



HAL
open science

Last shall be first: A field study of biases in sequential performance evaluation on the Idol series

Lionel Page, Katie Page

► **To cite this version:**

Lionel Page, Katie Page. Last shall be first: A field study of biases in sequential performance evaluation on the Idol series. *Journal of Economic Behavior and Organization*, 2009, 73 (2), pp.186. <10.1016/j.jebo.2009.08.012>. <hal-00728417>

HAL Id: hal-00728417

<https://hal.science/hal-00728417v1>

Submitted on 6 Sep 2012

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

Accepted Manuscript

Title: Last shall be first: A field study of biases in sequential performance evaluation on the Idol series

Authors: Lionel Page, Katie Page

PII: S0167-2681(09)00211-X
DOI: doi:10.1016/j.jebo.2009.08.012
Reference: JEBO 2442

To appear in: *Journal of Economic Behavior & Organization*

Received date: 30-6-2008
Revised date: 18-8-2009
Accepted date: 18-8-2009

Please cite this article as: Page, L., Page, K., Last shall be first: A field study of biases in sequential performance evaluation on the Idol series, *Journal of Economic Behavior and Organization* (2008), doi:10.1016/j.jebo.2009.08.012

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



Last shall be first:
A field study of biases in sequential performance
evaluation on the Idol series.

Lionel Page*
Westminster Business School
University of Westminster
35, Marylebone Road
NW1 5LS, London, UK
l.page@wmin.ac.uk

Katie Page
Heythrop College
University of London
Kensington Square
W8 5HQ, London, UK
k.page@heythrop.ac.uk

*Corresponding author. Tel: +44 2 07911 5000 (2706)

Last shall be first:
A field study of biases in sequential performance
evaluation on the Idol series.

Abstract

When performances are evaluated they are very often presented in a sequential order. Previous research suggests that the sequential presentation of alternatives may induce systematic biases in the way performances are evaluated. Such a phenomenon has been scarcely studied in economics. Using a large data set of performance evaluation in the Idol series (N=1522), this paper presents new evidence about the systematic biases in sequential evaluation of performances and the psychological phenomena at the origin of these biases.

JEL codes: D81, Z1

Keywords: order effects, memory, television show

We frequently make judgements and decisions about information which is presented to us in a sequential manner. This in particular is the case when we have to quickly assess the performance of individuals within a pool of contestants: job interviews, singing auditions, political debates, or even dating evenings.

The psychological literature suggests that sequential presentation of information may influence the way each piece of information is processed and recorded. Studies in economics (Neilson, 1998) and marketing (Novemsky and Dhar, 2005) have also found that a choice among situations of sequential choices may be dependent on the history of the sequence. This issue is of special importance for situations of performance evaluation. If there is any effect of the order in which people are assessed on the final evaluation of individual performances, it means that the evaluation process is biased. Stated simply, what should be completely irrelevant information (the passing order) plays a significant role in the evaluation process.

The issue of potential bias in performance evaluation raises two main concerns: efficiency and fairness. First, from the perspective of the assessor, any bias in the evaluation process results in a loss in terms of efficiency because the best options may eventually not be selected. Second, from the perspective of the contestant, any bias in the evaluation process raises the question of the fairness of the selection process: are some contestants disadvantaged relative to others for irrelevant reasons?

If there are biases in evaluation processes involving a sequential ordering of contestants/options, we need to be aware of these in order to design strategies to minimize their adverse effects and ensure that outcomes are as fair and efficient as possible.

Paradoxically, few studies have attempted to assess empirically the presence of systematic biases in the sequential evaluation of performance (Bruine de Bruin, 2005). In economics, the fairness and efficiency of performance evaluation procedures have mostly been studied relative to the possible biases arising from the judges' incentives (Prendergast and Topel, 1993; Clerides and Stengos, 2006) and from discriminating preferences (Goldin and Rouse, 2000; Segrest Purkiss et al., 2006). The economic literature has largely ignored the possible distortions arising from the pure cognitive biases in the evaluation of performance. Such biases, if significant and of practical importance, must however be studied carefully in order to limit their detrimental effects on the efficiency and fairness of selection procedures which rely on the evaluation of performances.

Using a unique dataset on the Idol series spanning competitions from eight countries (Australia, Brazil, Canada, Germany, India, Netherlands, United Kingdom, USA), this paper contributes to our understanding of order biases in performance evaluation in a naturalistic setting. Because of their generic format, the Idol shows provide a large set of identical situations where a group of individuals have to perform sequentially and are assessed by television viewers who vote for them.

The statistical analysis of this large dataset of 1,522 performances over 165 shows confirms some of the previous empirical literature on ordering effects and contributes to furthering our understanding of the underlying psychological phenomena of these effects. Our results suggest that systematic biases in sequential evaluation of performance arise through two parallel processes: the effect of the ordering on the propensity to remember each candidate, and the propensity to assess a contestant by comparing him or her to the previous contestant(s).

The remainder of the paper is organized as follows, Section 1 presents the literature on sequential biases in performance evaluation, Section 2 presents our dataset and Section 3 our results. Section 4 concludes.

1 Sequential biases in performance evaluation

There are two main reasons why biases may result from sequential ordering. The first is that judges may not remember equally well the different performances in the sequence, and second, the criteria/benchmark of the evaluation may change over time. For example, the evaluation of a performance may be dependent on the history of previous performance(s).

These potential caveats may produce two types of biases. First, ordering biases may result because your performance evaluation is conditional on your passing order. The second potential bias is that the evaluation of one's performance may directly depend on the quality of the previous performance(s). We will call these two types of biases respectively "sequential order bias" and

“sequential history bias”.

1.1 Sequential order bias

Few studies have addressed the effect of order on judgments of performance. Generally the research evidence indicates that later serial positions benefit from more positive evaluations. The evidence comes from several naturalistic studies on performance in competitions, including a study on international synchronized swimming competitions (Wilson, 1977), work on the Queen Elizabeth Contest for violin and piano (Flôres and Ginsburgh, 1996; Glejser and Heyndels, 2001), and studies of the Eurovision song contest (Bruine de Bruin, 2005) and ice skating competitions (Bruine de Bruin, 2005, 2006).

Wilson (1977) showed that there was a significant negative correlation between serial positions and final ranks in the 1973 World Championship synchronized swimming championships and an amateur meet held in the same year such that better rankings tended to be in later serial positions. An evaluation of the judgments by 15 experts in the Queen Elizabeth Contest for classical violin and piano by Glejser and Heyndels (2001) showed that musicians who performed on a later day in the competition received better judgments. Moreover, higher overall rankings were also given for performances scheduled later in the week and later in the evening (Glejser and Heyndels, 2001). Bruine de Bruin (2005) examined the effect of order in both the Eurovision song contest and ice skating judgments. She found an increasing linear trend such that contestants who were in the later serial positions had significantly higher ratings than those in the earlier positions. This effect was also found in her follow up study on ice skating with a larger data set (Bruine de Bruin, 2006).

Two potential explanations exist in the literature for this observed order bias. First, there is a well established literature on the effects of order on memory. The serial position effect is the phenomenon demonstrating that recall accuracy (usually for words) varies as a function of an item’s position within a list (Murdock, 1962). Specifically, there are two main effects: the primacy effect and the recency effect. When asked to free recall items from a list participants generally remember better those stimuli at both the beginning (primacy effect) and end (recency effect) of a sequence, resulting in a roughly u-shaped curve. The serial position effect is a robust well researched phenomenon in the cognitive psychological literature (Glanzer and Cunitz, 1966; Burgess and Hitch, 1999; Gershberg and Shimamura, 1994).

