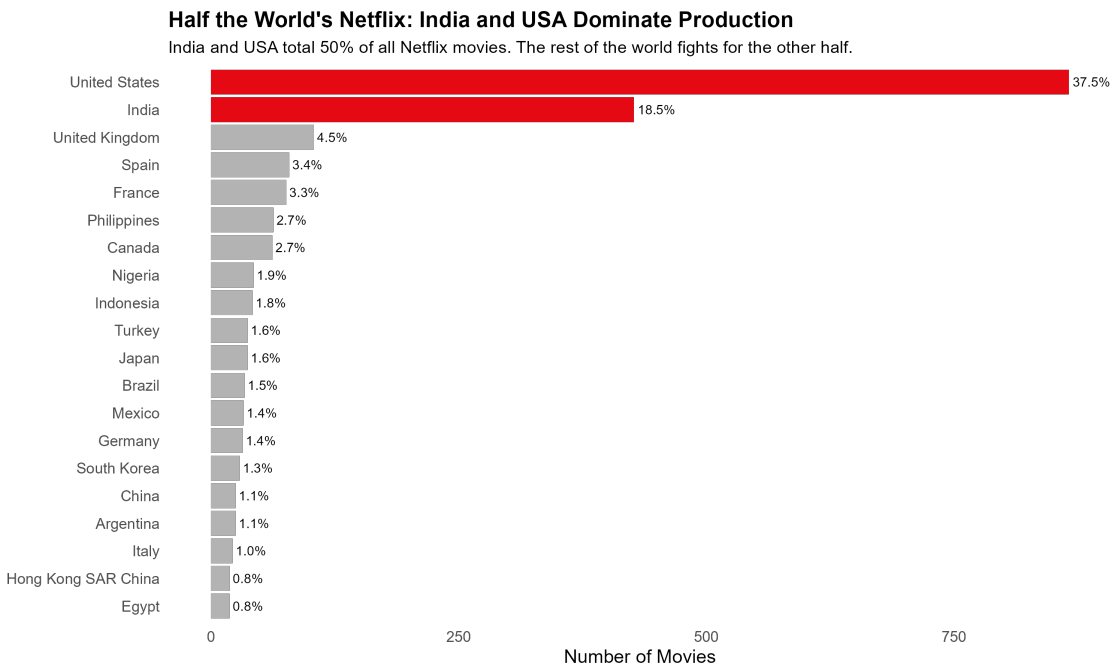


Figure 15

while action and comedy lag behind (mean = 6.2 and 6.5 respectively) (Figure 11). Crime and thriller categories also demonstrate strong viewer engagement.

**Implication:** Narrative depth and authenticity drive audience satisfaction.

**3. Optimal Runtime Windows.** Ratings peak for films with runtimes between 90–120 minutes (mean = 7.5). Shorter or longer formats tend to correlate with lower ratings (Figure 12). **Implication:** Viewer engagement is maximised within a standardised runtime threshold.

**4. Age and Longevity of Content.** Older films (>20 years) sustain higher average ratings (mean = 8.5)(Figure 7) compared to recent productions (1–3 years; mean = 6.0)(Figure 14). **Implication:** Legacy content benefits from curation and nostalgia; newer content faces saturation and discovery challenges.

## Strategic Recommendations for New Market Entrant

To position competitively against incumbent platforms:

**1. Targeted Content Acquisition.** Prioritise **high-quality, underrepresented niches**—such as Japanese documentaries and international dramas—to differentiate on quality and build a prestige brand image.

**2. Genre Investment Strategy.** Allocate production and licensing budgets towards **high-performing genres** (e.g., crime, thriller, drama). Avoid overserved categories like generic comedy unless uniquely positioned.

**3. Runtime Standardisation.** Anchor commissioned or licensed content within the 90–120 minute range to enhance user engagement and completion rates.

## Conclusion

A differentiated, quality-first strategy that leverages genre insights, runtime discipline, and curated international content provides a clear pathway for new entrants to

gain market share and investor confidence—without replicating the scale-driven pitfalls observed in Netflix’s recent trajectory.

