

《心理学报》审稿意见与作者回应

题目：系列决策任务中的策略转换：来自爱荷华赌博任务的证据

作者：胡馨允，沈悦，戴俊毅

第一轮

审稿人 1 意见：

本文汇报了两个研究，一个是基于之前研究的数据，另一个数据是作者自己收集的。两个研究类似，任务都是 IGT，都在比较相同的认知模型，得出的结论也比较一致。本文最大的贡献是提出并检验了决策者在 IGT 中出现策略转换的可能性，并发现考虑策略转换的模型可以拟合超过一半被试的数据。作者对决策任务、所检验的模型和模型对比结果都汇报的比较清晰，而且通过两个研究得到看似比较稳健的结果。但本人对模型对比方法和结论有些疑问，供作者参考。

意见 1：

AIC vs. BIC. 为了避免过分拟合，两种方法都惩罚参数多的模型，不同的是，BIC 对复杂模型的惩罚力度更大，在数据量不大、过分拟合更容易出现的情境下更为适用。本文研究中大多被试完成 100 个试次的任务，WSLS 有 2 个参数，试次-参数比为 50；SSO 有 7 个参数，试次-参数比大概为 14。按照这个过分简单的指标，SSO 更可能过分拟合数据。因此，从谨慎的角度出发，BIC 看上去是个更为合适的模型比较标准。此外，BIC 假设有唯一正确的模型，目的是最大可能找到这个模型；而 AIC 没有这个假设，其目的只是找到能最大程度解释数据的模型。因此，从寻求“真理”的角度出发，BIC 也更合适。

回应：

关于该选取何种评价指标作为模型比较和选择的依据，一直是统计学界争论的话题。有鉴于此，我们在本研究中同时考虑了 AIC 和 BIC 这两种最为常用的模型评价指标，并且使用模型复原测试的结果，来判断就本研究涉及的问题而言，何种指标更为合适。为了使得本研究的结论更为可信，我们在论文修改阶段，又对以往的相关文献进行了系统的回顾，尤其是以下两篇直接对比 AIC 和 BIC 且在心理学和社会科学界影响广泛的论文 (Burnham & Anderson, 2004; Vrieze, 2012)。综合这些论文的观点，我们认为在本研究中，使用基于信息论的 AIC (以及 AIC_C) 指标作为模型选择的依据，要比使用 BIC 指标 (虽然该指标也被冠以信息准则之名，但其实最初是在贝叶斯框架下作为贝叶斯因子的近似指标被提出的，Burnham & Anderson, 2004) 更为合理，主要理由包括：

a. 本研究考察的 IGT 是一个相当复杂的任务。一方面，任务本身带有随机性，主要表现为选择任意一个牌堆会引发的结果都具有不确定性；另一方面，个体的选择可能受到众多因素的影响，而且存在多种可能的决策策略 (及其转换)。因此，我们需要寻找或者逼近的所谓“真实情况”或者“真理”，大概率具有相当复杂的形式，甚至不能以参数化的模型来表达。因此，我们实际考虑的这些模型，只能是对于真实情况的近似，而几乎不可能是所谓的“真模型”。

b. 要让 BIC 能够确保选择出真模型需要满足两个重要前提，首先真模型需要包含在所考虑的模型集中，其次必须有足够大的样本量。就绝大多数现实条件下的模型选择任务而言，这两点都无法满足，对当前有关 IGT 的研究也是如此。此外，根据 Vrieze (2012) 的观点，

要让 BIC 能够以概率 1 选择出真模型，还需要真模型的模型参数维度随着样本量增加而保持不变，且真模型的参数个数是有限的。这两条就 IGT 而言也难以成立。

c. 当样本量较小，或者随着模型参数个数的增加，更加复杂的模型能够解释的额外效应逐渐减小时 (tapering effects)，使用 BIC 作为模型选择指标往往会导致欠拟合 (underfit) 的结果 (Burnham & Anderson, 2004)。

d. 当样本量与模型参数个数的比值较小 (样本量/参数个数 <40) 时，可以使用包含二阶偏差修正的 Akaike 信息准则 (AIC_C) 来弥补使用 AIC 可能导致的过拟合缺陷 (Burnham & Anderson, 2004)。

e. AIC 所要选取的模型，是在给定当前样本量的前提下，能够最好地平衡预测变异性 (variance) 和预测偏差 (bias) 的模型，而非真模型。正如评审专家所言，真实的策略转换，很可能要比我们假定的一次策略转换模型更加复杂。所以，我们当前分析的目的，并非找到能够完美表达这种真实的策略转换的模型，而是表明，包含策略转换假设的模型，要比单一策略模型能够更好地解释数据，从而为在 IGT 中存在策略转换这一观点提供支持。我们和评审专家一样，并不认为我们能够通过当前的分析就找到真模型。