These serial position effects have been demonstrated both in the laboratory (Singh and Cole, 1993; Snyder and Harrison, 1997) and in naturalistic settings (Terry, 2005; Pieters and Bijmolt, 1997). Different memory mechanisms have been proposed to underlie the primacy and recency effects, with primacy effects linked to long term memory and recency effects explained through short term memory mechanisms (Glanzer and Cunitz, 1966). Moreover, several factors have been found to influence or alter their effects, for example distinctiveness (Neath and Crowder, 1996), emotional content (Rubin and Friendly, 1986; Maratos et al., 2000; Snyder and Harrison, 1997), prolonged distraction (Glenberg et al.,

1980) and the length of the series (Anderson et al., 1998). Generally though, holding other factors constant, first and last items are remembered better.

The interest in the role of memory and its limitations in economic decision making has grown recently (Dow, 1991; Piccione and Rubinstein, 1997; Benabou and Tirole, 2002; Mullainathan, 2002; Devetag and Warglien, 2003; Bernheim and Thomsen, 2005; Devetag and Warglien, 2007). For instance Mullainathan (2002) proposes an exponential decay of recall probabilities, compatible with a recency effect:

$$E(f_k) = f \frac{1 - \rho^{-k}}{1 - \rho}$$

where f_k is the probability to forget the event and $\rho \in \{0, 1\}$ a parameter representing the propensity to remember an event from one period to the next. Recently, such a recency effect has been integrated by Sarafidis (2007) in a model where individuals can anticipate such biases and use them strategically. More generally a recency effect will be compatible with $\partial E(f_k)/\partial k \geq 0$ for $k \geq k'$, and a primacy effect with $\partial E(f_k)/\partial k \leq 0$ for $k < k''$.

These memory explanations have been seldom linked to the evaluation of sequential performance, and this issue has not yet been studied in economics. If we extend the idea suggested by Mullainathan (2002) to model memory limitations as the effect of time on the probability to remember an event, it is clear that contestants whose performance/qualities are less likely to be remembered are less likely to be positively selected. The primacy and recency effect would therefore suggest that contestants who are in earlier and later positions will benefit positively as a result of their performances being remembered better.

A second possible explanation for the empirical results is proposed by Bruine de Bruin (2005). They explain their results through a direction of comparison effect. Specifically, they posit that as each new option is presented judges search for unique features (positive or negative) in the performance and, if found, these influence upwardly (for positive unique features) and downwardly (for unique negative features) the judgments, because more weight is given to these unique aspects rather than any overlapping (similar) features of the performance. Overall, they conclude that the direction-of-comparison effect is most prominent in tasks that promote sequential judgment, and in options with unique positive features (Bruine de Bruin and Keren, 2003). They further speculate that the direction-of-comparison effect may have contributed to the linear order effects found in jury evaluations of world-level figure skating contests (Bruine de Bruin and Keren, 2003), international synchronized swimming competitions (Wilson, 1977), the Eurovision Song Contest for popular music (Bruine de Bruin, 2005), and the Queen Elizabeth Contest (Glejser and Heyndels, 2001). However, this would only be the case if the judges were focused on the unique positive features of each performance, which may or may not have been the case.

Fundamentally, the idea of a direction-of-comparison effect relies on a specific form of reference dependent preferences which is one of the most important hypotheses in modern behavioural economics (Bruni and Sugden, 2007). The

direction of comparison effect supposes that performances are evaluated for their differences relative to a previous set of performances. If contestants present positive and negative differences relative to previous ones, and if the judges focus on the positive ones, then a systematic positive trend in the evaluation of contestants should appear. It is, however, not clear why judges would focus on positive differences, while not taking into account negative differences.

1.2 Sequential history bias

The second possible bias in the sequential evaluation of performance is that each person's performance evaluation may depend on the performance of the previous person relative to whom they are often implicitly compared. For each judgment in a given sequence (with the exception of the first judgment), it is the case that the judge has already very recently evaluated another target on that same dimension. Therefore, the knowledge the judge has activated to make that previous judgment is highly accessible at the time the next judgment has to be made. Consequently, this knowledge of the previous judgment is likely to influence the subsequent judgment (Damisch et al., 2006). Thus, the evaluation of a target at almost any point of the sequence is likely to be affected by the information that was activated during the preceding judgment of another target on that dimension (Damisch et al., 2006, p. 167).

Mussweiler's 2003 selective accessibility model outlines two main comparison processes - contrast and assimilation - that take place during the assessment of two consecutive stimuli. These comparison processes are, from an economic point of view, two different forms of reference dependent evaluation. Contrast occurs when judges focus on differences in the stimuli, and assimilation occurs when the focus is on similarities. More precisely, the direction of the influence is determined by the perceived similarity between the two sequential stimuli. A priori it is not clear what phenomenon is likely to be at work in a sequential performance evaluation, but regardless of its nature it is likely to create biases in the individual evaluation of performances because the evaluation of a contestant's performance will be depend on the performances of the previous contestant.

Damisch et al. (2006) examined sequential performance judgments in both the 2004 Olympic Games and data gathered in a laboratory setting. Their aim was to apply the concepts in the selective accessibility model (Mussweiler, 2003; Mussweiler et al., 2004) to sequential judgments in sport. Their results demonstrated that the score of an athlete increases with increasing scores of his or her immediate predecessor and decreases with decreasing scores of his or her predecessor, showing assimilation rather than contrast. Moreover, this effect carries on after the first person such that the correlation between a target and subsequent targets, whom are not immediately after the target (but second third etc), are also significant. According to research by Mussweiler et al. (2004) and Gentner and Markman (1994) unless otherwise instructed judges tend to search for similarities in the performances of people, that is, assimilation often appears to be the default judgmental outcome, resulting in significant positive

correlations between performances.

This paper investigates these two potential form of biases in the sequential evaluation of performance in a large data set from a naturalistic setting. Its unique contribution is twofold. First, no previous work has evaluated these two biases concurrently, therefore this paper adds to the existing work by enabling a direct comparison of these two processes in sequential order effects on performance evaluation. This is extremely important because it will enable us to isolate what factors are contributing to an observed ordering effect in performance and provide clearer theoretical implications.

Second, this paper uses a large, multicultural dataset which has the advantage of ecological validity and generalisability. A large majority of the previous studies of these order biases tend to be laboratory based or naturalistic studies using much smaller or restricted datasets. Our paper is unique in this respect and hence provides a strong base for testing the theoretical predictions.

2 The data

Our data consist of observations of the ranking of contestants in live shows for several pop Idol series: Australia (Australian Idol: 2003, 2004, 2006, 2007; X-factor: 2005), Brazil (Ídolos Brazil: 2007), Canada (Canadian Idol: 2003, 2004, 2005, 2006, 2007), Germany (Deutschland sucht den Superstar: 2003, 2004, 2006, 2007), India (Indian Idol: 2006, 2007), Netherlands (Idols: 2005; X factor: 2006), UK (X-factor: 2004, 2005, 2006, 2007), and the USA (American Idol: 2002, 2003, 2004, 2005, 2006, 2007). All of these shows share the same format in their final stage, specifically, the final set of contestants (10 to 13 depending on the series) are progressively eliminated one by one after each show. In each show participants have to perform a new song. Their performance is then assessed by television viewers who can vote for their preferred performance by telephone (as many times as they want). The votes are tallied and one of the last two (or three) contestants who have received the fewest votes from the public is then eliminated (sometimes this last step is determined by a choice from the judges).