#Question4 Billionaires

## Data analysis

We used R (Version 4.4.3; R Core Team, 2025) and the R-packages *dplyr* (Version 1.1.4; Wickham, François, Henry, Müller, & Vaughan, 2023), *fastDummies* (Version 1.7.5; Kaplan, 2025), *forcats* (Version 1.0.0; Wickham, 2023a), *ggplot2* (Version 3.5.2; Wickham, 2016), *ggrepel* (Version 0.9.6; Slowikowski, 2024), *lubridate* (Version 1.9.4; Grolemund & Wickham, 2011), *papaja* (Version 0.1.3; Aust & Barth, 2024), *patchwork* (Version 1.3.0; Pedersen, 2024), *purrr* (Version 1.0.4; Wickham & Henry, 2025), *readr* (Version 2.1.5; Wickham, Hester, & Bryan, 2024), *stringr* (Version 1.5.1; Wickham, 2023b), *tibble* (Version 3.2.1; Müller & Wickham, 2023), *tidyr* (Version 1.3.1; Wickham, Vaughan, & Girlich, 2024), *tidyverse* (Version 2.0.0; Wickham et al., 2019), and *tinylabels* (Version 0.2.5; Barth, 2025) for all our analyses. ##Results

##Discussion

## Question 5 Health

### On Air Script

For the 9-to-5 gamer: 1 hour of extra sleep beats 1 hour on the treadmill. Data shows stress-free living is your real cheat code.

We’ve been told for years to ‘exercise more’—but new data reveals the real game-changers: sleep and stress management. Here’s the proof: Adults with poor sleep and high stress gained 7–10 pounds more than those with good sleep and low stress—based on regression analysis ( $p < 0.01$ ). The worst hit? Early-career professionals (ages 30–39), who

saw up to 10 pounds more weight gain when stressed—even if they exercised. Why this matters: Gym memberships won’t fix this. The solution? Protect your sleep: Aim for 7–9 hours, and keep a consistent schedule—it’s your metabolic lifeline. Tame stress: Short walks, 5-minute breathing exercises, or setting work boundaries can slash health risks. Bottom line? Small changes to sleep and stress beat marathon workouts. Start tonight—your body will thank you.

