# Local, General and Universal Prediction Strategies: A Game-Theoretical Approach to the Problem of Induction

Gerhard Schurz, University of Duesseldorf (Germany)

*Abstract:* In this paper I present a game-theoretical approach to the problem of induction. I investigate the comparative success of prediction methods by mathematical analysis and computer programming. Hume's problem lies in the fact that although the success of *object-inductive* prediction strategies is quite robust, they cannot be universally optimal. My proposal towards a solution of the problem of induction is *meta-induction*. I show that there exist meta-inductive prediction strategies whose success is universally optimal, modulo short-run losses which are upper-bounded. I then turn to the implications of my approach for the evolution of cognition. In the final section I suggest a *revision* of the paradigm of *bounded rationality* by introducing the distinction between local, general and universal prediction strategies.

## *1. Introduction: the Best-Alternative Approach to Induction*

In an *inductive inference*, a property, regularity, or frequency is transferred from the observed to the unobserved, or from the past to the future. How can we *rationally justify* inductive inferences? This is the famous *problem of induction*, or Hume's problem. David Hume has shown that all standard methods of justification fail when applied to the task of justifying induction. (1.) Obviously, inductive inferences cannot be justified by deductive logic, since it is logically possible that the future is completely different from the past. (2.) Induction cannot be justified by induction from observation, by arguing that induction has been successful in the past, whence − by induction − it will be successful in the future. For this argument is *circular*, and circular arguments are without any justificatory value (Salmon (1957, 46, has shown that also anti-induction may be pseudo-justified in such a circular manner). (3.) It is equally impossible to demonstrate that the conclusion of an inductive inferences is at least highly probable − for in order to show this, one must presuppose that the relative event frequencies observed so far can be transferred to the unobserved future, which is nothing but an inductive inference. These were the reasons which led Hume

to the skeptical conclusion that induction cannot be rationally justified, but is merely the result of psychological habit.

There have been several attempts to solve or dissolve Hume's problem, which cannot be discussed here. It seems that so far, none of these attempts has been successful in giving a *positive* solution to the problem of induction, which establishes in a *non-circular* manner that the inductive method is a superior prediction method in terms of its success frequencies. Given that it is impossible to demonstrate that induction *must* be successful (Hume's lesson), and that there are various *alternative* prediction methods, then it seems to follow that the only approach to Hume's problem for which one can at least uphold the *hope* that it *could succeed* if it were adequately developed is Reichenbach's *best alternative* approach (Reichenbach 1949, §91; Salmon 1974). This approach does not try to show that induction *must* be successful, but it attempts to establish that induction is an *optimal* prediction method – its success will be maximal among *all* competing methods in arbitrary possible worlds. Or in simplified words: if any method of prediction will work, then the inductive method will work (Rescher 1980, 207ff). It must be emphasized that in demonstrating optimality one must allow *all* possible worlds, in particular all kinds of *para-normal* worlds in which perfectly successful future-tellers do indeed exist. Restricting the set of worlds to 'normal' or uniform worlds would destroy the enterprise of justifying induction. For then we would have to justify inductively that our real world is one of these 'normal' worlds, and we would end up in exactly that kind of circle or infinite regress in which according to the *Humean skeptic* all attempts of justifying induction must end up.

Reichenbach did not succeed in establishing the best alternative argument with respect to the goal of *single event predictions*. He only demonstrated this argument with respect to the goal of *inferring* an event's frequency limit in the long run. With respect to that goal, his argument is *almost trivial*: if the sequence of events has a frequency limit, then inductive methods will long-run approximate this limit while other non-inductive methods may or may not approximate the limit; and if the sequence

does not have a frequency limit, then *no* method can find the limit (Reichenbach 1949, 474f). However, our ability to infer approximately correct frequency limits is practically not significant. What *is* of practical significance is our success in true *predictions*. In this respect, Reichenbach's approach fails: nothing in Reichenbach's argument excludes the possibility that a perfect future-teller may have perfect success in predicting random tossings of a coin, while the inductivist can only have a predictive success of 0.5 in this case (cf. also Reichenbach 1949, 476f; Skyrms 1975, ch. III.4).