The generic format of these shows, which is almost identical across countries and seasons, provides a unique opportunity to study the effects of ordering on the evaluation of individual performance. In addition, the variety of countries in our sample ensures that our results are not idiosyncratic to a given culture or to a given series.

For each season we observe the final shows where candidates have to perform one song one after the other before the public is allowed to vote for them. We do not analyse the very final stage of the competition, when four or five competitors are left and they each sing two or more songs. We therefore only observe shows where there are between 5 and 13 competitors singing one song and one or two competitors being voted off at the end of each show. These data have been collected from various online sources: wikipedia.org, tv.com and the shows' websites.

Table 1: Breakdown of the number of shows by country and number of contestants

Contestant	AUS	BRA	CAN	GER	IND	NED	UK	USA	Total
5	4	0	5	3	2	0	1	1	16
6	4	1	5	4	2	1	4	5	25
7	4	1	4	4	2	2	4	6	25
8	5	1	5	3	2	2	4	6	26
9	3	1	4	4	2	1	3	4	21
10	4	1	4	4	1	2	2	6	22
11	3	0	1	2	1	1	3	4	14
12	4	1	0	0	1	0	3	4	13
13	0	0	0	2	1	1	0	0	3
Total	31	6	28	26	14	10	24	36	165

Due to the marketing policy of the show, and in order to maintain the highest suspense during the competition the shows do not reveal the exact proportion of votes for each contestant. However, we do have some information about the rankings of the contestants because the bottom two, three or four competitors are revealed for each show.

3 Method

To assess the existence of a bias in the evaluation of contestants' performances, we compare the empirical probability to be "safe" during one show to the theoretical probability to be "safe" (when there are no biases from the sequential ordering).

Imagine a series of shows with a constant number N of contestants and suppose that these contestants have the same qualities (hence the same a priori probability to be safe). Let $b_k \in \{2, 3, 4\}$ be the number of individuals in the bottom tier for a show k , the probability to be safe for a contestant is:

$$p_k = 1 - \frac{b_k}{N}$$

Suppose now that the ordering of the performances in the live show has an impact on the evaluation of the performance by the television viewers. Some contestants will be favored by their position in the series and other disadvantaged. Let us call $bias(X, Z)$ the systematic departure from the theoretical probability of being safe where X is a set of variables characterising the position of the contestant in the passing order, and Z is a set of variables describing the characteristics of previous contestants. The probability to be safe for a participant in this position is

$$p_i = 1 - \frac{b_k}{N} + bias(X, Z)$$

Suppose that, in this simple situation, we want to estimate the bias linked with every position i in the order, $E(bias(i, Z)|i)$, we could compare the theoretical probability to be safe $p_T = 1 - b_k/N$ to the actual frequency of safe contestants in each position i , $\hat{p}_i = \sum \mathbb{1}_{\{i \text{ is safe}\}}/N_s$, where N_s is the number of shows observed:

$$E(bias(i, Z)|i) = \sum \frac{\mathbb{1}_{\{i \text{ is safe}\}}}{N_s} - \left(1 - \frac{b_k}{N}\right)$$

Our data are slightly more complex than this example because the number of contestants varies across the shows. To estimate $E(bias(X, Z)|X, Z)$, we calculate the variable $bias_{jk}$, which, for a participant j performing in the show k takes the value:

$$bias_{jk} = \mathbb{1}_{\{j \text{ is safe}\}} - \left(1 - \frac{b_k}{N_k}\right)$$

By definition, we have $E(bias(X, Z)|X, Z) = E(bias_{jk}|X, Z)$. We can then define the two biases found in the literature as:

Definition 1 (Sequential order bias) *There is a sequential order bias as soon as for any variable x_j characterising the position of a performance j in the passing order:*

$$E(bias_{jk}|x_j) \neq 0$$

Definition 2 (Sequential history bias) *There is a sequential history bias as soon as for any variable z characterising the previous candidates:*

$$E(bias_{jk}|z) \neq 0$$

The following sections will consecutively study these two possible biases.

4 Sequential order bias

A sequential order bias arises when a candidate is advantaged or disadvantaged because of his/her position in the order. To study this possible bias, we first look at the value of $E(bias_{jk}|i)$ which represents, for a given position in the order of appearance i , the difference in percentage points between the actual and theoretical probability to be safe. It therefore measures the advantage/disadvantage the position confers to a contestant in terms of the probability to be safe. Specifically, if $E(bias_{jk}|i)$ is positive then a contestant j in position i is more likely to be safe, and if $E(bias_{jk}|i)$ is negative he/she is less likely to be safe.

*** Figure 1: Bias in performance evaluation by position order ***

Figure 1 presents the mean bias per position over the whole set of positions. A clear pattern emerges which shows a positive trend as the order increases. However, this graph is slightly inaccurate because the relative location of each position may be different. For example, the 5th position will be the last one in some situations, while in other situations it will be located in the middle between the beginning and the end of the series. In this graph the last position also consists of different positions, for example, sometimes it is 5th, 9th or 11th. Figures 2 and 3 present the decomposition of the ordering effect for the shows which have between 5 and 12 contestants. The last contestants appear to benefit from a positive bias, while contestants in the middle of the order (especially closer to the beginning) seem to be disadvantaged.

*** Figure 2: Order effect for each type of show ***

*** Figure 3: Order effect for each type of show ***

In order to summarize the effects at the beginning and at the end of the sequence, Figure 4 compares the evolution from the beginning of the order to the evolution when looking at the reverse order. The last contestants appear to have a significant advantage relative to the contestants in other positions.

*** Figure 4: Bias in performance evaluation at the beginning and the end of the series ***

Overall, these results suggest that there seems to be an increasing linear trend such that contestants in the later positions have an advantage relative to those contestants in earlier positions. The worst positions in terms of bias seem to be positions two and three.

One potential caveat of the research concerns the allocation process of the contestants. The above analysis assumes the random ordering of contestants to positions. What if this is not the case? In fact, there are two main reasons to think that the ordering is not random.

First, the goal of the production is to maximise the entertainment value and, if there is not a strict rule about the random allocation of contestants, this could produce a spurious correlation between the ordering and the results. For example, better quality contestants could be more likely to be placed in some specific positions (like the beginning or the end) just for production purposes. This implies that even if there were no ordering effect at all, a selection bias could induce some differences between the probability of success of different positions.

Second, the production company could have an agenda regarding the contestants, and therefore be willing to keep good contestants longer because they will attract more viewers in later shows. So, if there is any ordering effect, they could use it to advantage/disadvantage some contestants. This implies that if there is an ordering effect, the magnitude of this effect could be biased by a selection effect. In order to control for this potential caveat, we implement fixed effect models and estimate the ordering effects while controlling for the ability of the contestant.

To analyse the effect of the ordering on the evaluation of the performance of contestants, it is possible to use a linear regression model with the variable *bias* as a dependent variable. Contestants in general perform more than once in the shows, we can therefore write:

$$bias_{jk} = \beta_0 + X_{jk}\beta + u_j + \varepsilon_{jk} \quad (1)$$

where X_{jk} is a vector of variables relative to the order i of the participant j in the show k . The term u_j is an individual effect specific to the individual j representing his/her ability. If the allocation of contestants is random, contestants performing at different positions in the order do not tend to have, on average, differences in ability: $E(u_j|X_{jk}) = 0$. In such a situation, the OLS estimator is unbiased but not efficient and a random effect model must be used instead.

