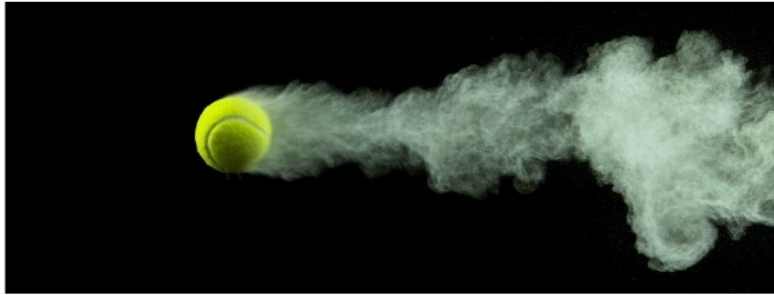


2024 MCM问题C:网球的动量



在2023年温布尔登男单决赛中，20岁的西班牙新星卡洛斯·阿尔卡拉兹击败了36岁的诺瓦克·德约科维奇。这是德约科维奇自2013年以来在温布尔登的首次失利，也终结了这位**大满贯**史上最伟大球员之一的辉煌战绩。

这场比赛本身就是一场非凡的战斗。[1]德约科维奇似乎注定要轻松获胜，因为他以6比1的比分控制了第一局(7局中赢了6局)。然而第二局气氛紧张，最终阿尔卡拉兹在抢七局中以7-6获胜。第三局与第一局相反，阿尔卡拉兹以6-1轻松获胜。这位年轻的西班牙人在第四盘开始时似乎完全控制了比赛，但不知何故，比赛再次改变了方向，德约科维奇完全控制了比赛，以6比3赢得了比赛。第五盘也是最后一盘，德约科维奇从第四盘开始保持优势，但再次改变方向，阿尔卡拉兹控制局面，以6-4获胜。本场比赛数据在提供的数据集中，“2023-温布尔登-1701”的“match_id”。你可以看到德约科维奇在第一盘领先时的所有得分，使用“set_no”列等于1。出现在看似有优势的球员身上的不可思议的挥拍，有时是多分甚至多局，往往被归结为“冲劲”。

字典上对动量的一个定义是“通过运动或一系列事件获得的力量或力量”。[2]在体育运动中，一支球队或一名球员可能会在比赛/比赛中感觉到他们有动量，或“力量/力量”，但这种现象很难测量。此外，如果动量存在的话，比赛中的各种事件是如何产生或改变动量的，这一点也不容易弄清楚。

2023年温布尔登男单前两轮之后的每一分数据。您可以自行选择加入额外的球员信息或其他数据，但您必须完整地记录来源。使用这些数据：

- 开发一个模型，捕捉得分发生时的比赛流程，并将其应用于一场或多场比赛。你的模型应该识别出哪位球员在比赛的特定时间表现得更好，以及他们的表现有多好。基于你的模型提供一个可视化来描述比赛流程。注意：在网球比赛中，发球的选手赢得分/局的概率要高得多。你可能希望以某种方式把这个因素考虑到你的模型中。
- 一位网球教练对“动量”在比赛中起作用持怀疑态度。相反，他认为一个球员在比赛中的摇摆和成功是随机的。用你的模型/指标来评估这一说法。

- 教练们很想知道是否有指标可以帮助确定何时比赛流程即将从有利于一名球员转变为另一名球员。
 - o 使用至少一场比赛提供的数据，开发一个模型来预测比赛中的这些波动。哪些因素似乎最相关(如果有的话)?
 - o 考虑到过去比赛中“动量”波动的差异，你如何建议一名球员进入一场与不同球员的新比赛?
- 在一场或多场其他比赛中测试你开发的模型。你对比赛中挥杆的预测有多好?如果模型有时表现不佳，你能识别出任何可能需要纳入未来模型的因素吗?你的模型对其他比赛(如女子比赛)、锦标赛、球场表面和其他运动(如乒乓球)的泛化程度如何?
- 用你的发现制作一份不超过25页的报告，并包括一到两页的备忘录，总结你的结果，并就“势头”的作用向教练提出建议，以及如何让球员准备好应对影响网球比赛过程的事件。

总页数不超过25页的PDF解决方案应包括: · 一页总结表。

- 目录表。
- 完整的解决方案。
- 一到两页的备忘录。
- 参考书目。
- [AI使用报告](#)(如果使用，不计入25页的限制)