By *object-induction* (abbreviated as *OI*) I understand methods of induction applied at the level of events – the 'object level'. Generally speaking, the problem of Reichenbach's account lies in the fact that it is impossible to demonstrate that object-induction is an approximately (optimal) prediction method – which is also a lesson of formal learning theory (see §2). In contrast to Reichenbach's approach, my approach is based in the idea of *meta-induction*. The meta-inductivist (abbreviated as *MI*) applies the inductive method at the level of competing prediction methods. More precisely, the meta-inductivist bases her predictions on the predictions and the observed success rates of the other (non-MI) players and tries to derive therefrom an 'optimal' prediction. The simplest type of meta-inductivist predicts what the presently best prediction method predicts, but one can construct much more refined kinds of meta-inductivistic prediction strategies.

One should expect that for meta-induction the chances of demonstrating optimality are much better than for object-induction. Is it possible to design a version of meta-induction which can be proved to be an optimal prediction method? The significance of this question for the problem of induction is this: if the answer is positive, then at least meta-induction would have a rational and non-circular justification based on mathematical-analytic argument. But this *analytic* justification of *meta-induction* would at the same time yield an *a posteriori* justification of *object-induction* in our real word: for we know by experience that in our real world, non-inductive prediction strategies have not been successful so far, whence it would be meta-inductively justi-

fied to favor object-inductivistic strategies. In other words: the common-sense argument in favor of object-induction which is based on its past success record would no longer be circular, given that we had a non-circular justification of meta-induction.

Note that the optimality of a prediction method alone is compatible with the existence of other equally optimal methods. Nevertheless, the optimality of meta-induction would already be sufficient for its rational justification, because as Reichenbach (1949, 475f) has pointed out, meta-induction is the *only* prediction strategy for which optimality can be *rationally demonstrated*. Of course, it would be desirable to extend an optimality result for meta-induction (if we had it) to a (weak) *dominance* result, i.e. to a result showing that meta-induction is the *only* optimal prediction method. But it is hard find a non-trivial version of dominance for meta-induction (cf. Schurz 2008, §6), and therefore I will concentrate on the question of optimality.

Let me finally emphasize that my notion of optimality is restricted to *accessible* prediction methods. There might be possible worlds in which alternative players do not *give away* their predictions but keep them secret. Indeed, this is possible, and so I have to restrict my claim to *accessible methods*. What I intend to show is that among all prediction methods (or strategies) who's *output* is accessible to a given person, the meta-inductivistic strategy is always the best choice. But I argue that this restriction is not a drawback. For methods whose output is not accessible to a person are not among her *possible actions* and, hence, are without relevance for the optimality argument.

*2. Prediction Games*

My prediction games shall cover *binary* as well as *real-valued* prediction games. A prediction game consists of:

(1.) An infinite sequence (e) := $(e_1, e_2,…)$ of events $e_n \in [0,1]$ which are coded or measured by elements of the unit interval [0,1]; hence $(\forall n \geq 1:) e_n \in [0,1]$. For example, (e) may be a sequence of daily weather conditions, stock values, or coin tossings.

(2.) A set of players $\Pi = \{P_1, P_2, \ldots, xMI (xMI_1, xMI_2,\ldots)\}$, whose task is to predict future events of the event sequence. $p_n(P)$ denotes the prediction of *player* P *for* time n, which is delivered *at* time n−1. Also the admissible predictions $p_n$ are assumed to be elements of [0,1]. The deviation of the prediction $p_n$ from the event $e_n$ is measured by a normalized loss function $l(p_n,e_n) \in [0,1]$. The *natural* loss-function is defined as the absolute difference between prediction and event, $l(p_n,e_n) := |p_n-e_n|$. However, my theorems will not depend on natural loss functions but hold for arbitrary (and in case of theorems 3+4 for convex) normalized loss-functions.