However, one may doubt the hypothesis of random allocation of contestants. One could suspect, for instance, the production company to select, on average, better contestants to perform at the end of the show. In this case, we have $E(u_j|X_{jk}) \neq 0$ and the random effect estimation will be biased. To control for such a possibility we use a fixed effect estimator to estimate equation (1). The fixed effect estimator is a *within* estimator which uses only the variations in results observed within each contestant when they perform in different positions¹.

Using hypothetical contestants, Figure 5 presents the intuition of this estimator and demonstrates how it corrects for a possible bias in the allocation of contestants. Part 1.1 and 1.2 of Figure 5 show that when there is no allocation bias, the fixed effect (FE) estimator is identical to the OLS or random effect (RE) estimator. If, on the contrary, there is an allocation bias such that strong contestants are allocated to better positions in the passing order (2.1 and 2.2) the FE estimator corrects appropriately for the selection bias. Part 2.1 of Figure 5 shows that if there is no ordering effect the FE estimator will accurately show that ordering does not impact on each contestant's results. Part 2.2 shows a situation where there is an ordering effect and an allocation bias. In this case the RE estimation is biased upward, while the FE estimation gives the correct estimate of the ordering effect.

*** Figure 5: Identification strategy: using within variations in results to eliminate a possible systematic bias in the allocation of contestants ***

If there is no order effect, no variable x from X_{jk} should have a significant coefficient. Given that the result for each contestant is not independent of the results of other contestants within a given show, these models are estimated with a clustered robust variance matrix with the shows as clusters.

For all shows the order variable was normalised between 0 (first) and 1 (last). A dummy variable was created to capture the difference between being the first

¹It is the equivalent of an ANCOVA in psychology and other social sciences with the contestant playing the role of the grouping variable. While psychologists use ANCOVA to study the between groups effect controlling for covariates, economists use the fixed effect model to study the effect of the covariates controlling for systematic differences between groups (here the contestants).

to perform (1) and all other positions (0). Table 2 presents the regression results. The first three columns are random effect estimations; they are more efficient and well identified if the ordering of candidates is not linked with their specific characteristics. The last three columns are fixed effect estimations, they are unbiased even if the ordering of contestants depends on their specific characteristics.

Table 2: Regression: the ordering effect on performance evaluation

	Dependent variable: bias			
	Random effects		Fixed effects	
	(1)	(2)	(3)	(4)
Order	0.202*** (6.25)	0.265*** (6.67)	0.181*** (5.07)	0.234*** (5.70)
First		0.111* (2.39)		0.092 (1.87)
Cons	-0.139*** (-8.69)	-0.182*** (-7.85)	-0.090*** (-4.58)	-0.128*** (-5.09)
R^2			0.022	0.026
N	1522	1522	1522	1522
Number of group	352	352	352	352
Hausman test p-value			0.263	0.492

t -statistics in brackets, * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Overall the order effect is very significant and implies that, with the exception of the first position, moving one position closer to the end of the show provides an additional 5 percentage point chance of being safe for a contestant. Therefore, ordering plays a major role in the competition, at least to discriminate between contestants close in ability (which is often the case in the latter rounds of such competitions).

The difference between the random effects and fixed effects model gives an indication about the existence of a selection bias of contestants for each position. The coefficients are very close indicating that the order effect is very unlikely to be driven by a selection bias. To test for a significant difference between the coefficients of the two types of models, we need to implement a generalised version of the Hausman test given that we use a matrix of variance robust to the clustering of data in our estimation of both models (Wooldridge, 2001, p. 291). In both cases this test indicates no significant difference in coefficients between the two models (p-values in the last row of Table 2). This result suggests that the random effects models are consistent and must be considered as the best estimation procedure available. Practically, this means that there is no reason to think that the results are driven by a non random allocation of the candidates.

Figure 6 presents the estimation of the parametric prediction from the fixed effect model and a non-parametric estimation using a local linear regression

for greater flexibility. The two curves match very well and this confirms the good calibration of the linear models. The results for the effect of ordering on performance evaluation show a J-shaped curve rather than a U-shaped curve indicating both primacy and recency effects, with a stronger recency effect.

*** Figure 6: Effect of the relative order on performance evaluation ***

5 Sequential history bias

Another bias possibly arising from the sequential ordering of contestants is that the evaluation of a contestant's performance may be influenced by the performance of the previous contestant to whom they may be compared. If there is an assimilation process, we would expect that contestants performing just after a good contestant are more likely to be highly evaluated and to be safe. On the contrary, if there is a contrast effect, we would expect it to be a disadvantage to perform after a good contestant as this is likely to negatively affect the evaluation of the contestant's performance.

It is possible to have an indicator of the quality of the contestant using the previous results of each contestant. We calculate the indicator *strong* which is a binary variable indicating if the candidate has always been safe in the previous shows. While some contestants are in the bottom only once when they are eliminated, there are other contestants who are in the bottom several times before being eliminated. For every show after the first one, there are two categories of contestants: those who have always been safe before and those who have been in the bottom tier in a previous show. Arguably, for a given show, a contestant who has never been in the bottom tier previously is less likely to be in the lower range of the ranking than contestants who have previously been in the bottom tier.

Using the variable *strong*, we examine the effect of being preceded by *strong* contestants on the probability to be safe. We therefore estimate the model:

$$bias_{jk} = \beta_0 + X_{jk}\beta + \sum_{h=1}^6 strong_{i-h} + u_j + \varepsilon_{jk} \quad (2)$$

Where $strong_{i-h}$ is the dummy variable indicating if the contestant who passed h positions before has always been safe in previous shows.

Table 3 displays the results of this model. The estimation of the random effect model does not indicate any effect of the quality of previous contestants. However, the fixed effects model suggests a strong effect of the previous contestant. The Hausman test indicates that the coefficients in the fixed effects model are significantly different from the coefficients in the random effects model. This suggests that the random effects model is inconsistent. This may be the case if, for instance, the producers of the shows tend to prevent placing two weak candidates consecutively. The effect estimated in the fixed effects model is then underestimated in the random effects model.

The results of the fixed effects model suggests a significant and important effect of the previous contestant's quality on the evaluation of the current contestant's performance. When the previous contestant has never once been in the bottom tier before, the current contestant has 10 percentage points more chance to be safe. The coefficients for other previous contestants are also negative but lower, and almost always non significant.

Table 3: Regression: the comparison effect relative to the previous contestant

	Dependent variable: bias							
	Random effects				Fixed effects			
	(1)	(2)						
Order	0.272*** (6.88)	0.288*** (6.06)	0.291*** (5.20)	0.310*** (4.61)	0.251*** (5.91)	0.249*** (4.96)	0.234*** (3.93)	0.239** (2.99)
$strong_{i-1}$	0.047 (1.84)	0.043 (1.51)	0.047 (1.53)	0.027 (0.82)	0.108*** (3.90)	0.102** (3.21)	0.092** (2.61)	0.056 (1.47)
$strong_{i-2}$		-0.008 (-0.30)	-0.015 (-0.49)	0.003 (0.09)		0.034 (1.08)	0.016 (0.48)	0.028 (0.70)
$strong_{i-3}$			0.026 (0.84)	0.014 (0.41)			0.069* (2.13)	0.062 (1.58)
$strong_{i-4}$				-0.033 (-0.97)				-0.012 (-0.30)
Cons	-0.225*** (-7.17)	-0.222*** (-4.92)	-0.241*** (-3.94)	-0.209** (-2.65)	-0.219*** (-6.50)	-0.239*** (-5.24)	-0.260*** (-4.04)	-0.229* (-2.56)
R^2					0.047	0.039	0.033	0.023
N	1339	1156	973	790	1339	1156	973	790
Nb of group								
Hausman p-value					0.001 <	0.001 <	0.001 <	0.001 <

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Given that these results suggest an assimilation effect, we tested if this effect was stronger for contestants of the same gender relative to contestants from different genders. We did not find any indication of a stronger assimilation effect for contestants with the same gender.

