

Top performers are not the most impressive when extreme performance indicates unreliability

Jerker Denrell^{a,1} and Chengwei Liu^b

^aSaïd Business School, University of Oxford, Oxford OX1 1HP, United Kingdom; and ^bWarwick Business School, University of Warwick, Coventry CV4 7AL, United Kingdom

Edited* by James G. March, Stanford University, Stanford, CA, and approved April 26, 2012 (received for review October 14, 2011)

The relationship between performance and ability is a central concern in the social sciences: Are the most successful much more able than others, and are failures unskilled? Prior research has shown that noise and self-reinforcing dynamics make performance unpredictable and lead to a weak association between ability and performance. Here we show that the same mechanisms that generate unpredictability imply that extreme performances can be relatively uninformative about ability. As a result, the highest performers may not have the highest expected ability and should not be imitated or praised. We show that whether higher performance indicates higher ability depends on whether extreme performance could be achieved by skill or requires luck.

regression to the mean | randomness | social learning | performance evaluation | ecological rationality

Extrême performance attracts people's attention. People tend to believe the most successful are the most skillful and that failures lack skill (1, 2). A tendency to imitate the most successful has also been argued to be a basic universal trait that is shaped by evolution and promotes adaptiveness (3, 4). However, is success necessarily an indication of skill and worthy of praise and imitation and failure an indication of lack of skill?

Clearly, observed performance is not always a reliable indicator of skill. Chance events outside the control of individuals often influence performance (5–7). Moreover, such chance events rarely average out over time. Instead, due to “rich-get-richer” dynamics and “Matthew effects” (8), success usually breeds success and failure breeds failure. For example, individuals with early success might be given more resources and instruction, or consumers may favor products with a high market share (9, 10). Prior research has shown how such processes can amplify chance events and produce a weak association between performance and ability (11–13), leading to a distribution of outcomes that is both unpredictable and highly unequal (14). In such settings, extreme success and failure are, at best, only weak signals of skill. The highest performers may be more able than others and the lowest performers less able than others, but one should not expect their skill level to be very far from the mean (15).

These prior contributions show that performance and skill may be weakly associated due to noise and rich-get-richer dynamics, but they do not challenge the idea that higher performers are likely more skilled and worthy of imitation. Even if the highest performers are only marginally more skilled than others, it makes sense to imitate them. In this paper, we show that noise and rich-get-richer dynamics can have more counterintuitive implications that go beyond the conventional understanding of regression to the mean. Noise and rich-get-richer dynamics not only introduce unpredictability but also change how much one can learn from extreme performances and whether higher performance indicates higher skill. In particular, we show that when noise and rich-get-richer dynamics can strongly influence performance, extreme performances can be relatively uninformative about skill. As a result, higher performance may not indicate higher skill. The highest performers may not be the most skilled and the lowest performers may not be the least skilled. The implication is that one should not

imitate the highest performers nor dismiss the worst performers. More generally, we show that whether higher performance indicates higher skill depends on whether extreme performance could be achieved by skill or requires luck.

The intuition behind our results is that an extreme performance may be more informative about the level of noise and the strength of rich-get-richer dynamics than about skill. People often have to infer the degree of skill from performance without knowing the extent to which performance is subject to noise or the extent to which past performance influences future performance. Extreme performance indicates that the level of noise is high and that past performance strongly influences future performance, because extreme performances are more likely then. In settings with high levels of noise and when past performance strongly influences future performance, however, observed performance is a less reliable indicator of skill because chance events and early success strongly influence performance. Because extreme performances are less informative about skill levels than moderate levels of performance, a rational person should regress more to the mean when observing extreme performances, implying that the association between performance and ability can be nonmonotonic.

We develop two models to formalize this intuition. The first model assumes that current performance depends on skill but also on past performance and evaluators are uncertain about how much past performance matters. The second model assumes that performance depends on skill and noise and evaluators are uncertain about the extent to which noise matters. For both models, we show that higher performance does not indicate higher skill if luck is essential to achieving extreme performance. On the other hand, when luck is unlikely to result in extreme performance, we show that extreme performances, high or low, can be especially informative about skill.

The implication of our models contradicts descriptive studies about how people evaluate performances: People tend to believe that higher performance indicates higher skill. We show, however, that even when this assumption is faulty, people may have little opportunity or incentive to correct this assumption, as predictions based on such an assumption can be very accurate and may even outperform a correct model when information is scarce. Although it leads to accurate predictions on average, widespread use of this heuristic to identify whom to learn from can lead to diffusion of very risky behavior, and “nudges” (16) may be necessary to help people resist the temptation to praise, blame, or learn from extreme performers. In the following, we show how our conclusion follows from two simple models and discuss the descriptive and normative implications of our findings.