注意:对于完整的MCM提交，没有特定的最低页数要求。您可以使用最多25页的总页数来完成所有解决方案工作和您想要包含的任何其他信息(例如:图纸，图表，计算，表格)。部分解决方案是可以接受的。我们允许谨慎地使用AI，如ChatGPT，尽管没有必要为这个问题创建一个解决方案。如果您选择使用生成式AI，则必须遵循[COMAP AI使用策略](#)。这将导致额外的AI使用报告，您必须将其添加到PDF解决方案文件的末尾，并且不计入解决方案的总页面限制25页。

提供的文件:

- [Wimbledon_featured_matches.csv](#) -温布尔登2023绅士单打第二轮后的比赛数据集。
- [data_dictionary.csv](#) -数据集的描述。
- [data_examples](#) -帮助理解所提供数据的示例。

术语表

大满贯:网球的大满贯是在一个日历年内赢得所有四个大满贯的成就。四项大满贯赛事是澳大利亚网球公开赛、法国网球公开赛、温布尔登网球公开赛和美国网球公开赛，每项比赛都持续两周。

关键术语/概念词汇表:

- **得分:[3]**
 - o **比赛:**五局四胜制(温布尔登绅士赛)
 - o **盘:**比赛合集;6局赢一局, 但选手必须先赢两局, 直到决胜局以6比6打平(见下文)
 - o **局:**收点;玩家达到4分时获胜, 但必须赢2分。见下文“一局得分”。
- **一场得分:[3]**
 - 0分=爱
 - 0 1分= 15分
 - 0 2分= 30分
 - 0 3分= 40分
 - o **平局=所有**(例如, “30 All”)
 - o 40 -40 =平分(玩家获得相同的分数, 每人至少3分)
 - o 服务器赢得一个平分点= Ad-in(或 “advantage in”)
 - o 接发球方赢得一分= Ad-out
- **发球:**玩家交替担任“发球者”(击球第一个点的玩家)和“接发球者”。在职业网球比赛中, 发球者往往有很大的优势。在每个点上, 球员有两次发球机会将球送入比赛(送入“发球箱”)。在两次发球尝试中发球失败是“双误”, 回发球的球员得分。
 - o **破发球**——当回击球员赢得一场比赛时。
 - o **破发点**——在这个点上, 如果接发球者获胜, 他们将赢得比赛。
 - o **守发球**——当发球球员赢得比赛时。
- **抢七:**每一局结束时, 只要一名球员赢了6局, 只要他们领先至少2局(即6 -4)。如果没有, 继续比赛, 直到6比6打成平手。此时进行决胜局。在温网比赛中, 抢七是先比7分(必须胜2分), 但在5th局比赛中, 先比10分(必须胜2分)除外。
- **休息时间/场边:**球员在第1场比赛后, 然后每两场比赛后更换场边。从第3场比赛开始, 每次换边时允许有90秒的休息时间。在决胜局中, 球员每6分换边。每盘结束后, 球员也要休息至少2分钟。允许医疗暂停和一次洗手间休息。

引用:

[1] Braidwood, J. (2023), Novak Djokovic has created a unique rival –is Wimbledon defeat the beginning of the end, The Independent, <https://www.independent.co.uk/sport/tennis/novak-djokovic-wimbledon-final-carlos-alcaraz-b2376600.html>.

[2] <https://www.merriam-webster.com/dictionary/momentum>

[3] Rivera, J. (2023), Tennis scoring, explained: A guide to understanding the rules terms & point system at Wimbledon, The Sporting News, <https://www.sportingnews.com/us/tennis/news/tennis-scoring-explained-rules-system-points-terms/7uzp2evdhbd11obdd59p3p1cx>.