6 Test of the random allocation of contestants

In the previous sections we have been careful to control for a possible non-random allocation of contestants in the show. It is, however, interesting to check if this allocation is random or not. While this is not necessary to assess the validity of our previous analyses which are robust to a possible selection bias in the allocation of contestants, it is interesting in itself to check if the allocation is random. It allows us, in particular, to check whether our choice of the random effect model was justified.

Given that the shows do not reveal the number of votes received for each contestant, it is not possible to directly assess if the allocation of contestants is random. There are, however, some ways to assess if the allocation is roughly random, or if it tends to be systematically biased. Some relevant information comes from the fact that for a small subset of shows in the American version, a website (Dialidol.com) provides estimates of the success of each contestant in term of votes. The website estimates the number of phone calls sent for each candidate (voters have to call a number specific to the candidate they want to support). This website has proved very successful in its estimations

Table 4: Test of the random allocation of the contestants

Conditional logit	
	Strong candidate
Order	-0.0299 (0.22)
First	-0.0939 (0.22)
Observations	1153
R-squared	<0.001

Robust standard errors in parentheses
 *** p<0.01, ** p<0.05, * p<.1

with rates of success in predictions of 87, 91 and 97% in the last three seasons respectively. Using these numbers we can see if over these three seasons there is a link between the results of candidates in previous shows and their place in the ordering sequence in a show. Using the sum of the results over the last shows as an indication of quality we estimated, using local linear regression, how the average quality of a contestant varies as a function of the order in a show. Figure 7 shows the results of this estimation and indicates that there is no link between the relative place in the ordering sequence and the average quality of the contestant.

*** Figure 7: Random allocation of the contestants in the American Idol shows ***

While worth noting, this result concerns only a subset of our sample ($N = 215$). Whilst we do not have complete information on the results of contestants for our whole dataset, the information on the performances of the contestants on previous shows provides us with a way to test more generally if there is a random allocation in the show. We can test if “strong” contestants who have never been in the bottom tier in previous shows are more likely to be at the end or the beginning of the show. To do so, we assess the probability that a contestant at a given order is strong depending on his/her order:

$$strong_{ik} = \beta_0 + X_{ik}\beta + \nu_k + \varepsilon_{ik} \quad (3)$$

Where ν_k is the fixed effect specific to the show k . This fixed effect approach is necessary as the proportion of candidates having been placed in the bottom tier previously may change from one show to the next, typically it can increase with the number of shows in the competition². Assuming that the term ε_{ik} represents an error with a logit distribution, this model is a conditional logit. Table 4 presents the results of the estimation of this model.

²Note that this does not bias the estimations presented in Table 3.

These results confirm what our previous analyses suggest: there is no systematic bias in the allocation of contestants relative to the passing order. That is, better contestants are not more likely to be toward the end or the beginning of the show.

7 Discussion

Our results indicate that in a competition the order of contestants may have a decisive role in the evaluation of their performance. Given the importance of job interviews or oral examination competitions in allocating positions and rewards, these results should raise concerns about the necessary awareness of these potential biases in the evaluation process.

More specifically, our analyses suggest that two mechanisms, memory and direct comparison, both play a role in the order bias. With respect to memory it appears that both primacy and recency effects are implicated when sequentially evaluating performance. Irrespective of ability, contestants who perform first are more likely to be positively evaluated than those who come in second and third positions, which provides evidence of a primacy effect. Contestants who perform in the later serial positions (particularly last position) have the largest advantage with respect to positive evaluations, implying a strong recency effect. The curve showing performance evaluation by serial positions is J-shaped for this dataset implying a much stronger recency effect. These results are partially consistent with those of Bruine de Bruin (2005) who found an increasing linear trend. However, we find evidence of a small primacy effect while Bruine de Bruin (2005) found no benefit to being in first position. A close reading of her results indicates that she actually found a positive effect for the first contestant but this was not significant at 5% ($t=1.72$). Given that her sample size included only 47 shows while ours includes 165 shows, the non-significance of her result was quite possibly due to a limited sample size.

Our results seems to indicate that memory limitations do play a role in the sequential evaluation of performance. In addition, they also suggest that the primacy effect could receive more attention in economics. The economic models of memory limitation like those of Mullainathan (2002) and Sarafidis (2007) only integrate a recency effect.

The second bias we demonstrate is a direct comparison effect with the previous contestant. Specifically, one's performance evaluation is influenced by the evaluation of the previous contestant. If you perform after a weak contestant there is a bias such that you are more likely to be evaluated poorly than if you perform after a strong contestant. Therefore, we find evidence for an assimilation effect with respect to sequential judgements. These findings lend further support to the selective accessibility model of Mussweiler (2003) and Mussweiler et al. (2004). Specifically, our results indicate that judges tend to assess performances based on similarities with the previous contestant and not differences. This also concurs with evidence from Damisch et al. (2006). Overall, we show that these two effects both operate and are important explanatory mechanisms

in the evaluation of sequential performance.

One factor which could influence these findings concerns the changing performance as a result of being privy to the performance of others. Specifically, it could be plausible that people change their performance (increase level of effort, motivation) after having witnessed the previous performance(s). This mechanism could work in one of two ways. If the task is novel the contestants could learn from the previous performances. However, this is not the case in most tasks which have been studied in the literature (sport and singing competitions) as the task is known in advance. Second, previous performances could act as a benchmark or goal for which the future contestant can aim. Exactly how this process works is unclear and not easy to predict. It could, however, be an explanation for the apparent dominance of assimilation over contrast because the actual performance is changing rather than the criteria of the judges. One way to test this idea would be to investigate these biases in cases where performances are not seen by the contestants, for example in job interviews or private auditions and compare these effects to those cases where the performances are able to be witnessed.

A limitation of the current study is that we do not have information about the number of people who are watching the shows throughout the broadcasts. It is possible, although unlikely in our opinion, that more people are watching the show toward the end of the program and these very same people who miss the beginning of the show also decide to vote. First, it seems likely that the people who are voting are the more ardent fanatics and are less likely to miss the beginning of the show. Second, even if there was a large enough proportion of people voting who miss the early performance(s) then this would mean that we should just see an increasing monotonic trend (assuming people do not vote for people they do not see). Having found a significant primacy effect this result is contrary to this prediction. If anything, these “late voters” would bias downwards the primacy effect which means our estimate of the initial memory effect is likely to be conservative.

Relatively speaking the magnitude of the effect is quite large and therefore is likely to have a significant impact on both the contestants and the judges. Specifically, it is significant enough to raise questions about the fairness of the process from the contestants’ perspective and to pose problems in relation to the efficiency of the process from the perspective of the judges. These findings have implications for the way in which performances should be evaluated. At the very least judges (and perhaps contestants) could be made aware of these effects. What they do with this information and how best they assimilate it into their judgments (performances) remains to be studied.

This work also suggests that future research is definitely needed in this area to study in depth these effects. For example, questions that need to be addressed include which is the stronger of these two mechanisms? Do these biases depend of the type of competition and the delay before judging? Also, does making people aware of these biases eliminate them? Moreover, future work needs to study the conditions under which assimilation and contrast are likely to occur in the evaluation of sequential performance. Are certain types of performances

(those that are judged on a tight set of criteria) more likely to lead to assimilation effects?