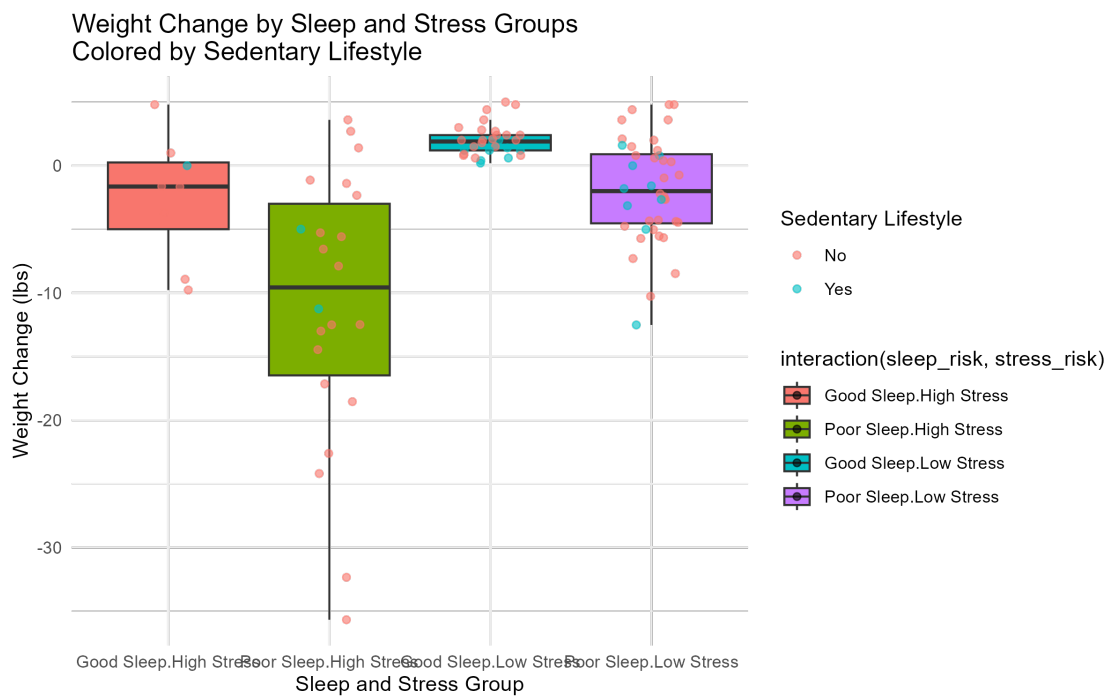


Figure 16

Table 1

(#tab:sleep-risk-pivot )Weight Change Summary

Gender	age_group	mean_weight_change	median_weight_change	n
F	Adults (20-29)	-1.0181151	-0.3680584	10
F	Early Career (30-39)	-2.4545587	-1.3943583	11
F	Middle-aged (40-49)	-3.4805916	0.7000000	12
F	Preretirement (50-64)	-1.5639095	0.0000000	9
F	Teenagers (<20)	-1.7940640	-1.7940640	1
M	Adults (20-29)	-2.3794779	-0.1679132	14
M	Early Career (30-39)	-5.5847844	0.6000000	12
M	Middle-aged (40-49)	-3.7051828	0.4500000	12
M	Preretirement (50-64)	-0.8648233	-0.3807078	14
M	Teenagers (<20)	-5.2528414	-5.7184960	5

	Sleep-Stress Interaction	Stress-Age Interaction
(Intercept)	1.758 (6.660)	-0.315 (1.510)
sleep_riskPoor Sleep	-7.168*** (2.704)	
stress_riskLow Stress	5.800** (2.717)	
physical_activityModerately Active	-1.142 (1.893)	
physical_activitySedentary	-0.539 (1.982)	
physical_activityVery Active	-1.076	

	Sleep-Stress	Stress-Age
	Interaction	Interaction
	(2.311)	
age_groupEarly Career (30-39)	-2.976	0.653
	(1.875)	(2.135)
age_groupMiddle-aged (40-49)	-2.226	0.107
	(1.766)	(2.103)
age_groupPreretirement (50-64)	-0.777	0.164
	(1.891)	(2.135)
age_groupTeenagers (<20)	-3.178	-1.867
	(2.861)	(3.094)
GenderM	-1.650	
	(1.597)	
Daily Caloric Surplus/Deficit	0.002	
	(0.003)	
BMR (Calories)	-0.001	
	(0.002)	
Duration (weeks)	-0.158	
	(0.175)	
sleep_riskPoor Sleep $\times$ stress_riskLow	2.980	
Stress		
	(3.290)	
High Stress RiskYes		-4.493*
		(2.615)
High Stress RiskYes $\times$ age_groupEarly		-10.049***
Career (30-39)		(3.785)



	Sleep-Stress Interaction	Stress-Age Interaction
High Stress RiskYes × age_groupMiddle-aged (40-49)		-7.112*  (3.767)
High Stress RiskYes × age_groupPreretirement (50-64)     1.248		(3.785)
High Stress RiskYes × age_groupTeenagers (<20)		-10.476  (7.112)
Num.Obs.	100	100
R2	0.455	0.402
R2 Adj.	0.366	0.342
AIC	655.5	654.9
BIC	697.2	683.5
Log.Lik.	-311.752	-316.433
RMSE	5.47	5.73
• p < 0.1, ** p < 0.05, *** p < 0.01		

Key Findings on Health Determinants.

1. Sleep Quality is Critical:

- Poor sleep correlates with **-7.2 lbs** weight gain ( $p<0.01$ ), surpassing the impact of physical activity (non-significant coefficients).
- Teens and early-career adults with poor sleep show the worst outcomes (**-7.7 to -8.0 lbs** mean weight change).Teens are abstracted from main analysis as their n is very small , n =5 .

Table 2

(#tab:sleep-risk-pivot )*Sleep Risk Pivot Summary*

Poor Sleep Risk	age_group	Mean_Weight_Change	Count
No	Adults (20-29)	0.2142665	9
No	Early Career (30-39)	1.5777778	9
No	Middle-aged (40-49)	2.9571429	7
No	Preretirement (50-64)	-0.0982816	11
No	Teenagers (<20)	1.9500000	2
Yes	Adults (20-29)	-3.0281494	15
Yes	Early Career (30-39)	-7.7298256	14
Yes	Middle-aged (40-49)	-6.2899584	17
Yes	Preretirement (50-64)	-2.0918011	12
Yes	Teenagers (<20)	-7.9895678	4

## 2. Stress Drives Negative Outcomes:

- High stress alone leads to **-4.5 lbs** weight gain ( $p < 0.1$ ).
- Stress combined with poor sleep exacerbates effects, especially in younger age groups (interaction terms up to **-10.0 lbs**,  $p < 0.01$ ).