基于以上五点，我们认为在本研究中，使用 BIC 并不能够帮助我们找到所谓的“真模型”，但却可能导致出现欠拟合的情况。另一方面，通过将 AIC_C 而非 AIC 作为模型选择和选择的标准，大概率可以避免至少缓解由于使用 AIC 而可能导致的过拟合的问题。因此，在修改过程中，我们使用 AIC_C 作为指标重新分析了研究数据。总体而言，使用 AIC_C 指标进行小样本修正后，模型比较的结果与使用 AIC 并无太大差别，SSO 模型仍然能够最好地拟合至少一半被试的数据，并且 SSO 模型的表现随着试次数的增加也会变得更好。(在修改稿中，参考另一位评审专家的建议，我们把假定不同转换方向的两个 SSO 模型，组合成了一个统一的模型，以便更好地将策略转换模型和单一策略模型进行比较) 此外，即便使用 BIC 作为模型选择的指标，本研究的两项核心结论 (即在完成 IGT 的过程中，有部分被试会发生策略转换，以及随着试次数的增加，发生策略转换的个体比例会有所上升) 依然成立。具体而言，在研究一中，SSO 模型最优拟合的被试比例由 95 试次的 0%，上升到 100 试次的 26.58%，再上升到 150 试次的 47.96%；在研究二中，SSO 模型最优拟合的被试比例由 100 试次的 28.76%，上升到 200 试次的 45.97%。因此，基于以上分析和论证，我们认为使用 AIC_C 作为模型拟合和比较的指标更为合理，且模型选择指标的选取，并不会从本质上改变我们的研究结论。

意见 2:

模型复原度测试。但为什么 BIC 在模型复原度测试中表现的如此之差，不能区分模型呢？这个结果比较少见，对我也比较 puzzling。结合前一点，我想可能是如下原因造成的：a) 用 AIC 选择出来的模型更容易过分拟合，这是为什么在 AIC 下，参数更多的 VPP 和 SSO 拟合结果明显更好 (表 2) b) 基于 AIC 标准筛选出的模型复原测试，是用可能存在过拟合的模型重新生成模拟数据，所生成的数据本身可能就偏离了真实过程，导致这些数据也更容易用该模型去拟合，而不是用其他模型，所以 AIC 对不同模型的区分度也更高；c) 基于 BIC 的模型复原测试，最简单的 WLS 对所有重新生成数据的拟合结果都最好，说明 WLS 是最 robust 的模型，其他复杂模型都存在一定程度的过拟合。如果真是这样的话，本研究的结论应该受到质疑，SSO 表现好可能不是因为它是那个更准确捕捉被试心理过程的模型，而且那个拟合数据程度最好的模型。

回应:

关于在模型复原测试中使用 BIC 作为模型选择指标所出现的情况，的确不太常见，也和作者以往针对其他任务的研究结果有所不同。但是，需要指出的是，模型复原测试会先使

用特定模型生成模拟数据，再使用模型集中的所有模型进行拟合。换言之，对于模拟数据而言，真模型是存在于所考虑的模型集之中的，即使这些产生模拟数据的真模型可能对于真实的观测数据存在过拟合的情况。因此，在较为理想的情况下（即样本量足够大时），针对这些模拟数据，使用 BIC 作为模型选择指标应当有很大的概率能够选择出产生数据的模型。然而，我们的模型复原测试结果却表明，此时使用 BIC 作为指标，当真模型是更为复杂的模型时，最有可能选取出的仍然是 WLS 模型，即出现了欠拟合的情况。正如上一条回复所言，这是当样本量较小或者更复杂模型能够解释的效应逐渐减小时，使用 BIC 常会出现的一种缺陷。换言之，使用 BIC 作为模型选择指标进行模型复原测试出现的结果，更可能表明的是在本研究中不太适合使用 BIC 来对观测数据进行分析，或者说，如果使用 BIC 指标来选取模型的话，很可能选取出的模型会过于简单，不能较好地反应产生观测数据的真实机制。