References

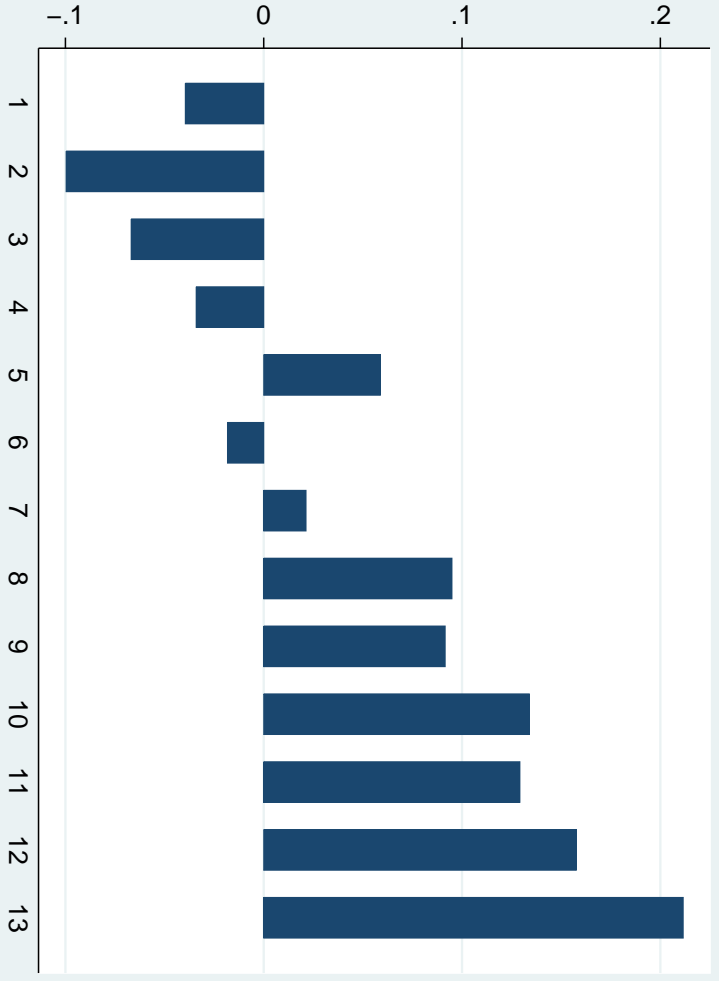
- Anderson, J., Bothell, D., Lebiere, C., Matessa, M., 1998. An Integrated Theory of List Memory. *Journal of Memory and Language* 38, 341–380.
- Benabou, R., Tirole, J., 2002. Self-Confidence and Personal Motivation. *Quarterly Journal of Economics* 117, 871–915.
- Bernheim, B., Thomsen, R., 2005. Memory and Anticipation. *The Economic Journal* 115, 271–304.
- Bruine de Bruin, W., 2005. Save the last dance for me: unwanted serial position effects in jury evaluations. *Acta Psychologica* 118, 245–260.
- Bruine de Bruin, W., 2006. Save the last dance II: Unwanted serial position effects in figure skating judgments. *Acta Psychologica* 123, 299–311.
- Bruine de Bruin, W., Keren, G., 2003. Order effects in sequentially judged options due to the direction of comparison. *Organizational Behavior and Human Decision Processes* 92, 91–101.
- Bruni, L., Sugden, R., 2007. The road not taken: how psychology was removed from economics, and how it might be brought back. *The Economic Journal* 117, 146–173.
- Burgess, N., Hitch, G., 1999. Memory for serial order: A network model of the phonological loop and its timing. *Psychological review* 106, 551–581.
- Clerides, S., Stengos, T., 2006. Love Thy Neighbour, Love Thy Kin: Strategy and Bias in the Eurovision Song Contest. Centre for Economic Policy Research.
- Damisch, L., Mussweiler, T., Plessner, H., 2006. Olympic Medals as Fruits of Comparison? Assimilation and Contrast in Sequential Performance Judgments. *Journal of Experimental Psychology Applied* 12, 166.
- Devetag, G., Warglien, M., 2003. Games and phone numbers: Do short-term memory bounds affect strategic behavior?. *Journal of Economic Psychology* 24, 189–202.
- Devetag, G., Warglien, M., 2007. Playing the wrong game: An experimental analysis of relational complexity and strategic misrepresentation. *Games and Economic Behavior* .
- Dow, J., 1991. Search Decisions with Limited Memory. *Review of Economic Studies* 58, 1–14.

- Flôres, R. G., Ginsburgh, V. A., 1996. The Queen Elisabeth Musical Competition How fair is the final ranking. *The Statistician* 45, 97-104.
- Gentner, D., Markman, A., 1994. Structural alignment in comparison: No difference without similarity. *Psychological Science* 5, 152–158.
- Gershberg, F., Shimamura, A., 1994. Serial position effects in implicit and explicit tests of memory. *Learning, Memory* 20, 1370–1378.
- Glanzer, M., Cunitz, A., 1966. Two storage mechanisms in free recall. *Journal of Verbal Learning and Verbal Behavior* 5, 1–360.
- Glejser, H., Heyndels, B., 2001. Efficiency and Inefficiency in the Ranking in Competitions: the Case of the Queen Elisabeth Music Contest. *Journal of Cultural Economics* 25, 109–129.
- Glenberg, A., Bradley, M., Stevenson, J., Kraus, T., Tkachuk, M., Gretz, A., et al., 1980. A two-process account of long-term serial position effects. *Journal of Experimental Psychology: Human Learning and Memory* 6.
- Goldin, C., Rouse, C., 2000. Orchestrating Impartiality: The Impact of “Blind” Auditions on Female Musicians. *The American Economic Review* 90, 715–741.
- Maratos, E., Allan, K., Rugg, M., 2000. Recognition memory for emotionally negative and neutral words: an ERP study. *Neuropsychologia* 38, 1452–1465.
- Mullainathan, S., 2002. A Memory-Based Model of Bounded Rationality. *Quarterly Journal of Economics* 117, 735–774.
- Murdock, B., 1962. The serial position effect of free recall. *Journal of Experimental Psychology* 64, 482–488.
- Mussweiler, T., 2003. Comparison Processes in Social Judgment: Mechanisms and Consequences. *Psychological Review* 110(3), 472–489.
- Mussweiler, T., Rüter, K., Epstude, K., 2004. The ups and downs of social comparison: Mechanisms of assimilation and contrast. *Journal of Personality and Social Psychology* 87, 832–844.
- Neath, I., Crowder, R., 1996. Distinctiveness and very short-term serial position effects. *Memory* 4, 1–18.
- Neilson, W., 1998. Reference Wealth Effects in Sequential Choice. *Journal of Risk and Uncertainty* 17, 27–48.
- Novemsky, N., Dhar, R., 2005. Goal Fulfillment and Goal Targets in Sequential Choice. *Journal of Consumer Research* 32, 396–404.
- Piccione, M., Rubinstein, A., 1997. On the Interpretation of Decision Problems with Imperfect Recall. *Games and Economic Behavior* 20, 3–24.