Author contributions: J.D. and C.L. designed research, performed research, contributed new reagents/analytic tools, analyzed data, and wrote the paper.

The authors declare no conflict of interest.

*This Direct Submission article had a prearranged editor.

¹To whom correspondence should be addressed. E-mail: jerker.denrell@sbs.ox.ac.uk.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1116048109/-DCSupplemental.

Model 1: Extreme Performance Indicates Strong Rich-Get-Richer Dynamics

One reason why extreme performance may not be a reliable indicator of skill is that an extreme performance indicates especially strong rich-get-richer dynamics. To formalize this, we develop a model in which success depends on skill but also on past success and where evaluators are uncertain about how much past performance matters.

Consider a game with 50 rounds, with each round being a success or a failure. The goal is to obtain as many “successes” as possible. Individuals differ in their skill levels: Some individuals have a higher probability of obtaining success in any given round. In addition to skill, we also assume that outcomes in consecutive rounds are dependent. The probability of succeeding increases if the previous outcome was a success. Specifically, after a success, the probability of success in the next period is $c_i(1 - w_i) + w_i$, whereas it is only $c_i(1 - w_i)$ after a failure. In the first period, the probability of success is c_i . Here, $c_i \in (0,1)$ is the skill of individual i , and $w_i \in (0,1)$ represents the extent to which success probabilities depend on the previous outcome. w_i would be high in industries where consumers want to buy the currently most popular product (due to network externalities, for example) and in careers where early success brings resources, training, and visibility that increase future success (9, 10). The level of dependency, w_i , is not the same for every player—some players are in contexts where dependency is relatively strong. We also assume that the level of dependence is not fully known to outside observers. It is often difficult to estimate the extent to which success depends on past success rather than on superior skill, especially when few data are available and when past success operates through difficult-to-observe processes such as consumer loyalty. Thus, it may not be clear to an observer whether a streak of successes is due to exceptional skills or to strong dependencies combined with the good fortune of being successful initially. Witness, for example, the debate about whether Microsoft’s success is due to their early lead or to superior quality (17).

We simulated this game with 5 million players. Each player has a different value of c_i (drawn from a beta distribution with parameters 10,10; the beta distribution is a flexible distribution and also a common choice for modeling heterogeneity in success probabilities) and w_i (drawn from a uniform distribution; i.e., a beta distribution with parameters 1,1). These assumptions about the distributions of skill and dependencies imply that the distribution of dependency (w_i) is less concentrated around 0.5 than the distribution of skills (c_i). Thus, extreme values are more likely for w_i than for c_i .

Based on the simulated data, we can examine how average skill levels (c_i) vary with the number of successes obtained. Intuitively, one might expect that players who achieved the most successes are the most impressive and have the highest value of c_i . However, as Fig. 1 shows, the association between success and skill level is nonmonotonic. The average value of skill reaches a maximum at about 40 successes out of 50 and then starts to decline. Players who achieved exceptional performance, that is, successes in 50 rounds, have an average value of skill lower than those with 40 successes. Stated differently, the most successful players are not the most impressive. Rather, moderately successful players are the most impressive ones. A similar pattern is observed for very low levels of success: The players with the lowest levels of success are not the least impressive.

The explanation for this result is that an extreme performance indicates that the level of dependency was strong (w_i high): Extreme outcomes are more likely then. When w_i is high, however, performance is less informative about skill, because outcomes are substantially influenced by chance events. Chance events can substantially influence outcomes when w_i is high, because initial outcomes then strongly influence subsequent outcomes and players

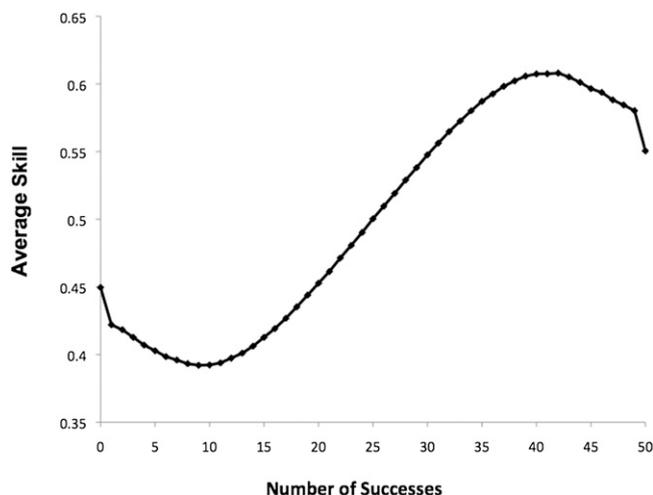


Fig. 1. The average value of c_i for players who obtained different numbers of successes in 50 rounds. Based on 5 million simulations.