对于评审专家提出的解释，我们也进行了认真的思考。但是，从我们的角度来看，评审专家提出的解释，似乎更适用于在训练集和测试集之间的比较场景，而非模型复原测试的场景。这两类场景的主要区别在于，当探讨不同模型在训练集和测试集上的表现差异时，训练集和测试集都是由真实的数据产生机制产生的，而在模型复原测试中，需要再次进行分析的数据（对应于前一场景中的测试集），是由特定的模型而非产生真实观测数据的机制产生的。换言之，训练集和测试集的数据产生机制，是我们所未知的。此时，如果某个模型在训练集上表现更好，可能是由于它不仅拟合了真实机制所包含的规律性的成分，也拟合了噪声成分，即出现了过拟合的情况。用这样的模型去预测测试集上的数据，则可能会出现表现较差的结果，因为该模型假设的机制和真实机制存在差异，前者的预测还反映了测试集中的噪声成分。因此，如果某个模型能够在测试集中表现得更好，那么它就有可能更能反映真实机制的模型。但是，对于模型复原测试而言，当使用更复杂的模型产生了模拟数据时，这些数据背后的真模型就是那些更复杂的模型。因此，如上一段论证的那样，不应该出现其他更简单模型能够更好拟合的情况。如果出现了这种情况，更可能的解释是所选用的模型选择指标会导致欠拟合的情况。

最后，如上一条回复所言，即便使用 BIC 作为模型选择指标，本研究的两条核心结论，仍然是成立的。因此，似乎不必太过纠结于是使用 AIC_c 会导致过拟合的情况，还是使用 BIC 会导致欠拟合的情况这一问题。

意见 3:

SSO。这是作者提出的新模型，主要参数是 sp ，那个策略转化的点。有几个问题：

a) 一个如此重要的参数，为什么不在文中汇报估出的参数值，让读者知道被试大概在哪个时间点转换策略？

b) 正如作者所指出的，假设一个策略转化点过于简单，现在我们只能知道在 sp 之前的所有试次中被试大体使用什么策略，之后的所有试次中大体使用什么策略，这个太过粗糙，不能体现真正的心理过程；

c) 为什么有的被试先用 RL 后用 WLS，而其它被试则相反？有没有可能也是过分拟合的结果？

回应:

感谢评审专家的建议，在修改稿中，我们对 sp 的分布情况进行了介绍，并且比较了不同试次数下 sp 分布的差异，详细情况请见修改稿中标定的内容。

我们当前的研究通过构建一次策略转换模型的方式探究了人们在 IGT 任务中发生策略转换的可能性。和评审专家一样，我们并不认为这样的模型就能够反映在 IGT 中所发生的策略转换的全部情况。但是，如果这样的一次策略转换模型能够比假设单一策略的模型更好

地解释数据，那么还是可以被认为是支持了在 IGT 中存在策略转换的。正如我们在总讨论部分中提到的那样，这一模型所假设的仅发生一次转换以及以突变的方式发生转换这两点都存在局限性。在后续研究中，我们将会使用更加前沿更加合理的方式，例如贝叶斯方法和机器学习方法等，探究包含渐变的、多次策略转换的认知计算模型，以期对于 IGT 以及类似任务中的策略转换状况产生更加深入的认识。

关于为什么有的被试会先使用强化学习策略，后使用启发式策略，而另一些被试则正好相反这一点，我们也进行了一定的思考。如我们在初稿中暗示的那样，产生前一种转换的原因，可能是人们在初始阶段使用了强化学习策略，并随着任务的进行，因为疲劳、倦怠或者降低认知负荷的需求，转而采用启发式策略。产生后一种转换的原因，则是人们在任务开始阶段由于缺乏信息而使用对信息依赖度较低的启发式策略，并在对各牌堆有了更多了解之后，转而使用更为复杂更为精细的强化学习策略。对于评审专家提出的，这一分析结果是由过拟合所致的这种可能性，我们认为的确存在，即本质上并不存在策略转换，或者只存在一种方向的策略转换，但由于数据中的噪声，导致了特定的模型选择结果。但是，如上文所述，使用 AIC 可能产生的过拟合的问题，应该可以通过使用 AIC_C 加以克服，但我们使用 AIC_C 的结果与使用 AIC 的结果高度一致。鉴于此，我们认为相关分析所得的结论，更有可能反映的是真实发生的过程，而非仅仅是过拟合的结果。未来的研究可以通过实验操控的方式，考察在特定任务条件下，模型比较的结果和所预期的策略转换类型是否一致，从而为上述解释提供进一步的支持证据。

意见 4:

虽然有这些问题，但我认为本研究的目的还是可取的。假设被试在一个实验的所有任务中使用同一个策略是个不现实的假设，早就应该想办法纠正。除了 Lee 等人的研究，也可以看一下这篇：<https://link.springer.com/article/10.3758/s13428-022-01828-1>

回应:

非常感谢审稿专家的建议，我们已阅读了相关文献，并且在总讨论中，将该文献中提出的新的研究范式作为将来研究的一种可能方向加以介绍。详细情况请见修改稿中标定的内容。

审稿人 2 意见:

当前研究在爱荷华决策任务的情境中借助计算模型探讨了一个非常重要的问题：即使是在环境不变时，同一个体在决策过程中采用的认知策略是否也可能随着时间的推移或经验的累积发生改变？作者们为爱荷华决策任务中人们的决策行为构建了五个不同的计算模型。其中，前三个模型前人研究中所提出的代表性模型，它们都假定同一被试采用的认知策略在实验过程中始终不变，分别对应着启发式策略（WLSLS 模型）、强化学习（PVL2 模型）和二者的混合（VPP 模型）。后两个模型为本文新提出，假定同一被试采用的认知策略在实验过程中会在某个时间点发生一次转变，从启发式策略转向强化学习（RL-H 模型），或者反之（H-RL 模型）。作者们首先重新分析了总共包含 600 余名被试的多个前人数据集，比较了五个模型在每个被试的决策行为数据上的拟合表现，发现 RL-H 和 H-RL 模型加起来在超过一半的被试中表现超过其他模型；他们之后又采集了一个包含 300 余名被试的新实验对结论进行了进一步检验。作者们的另一发现是，“随着试次数的增加，发生策略转换的可能性也会上升。”

这篇文章有着独到见解和创新视角，研究逻辑清晰，方法严谨，语言简洁优美，近乎完美，让我在评阅过程中也收获良多。

不过，文章的一个重要结论（“随着试次数的增加，发生策略转换的可能性也会上升”）

所依赖的证据似乎难以排除 artefacts，需要再斟酌用词削弱结论并进行相关讨论。作者们在得出这个结论时依据的实验证据是，当实验包含 200 个试次时，RL-H 和 H-RL 模型占优的被试比例高于 100 个试次条件(见表 6，总共高了 13.28%)。我觉得这个结果可能包含 artefacts，有如下两点依据。

意见 1:

(1) 作者们同时也发现，100 试次和 200 试次条件下模型复原的结果之间存在差异，在 200 试次条件下 RL-H 和 H-RL 模型被复原的比例也高于 100 试次条件(见表 7，总共高了 25.87%，甚至比实际观察到的 13.28% 的差异高)。换言之，即使两种实验条件下发生策略转换的被试比例没有差异，因为模型复原的差异，也会观察到 200 试次条件采用 RL-H 或 H-RL 的被试比例比 100 试次条件高。

回应:

感谢评审专家指出我们初稿中存在的这个结果解读方面的问题。诚如评审专家所言，之所以在不同试次数下 SSO 模型能够最好解释的个体比例有所差异，除了我们提供的解读外，也可能是由于当使用特定模型选择指标时，最接近于个体真实数据产生机制的模型，未必会被正确地选取，且在不同试次数下出现此类错误的可能性也不相同。为了更好地认识这一问题，我们重新分析了观测和模拟数据。具体而言，我们首先如评审专家下文中建议的那样，将假设不同转换方向的两个 SSO 模型，合并为一个统一的一次策略转换模型，以便更好地将策略转换模型和单一策略模型进行比较。以下我们会称这一模型为 SSO 模型，该模型共包含 8 个参数，多出的一个参数表达的是转换类型。然后，我们将该 SSO 模型和 3 个单一策略模型 (WSLS, PVL2 和 VPP) 组成模型集并进行了相应分析。考虑到 AIC 可能出现另一位评审专家指出的过拟合的问题，我们在分析时使用了更加适合于当前数据情况的包含二阶修正的 AIC_C 替代 AIC 作为模型比较和选择指标(根据 Burnham 和 Anderson (2004)，当样本量与模型集中包含最多参数模型的参数个数的比值小于 40 时，需要使用包含二阶修正的 AIC_C 代替 AIC)。采用新的模型集和模型选择指标后，在研究 2 的 100 试次条件下，SSO 模型在 50.00% 的个体数据上有最好的表现，当试次数上升为 200 时，这个比例上升到了 65.22%，仍然显著高于 100 试次下的比例。另一方面，在模型复原测试中，SSO 模型的复原比例在 100 试次条件下是 77.29%，在 200 试次条件下是 83.85%，差值为 6.56% (见更新后的表 7)，明显小于观测数据中的比例差值 (即 15.22%)。因此，观测数据中的差异，不太可能主要是由于模型复原能力的差异所致。

意见 2:

(2) 从表 6 来看，200 试次条件比 100 试次条件的 H-RL 被试比例高了 11%。如果是因为随着实验试次的增多，原先采用启发式策略 (WSLS) 的被试逐渐改为采用强化学习，那么，我们会预期这部分增加的比例来源于 WSLS 的比例的减少。然而，从 100 到 200，WSLS 的比例只减少了 4.4%，远低于 11%。这再次指向 100 和 200 间的差异可能是跟模型复原有关的 artefacts。

回应:

基于更新后的数据分析，100 试次下 WSLS 模型拟合最佳的个体数据比例为 10.63%，200 试次下为 2.48%，下降了 8.15%。100 试次下 PVL2 模型拟合最佳的个体数据比例为 16.88%，200 试次下为 9.32%，下降了 7.56%。100 试次下 VPP 模型拟合最佳的个体数据比例为 22.50%，200 试次下为 22.98%，基本持平 (见更新后的表 6)。因此，当试次数由 100 变为 200 时，SSO 模型的比例上升，几乎正好对应于 WSLS 和 PVL2 模型的比例下降。而且，两类转换人数比例的上升幅度 (WSLS→PVL2 上升 9.81%，PVL2→WSLS 上升 5.41%) 也与对应单一

策略模型人数比例下降的幅度大致相当。因此，之前使用两个 SSO 模型以及 AIC 所得结果会暗示存在的问题，在当前更加合理的分析方式下，已不复存在。换言之，试次上升后 SSO 表现的提升，不太可能是由 artefacts 所致。

意见 3:

此外，我还有一两个小建议如下，供参考：

因为文章希望得出的结论是关于策略转换的，RL-H 和 H-RL 属于同一类模型，如果能给出一个关于这类模型的统一结果而不只是将二者的最好模型比例加起来，结论可能会更具说服力。具体做法是，将 RL-H 和 H-RL 作为一个模型，增加一个 dummy variable 来控制 RL-H or H-RL。这个统一的策略转换模型可以与 WSLs、PVL2 和 VPP 直接进行比较，还可以应用 group-level Bayesian model selection 计算 protected exceedance probability。（这个模型不需要重新拟合，计算 AIC 或 BIC 时将 RL-H 和 H-RL 的 LL 取最大值，并多考虑补偿一个参数就可以。）

回应:

非常感谢评审专家关于合并模型的建议，在修改稿中我们已像建议的那样，将之前的两个 SSO 模型统一为一个模型，并相应地增加了一个代表转换类型的模型参数。

关于计算 protected exceedance probability (PEP) 的建议，我们也仔细阅读了 Rigoux 等人 (2014) 这一相关文献。根据该文提供的定义，PEP 指的是在群体层面，某个模型优于其他所有模型（即有最大比例的个体数据由该模型产生）的概率。我们的研究主要考察的是策略转换在 IGT 中发生的可能性，因而我们主要关注的是策略转换模型在人群中占优的比例，但并不假设策略转换模型一定是所有考虑的模型中占优比例最高的模型。基于以上分析，我们认为 PEP 并不是我们所需要的指标，但出自同一分析框架的群体层面的模型后验概率，则能够为我们的分析提供更加合理的指标。因此，我们使用 Rigoux 等人介绍的方法计算了群体层面的模型后验概率，结果见下表。

模型	研究一			研究二	
	95 试次	100 试次	150 试次	100 试次	200 试次
WSLS	0.2916697	0.3317703	0.2374186	0.3893717	0.1813083
PVL2	0.6765860	0.3698538	0.2637892	0.3360015	0.3867938
VPP	0.0158806	0.1568323	0.2366746	0.1087029	0.1124763
SSO	0.0158637	0.1415436	0.2621176	0.1659239	0.3194217

从表中的结果可以看出，在各种试次数条件下，SSO 模型的后验概率都大于 0，而且随着试次数的增加，SSO 模型的后验概率也在变大，即有更高比例的被试可能使用了 SSO 模型假设的机制产生了个体数据。这一结论同正文中报告的模型分析结果一致。同时我们也发现，SSO 模型最大的后验概率仅仅只有 31.94%（研究二 200 试次组），该比例低于以 AIC_C 作为模型选择指标所得的结果。产生这一结果的可能原因是使用了以下公式通过 BIC 计算了模型证据比（即 BF）：

$$BF_{01} \approx \frac{\Pr_{BIC}(D|H_0)}{\Pr_{BIC}(D|H_1)} = \exp\left(\frac{\Delta BIC_{10}}{2}\right) \text{ (Wagenmakers, 2007)},$$