- Pieters, R., Bijmolt, T., 1997. Consumer Memory for Television Advertising: A Field Study of Duration, Serial Position, and Competition Effects. *Journal of Consumer Research* 23, 362.
- Prendergast, C., Topel, R., 1993. Discretion and bias in performance evaluation. *European Economic Review* 37, 355–65.
- Rubin, D. C., Friendly, M., 1986. Predicting which words get recalled: measures of free recall, availability, goodness, emotionality, and pronunciability for 925 nouns. *Memory and cognition* 14, 79–94.
- Sarafidis, Y., 2007. What Have you Done for me Lately? Release of Information and Strategic Manipulation of Memories. *The Economic Journal* 117, 307–326.
- Segrest Purkiss, S., P. Perrewé, T. Gillespie, B. Mayes, and G. Ferris, 2006. Implicit sources of bias in employment interview judgments and decisions, *Organizational Behavior and Human Decision Processes* 101, 152–167.
- Singh, S., Cole, C., 1993. The Effects of Length, Content, and Repetition on Television Commercial Effectiveness. *Journal of Marketing Research* 30, 91–104.
- Snyder, K., Harrison, D., 1997. The affective auditory verbal learning test. *Archives of Clinical Neuropsychology* 12, 477–482.
- Terry, W., 2005. Serial Position Effects in Recall of Television Commercials. *The Journal of General Psychology* 132, 151–164.
- Wilson, V., 1977. Objectivity and effect of order of appearance in judging of synchronized swimming meets. *Perceptual and Motor Skills* 44, 295–298.
- Wooldridge, J., 2001. *Econometric Analysis of Cross Section and Panel Data*: MIT Press.