帮助理解数据集的例子

例1:第5行

列(年代)	值(年代)	描述
match_id	“2023 -温布尔登- 1301”	“1301” 中的3表示第3轮匹配，“01” 表示从该轮列出的第一个匹配。本场比赛第一分开始后1分31秒，发球开始。
elapsed_time	“0:01:31”	
Point_no, game_no, set_no(“no” 是number的缩写)	4, 1, 1	本场比赛第一局第一局的第4分为该场比赛第一局的第4分。
Pl_sets, p2_sets, pl_games, p2_games	0, 0, 0, 0	由于这是比赛的第一局，双方都还没有赢过一局。
pl_score, p2_score	15, 30	玩家1的得分为15分，玩家2的得分为30分。因此，玩家1赢得1分，玩家2赢得2分。
服务器	1	玩家1 (Alcaraz)在这个点发球。
没有任何	1	这一分是在第一次发球中打出的，这意味着阿尔卡拉斯在比赛中打出了他的第一次发球。
point_victor	1	Alcaraz赢得这一点(玩家1)。
pl_points_won, p2_points_won	2, 2	玩家1 (Alcaraz)是积分胜利者，所以他的总分现在是2(之前是1)。对于玩家2来说，价值仍然是2，因为玩家2失去了积分。
game_victor, set_victor	0,0	阿尔卡拉斯赢得这一分使比赛得分为30比30(各2分)，因此双方在这一分上都没有赢过一局或一局(均为0)。
U -AC列		允许我们判断这一分是如何赢得的：
pl_winner	1	Alcaraz以一记“不可触碰”的击球拿下了这一分。
pl_ace	0	这个球不是发球(since=0)。
winner_shot_type	F	
p2_net_pt	1	击球是正手(而不是反手)。球员2(贾利)在比赛中靠近球网。由于阿尔卡拉斯赢得了这一分，虽然贾利在这一分时处于网前，但这个值为0。即使玩家2赢了这个点，比赛也不会结束，所以这个点不是“断点”，这些都是0。
p2_net_pt_won	0	
列AH-AM	均= 0	
Pl distance run, p2_distance_run	51.108,75.631	每个运动员在这一点上跑的距离(以米为单位)。
rally_count	13	两名球员在该点的总击球数。
Speed_mph, serve_width, serve_depth, return_depth	130, bw, ctl, d	Alcaraz(发球方)打出了接发球者的130发球 “Body/Wide”(我们之前看到这是第一次发球)，接近表示入局或出局的线。贾里(接发球者)在场上“深”回球(所以离球场另一端很近)。

例2:第8 -12行

第一局的最后4分说明了平局(deuce)和优势(ad)的概念。每一行都是比赛中随后的一个时间点。

行	列(年代)	值(年代)	描述
行8	pl_score, p2_score	40, 40	比分是40-40, 这意味着每个玩家都赢得了3分(这也被称为“平局”)。
行9	point_victor pl 分数, p2 分数	- 广告,40	阿尔卡拉斯赢得第7分(第8行)。 由于阿尔卡拉斯赢得了前一分(第7点), 第8点的得分现在是阿尔卡拉斯的“AD”, 杰瑞的分数是“40”, 这意味着阿尔卡拉斯又赢了一分, 可能会在下一分赢得比赛。
行10	point_victor pl_score, p2_score	2 40, 40	杰瑞(玩家2)赢得第8点(第9行)。 分数回到40-40(“平分”), 这意味着每个玩家在之前赢得了相同数量的分数, 尽管现在是每人4分。
行11	point_victor pl 分数, p2 分数	- 广告,40	Alcaraz赢得第9点(第10行)。 阿尔卡拉斯赢得了第9分, 再次占据优势。
行12	point_victor game_no pl_games	- 3 -	Alcaraz赢得第10点(第11行), 这意味着他赢得了比赛(现在又多了2分)。 这是第二局的第一分。 阿尔卡拉兹赢了第一场。

例3:第51行

比赛的第51点表示“破发点”, 即不发球的球员(正在回发球的球员)有机会赢得比赛的点。

行	列(年代)	值(年代)	描述
行51	Pl_score, p2得分	40, 30	比分是40比30, 意味着玩家1 (Alcaraz)领先。
	服务器	2	贾利(玩家2)发球。
	pl_break_pt	1	如果Alcaraz赢得了这一分, 他将赢得比赛;由于他没有发球, 这是一个“破发点”。
	点维克多	1	Alcaraz赢得了分数(因此赢得了游戏)。
	pl_break_pt_won	1	阿尔卡拉斯赢得了比赛, 并没有发球。

在COMAP竞赛中使用大型语言模型和生成式AI工具

这一政策的动机是大型语言模型(法学硕士)和生成AI辅助技术的兴起。该政策旨在为团队、顾问和评委提供更大的透明度和指导。这项政策适用于学生工作的各个方面,从模型的研究和开发(包括代码创建)到书面报告。由于这些新兴技术正在迅速发展,COMAP将适当地完善这一策略。