再借助计算所得的模型证据比转换为（数学上等价的）模型证据输入相应的分析程序，进而计算得出模型的后验概率。因为在这一分析中使用了 BIC 指标，因此对应结果与使用 BIC 作为模型选择指标的结果类似。需要指出的是，通过 BIC 来计算 BF 是一种近似做法，严格说仅在大样本条件下成立，在当前研究中并不是十分适用，因此该结果的可靠性还需要进一步考量。当然，无论是当前的分析，还是正文中报告的以 BIC 或者 AIC_C 为模型选择指标进行的分析，都表明 SSO 在部分被试身上有最好的表现，而且对应的个体比例随着试次数的上升会增大。换言之，无论使用何种分析思路，本论文的核心结论不变。因此，我们在修改稿中，仍然

以使用 AIC_C 得到的结果为论述重点。

意见 4:

第一段中说,“存在一系列不同的补偿式和非补偿式策略”。不太明白“补偿式和非补偿式”分别是什么。可以用几个字解释一下吗?

回应:

当人们使用补偿式策略进行判断和决策时,选项在某一属性上的劣势可以被该选项在其他属性上的优势所补偿。例如,根据 WADD 策略,人们会将每个选项的各个属性根据重要性进行加权平均从而做出选择。在这一策略下,选项在某些属性上的劣势,可能被其他属性上的优势所抵消。相反,当使用非补偿式策略时,人们往往只会关注部分信息而忽略其他信息,因此某些属性上的劣势,无法被其他属性上的优势所补偿。例如,根据 take-the-best 策略,人们只会比较两个选项最重要的线索,当这一线索在选项间存在差异时,就会做出判断或者决策。因此,其他未被考虑的属性上的优势,无法补偿有区分度且已被考虑的属性上的劣势。由于相关内容不是本文探讨的重点,所以在修改稿中我们仅增加了一些简单的文字对这两类策略做出解释,详细内容请见修改稿。

参考文献

- Burnham, K. P., & Anderson, D. R. (2004). Multimodel inference: Understanding AIC and BIC in model selection. *Sociological Methods & Research*, 33(2), 261–304. <https://doi.org/10.1177/0049124104268644>
- Rigoux, L., Stephan, K. E., Friston, K. J., & Daunizeau, J. (2014). Bayesian model selection for group studies—revisited. *Neuroimage*, 84, 971–985.
- Vrieze, S. I. (2012). Model selection and psychological theory: A discussion of the differences between the Akaike information criterion (AIC) and the Bayesian information criterion (BIC). *Psychological Methods*, 17(2), 228–243.
- Wagenmakers, E. J. (2007). A practical solution to the pervasive problems of p values. *Psychonomic Bulletin & Review*, 14, 779–804. <https://doi.org/10.3758/BF03194105>

第二轮

审稿人 1 意见:

意见 1:

感谢作者对我上一轮建议和问题的回复,以及对文章进行的相应修改。文章修改后质量提高了很多,由其是用 AIC_C 替代 AIC 作为模型比较的指标部分,让研究的结果更为可信可靠。现在我只有一个问题,关于新增的有关参数 sp 的结果:不论是研究一还是研究二,当实验试次数增加时(如从 100 到 200),被估的策略转换点(sp)的值也相应提高。虽然作者给了一个简单的解释:“之所以会出现这一状况,可能是由于随着 IGT 试次数的增多,有更高比例的被试在整个任务过程中发生了策略转换,且新增的发生策略转换的被试的转换节点较晚。”且不说这个解释是否合理,让我更为好奇的是,为什么策略转换点的均值与实验试次数的中间值非常接近?在研究一中,当总试次数为 95 时,sp 估计值的均值为 48.5;当总试次数为 100 时,sp 估计值的均值为 48.92;当总试次数为 150 时,sp 估计值的均值为 81.42;在研究二中,当总试次数为 100 时,sp 估计值的均值为 47.03;总试次数为 200 时,sp 估计值的均值为 95.38。出现一次可能是巧合,但在所有的条件下都出现,巧合的概率应该不大。目前我自己还想不出一个好的理由来解释这个结果,希望作者能更详尽地看一下具体的数据,

进而给一个合理的解释。不然的话,可能会让人质疑这个 sp 参数值是个有心理意义的数值,还是模型拟合中的一个 *statistical artifact*?