Figure 1

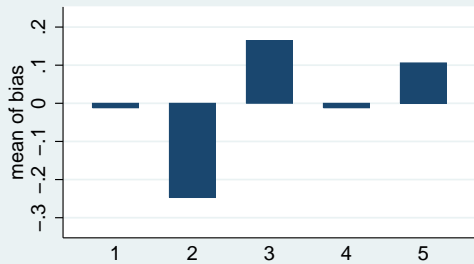
mean of bias



Preprint
Manuscript

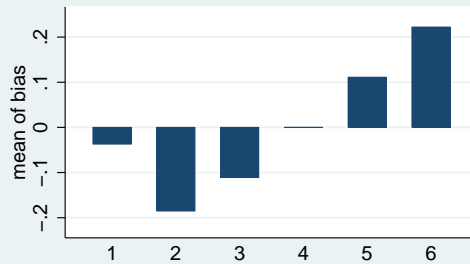
5 contestants

Nb of sessions=17



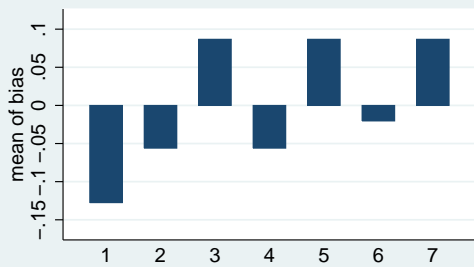
6 contestants

Nb of sessions=27



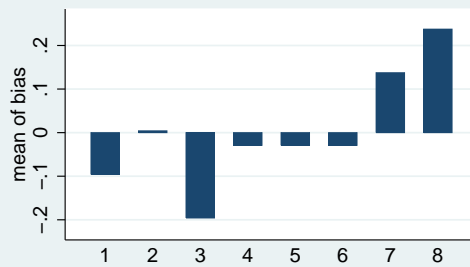
7 contestants

Nb of sessions=28



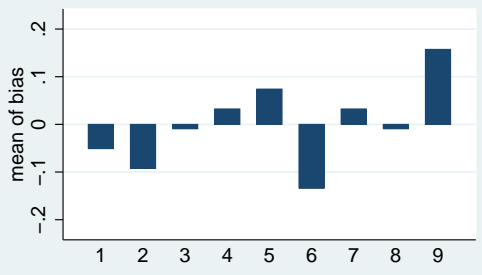
8 contestants

Nb of sessions=30

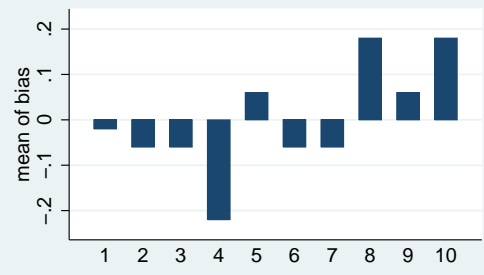


Accepted Manuscript

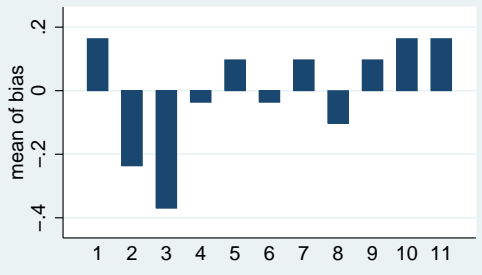
9 contestants
Nb of sessions=24



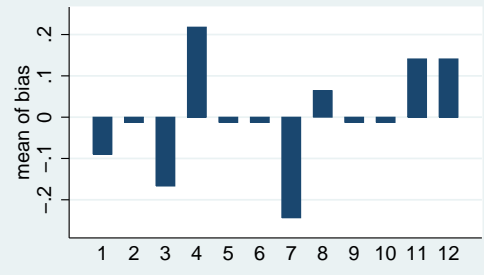
10 contestants
Nb of sessions=25



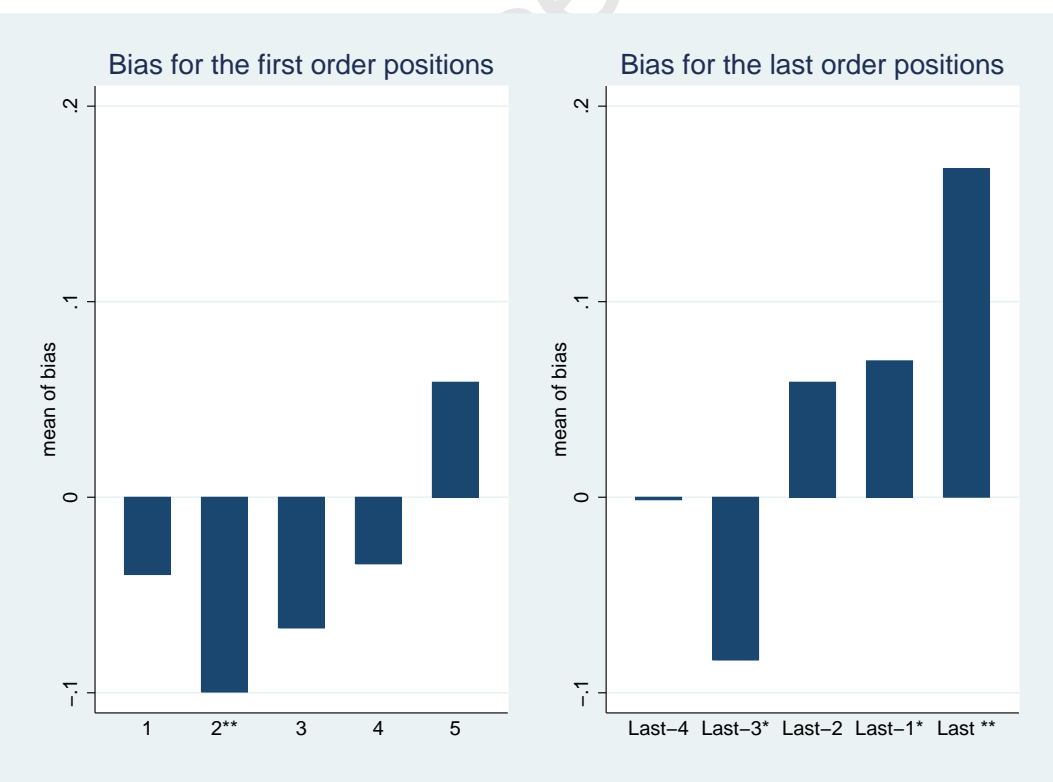
11 contestants
Nb of sessions=15



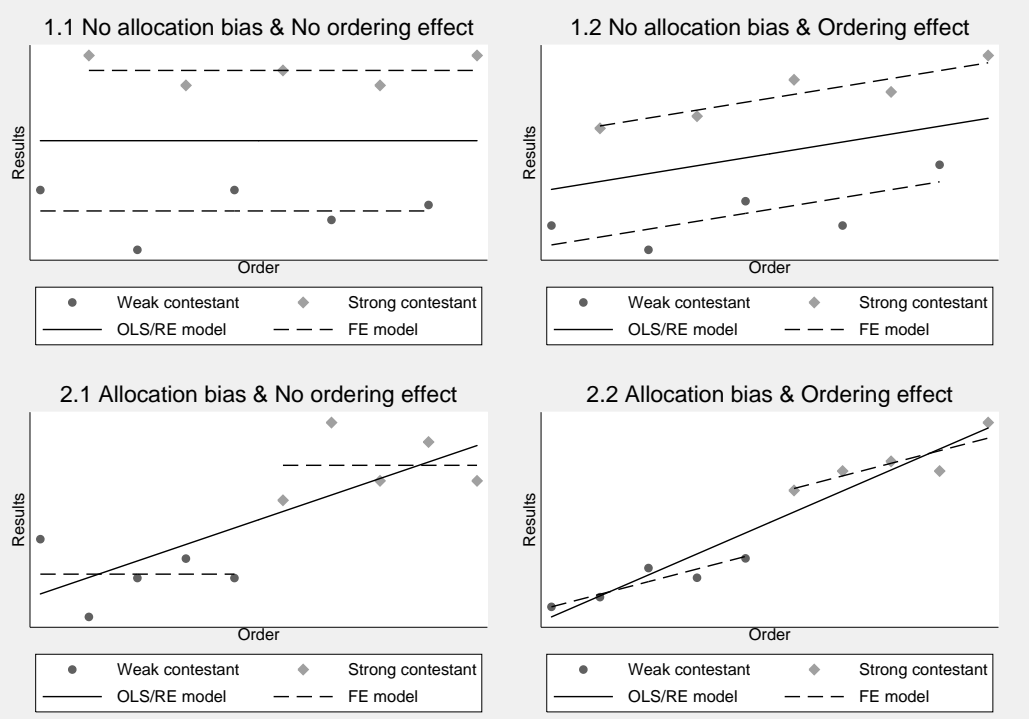
12 contestants
Nb of sessions=13



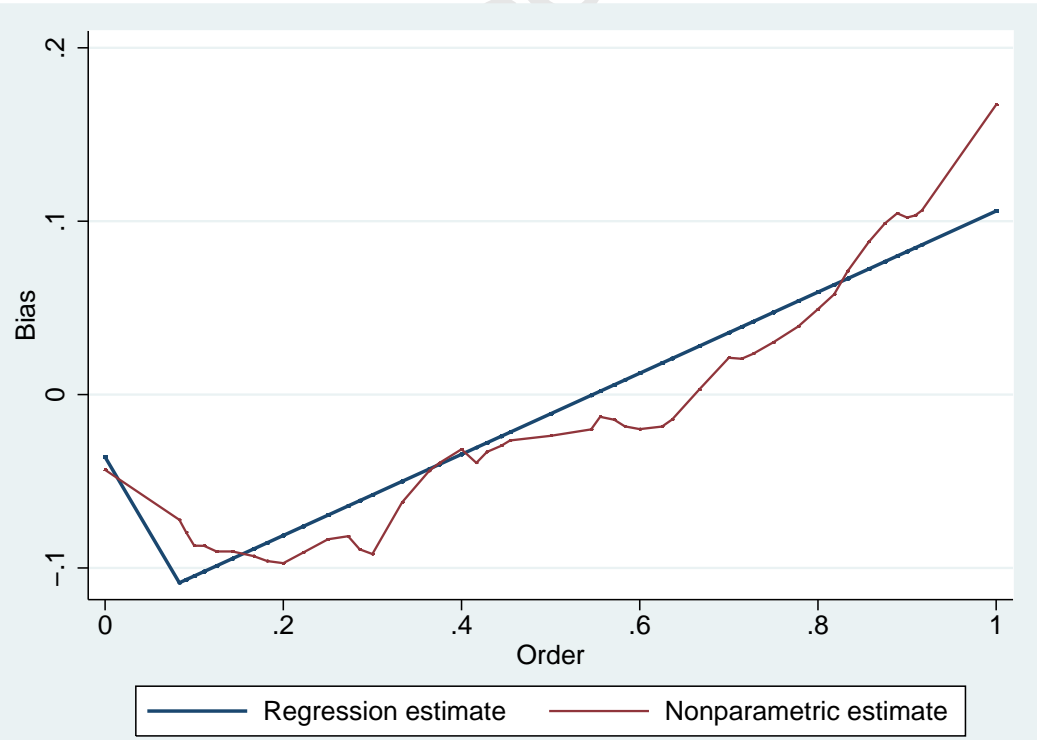
Accepted Manuscript

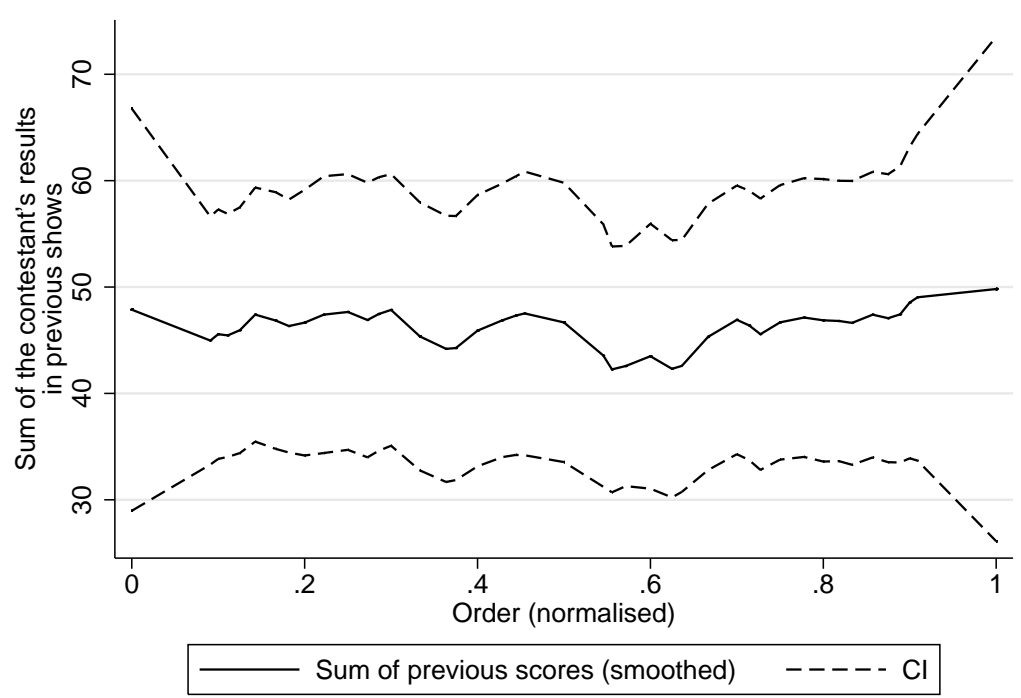


Accepted Manuscript



Accepted Manuscript





Note: the sum of the last scores is smoothed using a local linear regression