## 3. Sedentary Lifestyle Modifies Risk:

- Sedentary individuals (red dots, Figure 16) show higher weight variability, but sleep/stress dominate overall trends.

## Practical Recommendations.

- **Prioritise Sleep Hygiene:** Consistent sleep schedules improve metabolic health more than moderate exercise.
- **Stress Management:** Mindfulness or flexible work policies could mitigate high-stress impacts, particularly for younger adults.

Table 3

(#tab:sleep-risk-pivot )Stress Risk Pivot Summary

Poor Sleep Risk	age_group	Mean_Weight_Change	Count
No	Adults (20-29)	0.2142665	9
No	Early Career (30-39)	1.5777778	9
No	Middle-aged (40-49)	2.9571429	7
No	Preretirement (50-64)	-0.0982816	11
No	Teenagers (<20)	1.9500000	2
Yes	Adults (20-29)	-3.0281494	15
Yes	Early Career (30-39)	-7.7298256	14
Yes	Middle-aged (40-49)	-6.2899584	17
Yes	Preretirement (50-64)	-2.0918011	12
Yes	Teenagers (<20)	-7.9895678	4

- **Targeted Interventions:** Early-career professionals need tailored programs addressing sleep and stress.

#### Slide Deck Outline.

1. **Title Slide:** Sleep Over Squats
2. **Key Chart:** Figure 16 (weight change by sleep/stress groups, colored by sedentary status).
3. **Regression Snapshot:** Highlight significant coefficients (sleep, stress, age interactions).
4. **Call to Action:** “Policy Focus: Sleep & Stress > Gym Memberships.

## References

- Aust, F., & Barth, M. (2024). *papaja: Prepare reproducible APA journal articles with R Markdown*. <https://doi.org/10.32614/CRAN.package.papaja>
- Barth, M. (2025). *tinylabels: Lightweight variable labels*.  
<https://doi.org/10.32614/CRAN.package.tinylabels>
- Grolemund, G., & Wickham, H. (2011). Dates and times made easy with lubridate. *Journal of Statistical Software*, 40(3), 1–25. Retrieved from <https://www.jstatsoft.org/v40/i03/>
- Kaplan, J. (2025). *fastDummies: Fast creation of dummy (binary) columns and rows from categorical variables*. Retrieved from <https://CRAN.R-project.org/package=fastDummies>
- Müller, K., & Wickham, H. (2023). *Tibble: Simple data frames*. Retrieved from <https://CRAN.R-project.org/package=tibble>
- Pedersen, T. L. (2024). *Patchwork: The composer of plots*. Retrieved from <https://CRAN.R-project.org/package=patchwork>
- R Core Team. (2025). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from <https://www.R-project.org/>
- Slowikowski, K. (2024). *Ggrepel: Automatically position non-overlapping text labels with 'ggplot2'*. Retrieved from <https://CRAN.R-project.org/package=ggrepel>
- Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York. Retrieved from <https://ggplot2.tidyverse.org>
- Wickham, H. (2023a). *Forcats: Tools for working with categorical variables (factors)*. Retrieved from <https://CRAN.R-project.org/package=forcats>
- Wickham, H. (2023b). *Stringr: Simple, consistent wrappers for common string operations*. Retrieved from <https://CRAN.R-project.org/package=stringr>
- Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., . . .

- Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. <https://doi.org/10.21105/joss.01686>
- Wickham, H., François, R., Henry, L., Müller, K., & Vaughan, D. (2023). *Dplyr: A grammar of data manipulation*. Retrieved from <https://CRAN.R-project.org/package=dplyr>
- Wickham, H., & Henry, L. (2025). *Purrr: Functional programming tools*. Retrieved from <https://CRAN.R-project.org/package=purrr>
- Wickham, H., Hester, J., & Bryan, J. (2024). *Readr: Read rectangular text data*. Retrieved from <https://CRAN.R-project.org/package=readr>
- Wickham, H., Vaughan, D., & Girlich, M. (2024). *Tidyr: Tidy messy data*. Retrieved from <https://CRAN.R-project.org/package=tidyr>