回应:

感谢评审专家对我们所做修改的肯定,以及提出的这一重要的新问题。针对这一问题,我们进行了一系列的分析,以便对于当前考虑的 *SSO* 模型的性质有更加深入的认识。具体而言,对于修改稿中报告的存疑结果,即 *SSO* 模型拟合最优的个体数据对应的 sp 估计值的均值接近试次数中间值这一点,我们设想了两种可能的 *statistical artifact*,以及另一种相对合理地产生这一结果的机制。首先,导致估计均值接近中间值的一种可能的 *artifact*,是当使用 *SSO* 模型拟合任意个体数据时,所得的 sp 估计值并不能反映个体真实的策略转换节点位置,而是一个完全随机的变量,且该随机变量在该参数的允许范围内总是满足均匀分布。为了检验这种可能性,我们分析了当使用 *SSO* 模型拟合 *PVL2* 或 *WSLS* 模型拟合最优的个体数据时的 sp 估计值分布。如果如上所述,使用 *SSO* 模型拟合任意个体数据时,所得的 sp 估计值是一个在允许范围内满足均匀分布的随机变量,那么无论是 *PVL2* 拟合最优的被试还是 *WSLS* 拟合最优的被试(即更有可能使用了对应的单一策略的被试),当使用 *SSO* 模型拟合他们的数据时,在总体层面 sp 的估计值的分布也应该是均匀的。因此, sp 估计值有 50% 的可能高于中间值,有 50% 的可能低于中间值。相反,如果 *SSO* 模型拟合这些个体数据时得到的 sp 估计值能够较好地反映策略转换节点位置,那么对于 *WSLS* 模型拟合最优的个体,当 *SSO* 模型拟合的结果表明策略转换类型为强化学习到启发式策略(即 $st=1$)时, sp 的估计值应该在允许的下边界附近(即低于中间值),从而使得大多数试次是在 *WSLS* 策略下完成的。相反,如果 *SSO* 模型的拟合结果表明,策略转换类型为启发式到强化学习策略(即 $st=2$)时, sp 的估计值应该在允许的上边界附近(即高于中间值)。

基于这一推断,我们根据用 *SSO* 模型拟合 *WSLS* 模型拟合最优的个体数据得到的 sp 的估计值,生成了一个新的变量 S 。具体而言,当 $st=1, sp < \text{中间值}$ 或者 $st=2, sp > \text{中间值}$ 时,变量 S 赋值为 1,当 $st=1, sp > \text{中间值}$ 或者 $st=2, sp < \text{中间值}$ 时,变量 S 赋值为 0。如果对 sp 的估计值是由上述 *artifact* 所致,那么 S 取 1 的可能性应当等于 S 取 0 的可能性。相反,如果 sp 的估计值具有合理的心理意义,那么 S 取 1 的可能性应该高于 S 取 0 的可能性。对于 *PVL2* 模型拟合最优的那些被试,则应当出现相反的模式。我们的分析结果表明,对于 *WSLS* 拟合最优的那批被试的个体数据, S 取 1 的样本比例为 55%,高于 50%,单侧 $p=0.259$ 。对于 *PVL2* 拟合最优的那批被试的个体数据, S 取 1 的样本比例为 23.08%,低于 50%,单侧 $p < 0.001$ 。换言之,样本统计量不支持上述 *artefact*,且对于 *PVL2* 拟合最优的个体数据,统计检验的结果支持 *SSO* 模型的 sp 参数估计值具有心理意义。因此,可以排除用 *SSO* 模型拟合任意个体数据得到的 sp 是一个在允许范围内满足均匀分布的随机变量这种 *statistical artefact* 的可能性。

第二种可能导致 *SSO* 模型拟合最优的个体数据对应的 sp 估计值的均值接近试次数中间值的 *statistical artifact*,是仅当 *SSO* 模型拟合最优时,对应的 sp 的估计值满足均匀分布。针对这种可能性,我们使用 *KS* 检验分别考察了研究 1 和研究 2 中各种试次数条件下,*SSO* 模型拟合最优的个体对应的 sp 估计值的分布是否为均匀分布。当试次数为 100 时,无论是研究 1 还是研究 2,对应的 sp 估计值的分布都在统计上显著地偏离均匀分布(研究 1, $p=0.009$,研究 2, $p=0.002$)。对于研究 1 中 150 试次的的数据, $p=0.187$,对于研究 2 中 200 试次的的数据, $p=0.071$ 。虽然两者未达到统计显著程度,但综合上述分析的结果,可以认为至少在 100 试次下,*SSO* 拟合最优的个体数据对应的 sp 的分布不是均匀的,也就是说上述第二种 *statistical artefact* 存在的可能性也不大。

虽然上述分析排除了 sp 估计值不能反映真实的策略转换节点,而只是一个满足均匀分布的随机变量这一可能,但是我们仍需解释为什么会出现,*SSO* 模型拟合最优的个体数据