团队必须公开和诚实地使用AI工具。一个团队及其提交的内容越透明,他们的工作就越有可能得到他人的充分信任、赞赏和正确使用。这些披露有助于理解智力工作的发展和对贡献的适当承认。如果没有对AI工具作用的公开和清晰的引用和参考,那么有问题的段落和工作更有可能被认定为抄袭并被取消资格。

解决这些问题不需要使用AI工具,尽管允许负责任地使用它们。COMAP认识到法学硕士和生成AI作为生产力工具的价值,可以帮助团队准备提交;例如,为一个结构产生初步的想法,或者在总结、释义、语言润色等时。在模型开发的许多任务中,人类的创造力和团队合作是必不可少的,对AI工具的依赖会带来风险。因此,我们建议在将这些技术用于模型选择和构建、协助创建代码、解释模型的数据和结果以及得出科学结论等任务时要谨慎。

值得注意的是,法学硕士和生成式AI有局限性,无法取代人类的创造力和批判性思维。COMAP建议团队在选择使用法学硕士时要意识到这些风险:

- 客观性:法学硕士生成的文本中可能出现先前发表的包含种族主义、性别歧视或其他偏见的内容,一些重要观点可能未被代表。
- 准确性:法学硕士可能会产生“幻觉”,即产生虚假内容,特别是在他们的领域之外使用或处理复杂或模棱两可的主题时。他们可以生成语言上但科学上不合理的内容,他们可以错误地获取事实,并且他们已经被证明可以生成不存在的引用。一些法学硕士只接受特定日期之前发布的内容的培训,因此呈现的是不完整的画面。
- 语境理解:法学硕士不能将人类的理解应用到一篇文章的语境中,特别是在处理习惯用语、讽刺、幽默或隐喻语言时。这可能会导致生成的内容出现错误或误解。
- 训练数据:法学硕士需要大量高质量的训练数据来达到最佳性能。然而,在某些领域或语言中,这样的数据可能并不容易获得,从而限制了任何输出的有用性。

对团队的指导

参赛队伍需要:

1. 在**报告中明确指出使用了法学硕士或其他人工智能工具**，包括使用了哪个模型以及用于什么目的。请使用内联引文和参考文献部分。在你的25页解决方案之后，还要附上人工智能使用报告(如下所述)。
2. **验证内容的准确性、有效性和适当性**以及由语言模型生成的任何引用，并纠正任何错误或不一致之处。
3. **提供引用和参考文献，遵循这里提供的指导**。仔细检查引文，以确保它们是准确的，并被正确引用。
4. **要注意抄袭的可能性**，因为法学硕士可能会从其他来源复制大量文本。检查原始来源，以确保你没有抄袭别人的作品。

COMAP将采取适当的行动，当我们确定提交可能准备与未公开使用这些工具。

引文和参考说明

仔细考虑如何记录和引用团队可能选择使用的任何工具。各种风格指南开始纳入引用和参考人工智能工具的政策。在你的25页解决方案的参考部分，使用内联引用并列出所有使用的人工智能工具。

无论团队是否选择使用人工智能工具，主要解决方案报告仍然限制在25页。如果一个团队选择使用人工智能，在你的报告结束后，添加一个名为人工智能使用报告的新部分。这个新章节没有页数限制，不会被计入25页的解决方案中。

例子(这不是详尽的-根据你的情况调整这些例子):

人工智能使用报告

1. OpenAI *ChatGPT* (Nov 5, 2023 version, ChatGPT-4) Query1: *<insert the exact wording you input into the AI tool>* Output: *<insert the complete output from the AI tool>*
2. OpenAI *Ernie* (Nov 5, 2023 version, Ernie 4.0)
Query1: *<insert the exact wording of any subsequent input into the AI tool>*
Output: *<insert the complete output from the second query>*
3. Github *CoPilot* (Feb 3, 2024 version)
Query1: *<insert the exact wording you input into the AI tool>* Output: *<insert the complete output from the AI tool>*
4. Google *Bard* (Feb 2, 2024 version) Query: *<insert the exact wording of your query>* Output: *<insert the complete output from the AI tool>*