with low levels of skill who get lucky initially may have many successes. Similarly, players with high levels of skill could get unlucky initially and may have many failures. As a result, the association between skill levels and eventual outcomes will be weak when w_i is high. For example, the correlation between skill and success is only 0.37 when $w_i = 0.9$ (Fig. 2A). Thus, when dependency is strong, achieving extreme performance is a less informative indicator of skill. Nevertheless, strong dependency implies that extreme outcomes are more likely compared with when dependency is weak.

In settings when dependency is weak, skill matters more and chance events less in determining the outcome. For example, the correlation between skill and success is 0.82 when $w_i = 0.1$ (Fig. 2B). Because skills matter more in determining outcomes, outcomes are also more informative about skill levels. In particular, obtaining 50 consecutive successes (or failures), when dependency is weak, is a reliable indicator of high (or low) skill, because it is very unlikely that a player without very high (or low) skill would obtain such an extreme result. However, as Fig. 2B illustrates, exceptional performance is less likely when dependency is weak compared with when it is strong.

Overall, our basic result emerges, because an extreme performance indicates that it was achieved in a context in which chance events can substantially influence outcomes, and performance is then an unreliable indicator of skill. For example, achieving 50 successes out of 50 rounds indicates that the degree of dependency must have been very high, and in such settings achieving exceptional performance is not so impressive.

This nonmonotonic association between success and skill emerges only when w_i is less concentrated around 0.5 than c_i is, namely when extreme values are more likely for w_i than for c_i . The intuition is that when very high values of dependency are more likely than very high values of skill, extreme results are likely due to a high value of w_i rather than a high value of c_i . When w_i is more concentrated around 0.5 than c_i (for example, when w_i is a known constant), expected skill is increasing in the number of successes.

Model 2: Extreme Performance Indicates Extreme Noise

Another reason why extreme performance may not be a reliable indicator of skill is that an extreme performance indicates that performance is subject to an especially high level of noise. To illustrate this, we now consider a static model in which performance depends on skill and “luck” and evaluators are uncertain about the impact of luck. Consistent with standard models in

data on observed performance levels (see *SI Text* for how the experiments were conducted). Participants in these experiments had ample time to learn the nonmonotonic relationship between ability and performance. The results show that despite clear feedback and incentives to be accurate, 69 out of 119 participants never predicted higher performers to be less skilled than those with moderately high performance. In other words, a majority of participants assumed the most successful were the most skilled and thus mistook luck for skill.

Why do people misinterpret extreme performance, and what are the consequences? We suggest that people misinterpret extreme performance partly because they rely on the assumption that higher performers are more skilled. Relying on such an assumption is not irrational. Rather, it is often true. Perhaps more interesting, even when this assumption is false, a model built on this assumption may outperform a correct model that allows for nonmonotonicity. We use a simulation model to illustrate this point (see *SI Text* for more information). In the simulation, we examined whether a third-degree polynomial model or a linear model predicted skill levels more accurately based on a sample of performance–skill pairs. The simulation assumed that the association between performance and skill was nonmonotonic (as in Fig. 3). In every period of the simulation, the performance of an individual with unknown skill was observed and the task was to predict her skill. After the prediction, the skill was revealed. Although a third-degree polynomial model can better fit the performance–skill association, the linear model made more accurate predictions for small sample sizes. For example, suppose only four observations were available, that is, four performance–skill pairs. In this case, the linear model made predictions closer to the actual value of skill than the third-degree polynomial model 76% of the time (based on 1 million simulations). Only if 20 or more observations were available was the third-degree polynomial model more likely than the linear model to make accurate predictions.

The intuition for this result is related to the bias–variance dilemma (22): Fitting a third-degree polynomial introduces variance that degrades predictive accuracy. Because people often have to evaluate performances relying on small samples (23), using a linear model as a frugal heuristic may be ecologically rational (24, 25). Even in settings when the extreme performance is not a reliable indicator of skill, decision makers may be served well by initially assuming that performance indicates skill. Moreover, unless large samples are available, people will also have little opportunity to detect that their assumption is incorrect.

Although the assumption that higher performers are more skilled leads to accurate predictions on average, assuming that extreme performances indicate extreme skill can lead to undesirable consequences if people rely on this assumption for identifying outliers especially worthy of praise and blame.