对应的 sp 估计值的均值接近试次数中间值这一观测结果。除了均匀分布这一可能，对称分布的 sp 估计值也可能导致这一结果。由于主流的统计分析软件无法检验满足对称分布这一假设，我们设计了一个新的统计检验方法来考察这一问题。具体而言，我们将 sp 的允许取值范围划分为若干个各包含 5 个相邻可能取值的区间，且保证在中间值以下的区间数等于在中间值以上的区间数。随后，我们计算了 sp 估计值在每对关于中间值对称的区间中实际出现次数的差值（记作 D ， $D > 0$ 代表左侧频次高于右侧频次）。如果 sp 在总体中满足对称分布，那么 D 的取值应当为 0。考虑到 D 的样本值会包含误差， D 的样本值不一定为 0，但其抽样分布的平均值应该为 0。研究 1 包含了来自过往多个研究的数据，且这些研究的具体实验设定各有不同，这可能导致转换节点的分布具有异质性，因而不适合进行汇总分析。此外，研究 1 中完成 95 试次（15 人）和 150 试次（98 人）的被试人数较少，由 SSO 模型最优拟合的个体数据更少，这可能导致统计检验效力不足。相反，研究 2 是我们新开展的研究，除了试次数外，不同被试接受的处理是一致的，且被试人数较多（ ≥ 160 人），更为适合进行当前分析。因此，我们对研究 2 的数据进行了分析，发现无论是 100 试次还是 200 试次， D 的平均值都大于 0，意味着 sp 的估计值更有可能小于中间值（研究 1 中完成 100 试次 IGT 的被试的 D 的平均值也大于 0）。针对 D 的平均值为 0 这一零假设，单样本 t 检验的结果为 $m = 1.136$ ， $t = 1.869$ ，Cohen's $d = 0.394$ ，单侧 $p = 0.04$ 。由此，我们推测 sp 的估计值在中间值以下的可能性略高于在中间值以上的可能性，这将导致对于 SSO 拟合最优的个体数据， sp 的平均值接近于中间值，但前者略小于后者。

之所以 sp 的估计值在中间值以下的可能性略高于在中间值以上的可能性，有可能是由于当试次数足够多时，个体转换节点的分布满足单峰形态，但是众数位置的可能性，并不显著高于其他位置的可能性。当 IGT 的试次数有限时，仅有那些转换节点在 SSO 模型允许范围内或者附近的个体，会被 SSO 模型最优地拟合，也就是说，所得的 sp 的分布是一个允许范围内的截断分布。当 sp 估计值的整体分布的众数小于允许范围的中间值时，会造成 sp 的估计值在中间值以下的可能性略高于在中间值以上的可能性。同时，随着试次数的增多，我们设定的 SSO 模型允许的 sp 的范围的下限不变，但上限会变大，因而会导致估计值的均值上升（因为接近范围的中间值），也就是我们在上一轮修改稿中指出的，随着试次数的增加， sp 估计值的均值会有所上升。总的来看，出现评审指出的这种情况，不太可能是 statistical artefact 的结果，而更有可能是因为当试次数足够多时， sp 估计值的分布满足单峰形态，且众数位置的可能性，并不显著高于其他位置的可能性。当然，以上结论只是基于目前有限的数据和特定分析得出的推测。要进一步检验上述解释，还需要更多的实证研究。需要指出的是，对于 sp 的估计并非本文关注的重点，而且考虑到多次转换的可能性，目前对 sp 的估计也可能存在偏差。此外，在论文中添加上述分析内容，可能会影响阅读的流畅性，并且分散对于核心问题即策略转换可能性的关注。因此，我们并没有对论文进行相应的修改，仅在相关文本位置添加了一个脚注，希望评审专家能够理解。

.....

审稿人 2 意见：

感谢作者们对我之前建议的考量和回复，我没有其他意见了。另外，我赞同作者们对于第一位评审专家的 AIC vs. BIC 问题的回复。因为所有模型都可能是错的，BIC 应用的前提条件（模型集中有一个模型是真模型）通常不能被满足，因此，采用 AIC 进行模型比较是更为合理的选择。

回应：

感谢审稿专家对我们修改稿和之前回复的肯定。

第三轮

审稿人 1 意见:

感谢作者对 sp 值问题的细致分析和详尽描述。虽然产生这个问题的真正原因还不很确定（可能也无法确定），我同意作者最后的观点：“对于 sp 的估计并非本文关注的重点，而且考虑到多次转换的可能性，目前对 sp 的估计也可能存在偏差。”因此，我觉得现在讨论的深度是足够的。建议文章发表。

编委意见:

两位审稿人对于模型验证指标的选取讨论不仅对于作者很有参考意义，对于编委来说也同样有所获益。认同审稿人的评价，建议接受该稿并发表。

主编意见:

同意外审和编委意见，建议录用。