Consider blame for failures. People often attribute catastrophes to the leader in charge, but a catastrophe might be more informative about the character of the system that experienced failure than about the leader's ability. Complex systems in which components are tightly coupled are sensitive to chance events and external shocks. As a result, extreme failures are more likely for such systems than for systems in which components are loosely coupled. A catastrophe indicates that the underlying system is complex and failure-prone, and in such cases firing the

leader might be misguided. Moderate failures could provide more reliable evidence of low ability.

Consider, next, learning from successes. Imitating the practices of successful others has been argued to be an universal trait and beneficial for society (3). Our model also implies that imitating others with high, but not exceptional, performance is likely to be beneficial, whereas imitating exceptional performance could be detrimental. As our models show, the highest performers may both be less skilled and use methods with higher levels of risk. When exceptional performance is due to self-reinforcing processes and initial success, the exceptional performers may continue to perform well but imitators will likely be disappointed (26), because they can at best only replicate the practices, and thus the skill levels, of the high performers, but not their initial good fortune.

More generally, our results suggest a different perspective on when imitating top performers is beneficial for society. On the one hand, when an extreme performance is unlikely unless an agent is lucky, the highest performers might not be the most skilled. Nevertheless, imitation may be largely beneficial if people recognize that the highest performer is not the best or when they are only aware of a few others so the best observed performance is not exceptional. In a society where exceptional performers are highly visible, however, and their practices are covered in business magazines, imitation of the best may not increase skill levels as much but will lead to the diffusion of more risky practices that generate a variable outcome at the societal level (14). On the other hand, when extreme performance requires extreme levels of skill, imitation of and knowledge of the best performers can be especially beneficial, because extreme performances are especially informative.

Finally, consider rewards. The highest performers often receive the highest rewards in organizations. Our results suggest that one should suspect that extremely high performance could be due to excessive risk taking rather than prudent strategy and exceptional skill. Imitation of such highly rewarded performers may further diffuse such risky practices. Moreover, high rewards for exceptional performance may tempt other people to deliberately take risks or to cheat because they are unlikely to achieve extreme performance otherwise [as happened in Barings Bank (27)]. This observation may be relevant for the recurrent financial crises: Rewards for exceptional performance might have led to diffusion of risky practices that eventually resulted in very poor returns. To avoid this, reward systems would need to be redesigned to reflect not just actual performance but also the level of risk (due to leverage but also to focus on products or markets with strong self-reinforcing dynamics). More important, because a nonmonotonic relationship between performance and skill is counterintuitive and a reward system that reflects this relationship may not be perceived as fair, nudges may need to be developed to help people resist the temptation to praise or blame extreme performers.

ACKNOWLEDGMENTS. We thank Pamela Sammons, Stefan Scholtes, Daniel Ralph, Daniel Levinthal, Anne Miner, Daniel Read, and Thomas Powell for discussions; Lance Bai and Cheeven Tsai for technical assistance; three anonymous reviewers and the editor for their constructive comments; and Jesus College, Saïd Business School at the University of Oxford, and The Saïd Foundation for financial support.

- Baron J, Hershey JC (1988) Outcome bias in decision evaluation. *J Pers Soc Psychol* 54(4):569–579.
- Gilbert DT, Malone PS (1995) The correspondence bias. *Psychol Bull* 117(1):21–38.
- Richerson PJ, Boyd R (2005) *Not by Genes Alone: How Culture Transformed Human Evolution* (Univ of Chicago Press, Chicago).
- Rogers AR (1988) Does biology constrain culture? *Am Anthropol* 90(4):819–831.
- Thorngate W, Dawes R, Foddy M (2008) *Judging Merit* (Psychology Press, New York, NY).
- March JC, March JG (1977) Almost random careers: The Wisconsin school superintendent, 1940–1972. *Adm Sci Q* 22(3):377–409.

- Musch J, Grondin S (2001) Unequal competition as an impediment to personal development: A review of the relative age effect in sport. *Dev Rev* 21(2):147–167.
- Merton RK (1968) The Matthew effect in science: The reward and communication systems of science are considered. *Science* 159(3810):56–63.
- Frank R, Cook P (1995) *The Winner-Take-All Society* (Free Press, New York).
- DiPrete TA, Eirich GM (2006) Cumulative advantage as a mechanism for inequality: A review of theoretical and empirical developments. *Annu Rev Sociol* 32:271–297.
- Barabasi AL, Albert R (1999) Emergence of scaling in random networks. *Science* 286:509–512.