In *binary* prediction games, predictions as well as events must take one of the two values 0 and 1 which code instantiations of a binary event-type E ('1' for 'E obtains' and '0' for 'E does not obtain'). *Further* notation: The *score* $s(p_n,e_n)$ obtained by prediction $p_n$ given event $e_n$ is defined as 1 minus loss, $s(p_n,e_n) := 1 - l(p_n,e_n)$; the *absolute* success $a_n(P)$ achieved by player P at time n is defined as P's sum of scores until time n, $a_n(P) := \Sigma_{1 \leq i \leq n} s(p_n(P),e_n)$, and the *success rate* $suc_n(P)$ of player P at time n is defined as $suc_n(P) := a_n(P)/n$. $\bar{e}_n := (\Sigma_{1 \leq i \leq n} e_n)/n$ denotes the event's *mean* value at time n, and $\bar{e} := \lim_{n \to \infty} \bar{e}_n$ denotes the event's *limit* mean value, provided the mean values converge to a limit. For binary prediction games, (i) $suc_n(P)$ coincides with the relative frequency of P's correct predictions until time n, (ii) $\bar{e}_n$ with the relative frequency $f_n(E)$ of event E at time n, and (iii) $\bar{e}$ with E's limiting frequency.

The players in $\Pi$ include:

(2.1) One or several object-inductivists, abbreviated as OI $(OI_1,\ldots,OI_r)$. They have informational access to past events; their first prediction (at n=1) is a guess.

(2.2) A subset of alternative players $P_{r+1}, P_{r+2},\ldots$; for example, persons who rely on their instinct, God-guided future-tellers, etc. In para-normal worlds, these alternative players may have any success and any information you want, including information about future events and about the meta-inductivist's favorites. − Players of type (2.1) or (2.2) are called *non-MI-players*.

(2.3) One or several meta-inductivists, whose denotation has the form 'xMI', where 'x' is a variable (possibly empty) expression specifying the *type* of the meta-inductivist. The meta-inductivist has access to the past events and the past and present predictions of the non-MI-players.

The simplest type of meta-inductivist from which I start my inquiry is abbreviated as MI. At each time, MI predicts what the non-MI-player with the presently highest predictive success rate predicts. If P is this presently best player, then I say that P is MI's present *favorite*, or simply that MI *favors* P. If there are several best players, MI chooses the first best player in an assumed ordering. MI changes his favorite player only if another player becomes *strictly* better; otherwise he sticks to his present favorite. $fav_n(MI)$ denotes MI's favorite *for* time n, that is, the player with the first-best success-rate *at* time n−1 among the non-MI-players; observe that MI's favorite for time n is determined at time n−1. MI's first favorite is OI. I assume that MI has always access to OI: even if no person different from MI plays OI, MI may constantly simulate OI's predictions and use them their success rate is in favorite-position.

MI belongs to the class of so-called *one-favorite* meta-inductivists, which choose at each time a non-MI-player as their favorite for the next time and predict what their favorite predicts. In contrast, *multiple-favorite* meta-inductivists base their predictions on the predictions of *several* 'attractive' non-MI-players.

The simplest object-inductive prediction method, abbreviated as OI, is based on the already mentioned *straight rule*. In the case of *real-valued* events, this inductive rule transfers the observed mean value to the next event, i.e. $p_{n+1}(OI) = \bar{e}_n$. In the case of *binary* events, the straight rule is merely used for conjecturing the frequency limit (cf. Salmon 1974, 89-95; Rescher 1980, ch. VI.3). For the purpose of single event predictions one needs in addition the so-called *maximum rule*, which requires to predict an event with maximal conjectured frequency (cf. Reichenbach 1938, 310f). For binary events, this generates the prediction rule $p_n(OI) = 1$ if $f_n(E) \geq 1$, and else $= 0$, which can be summarized by saying that OI predicts the *integer-rounding* $[\bar{e}_n]$ of $\bar{e}_n$.

OI's prediction rule is appropriate as long as the event sequence is a *random* sequence; in this case OI's success rate converges in the binary case against the maximum of P(E) and P(¬E), and in the real-valued case against the limiting mean value of the absolute deviations, $\lim_{n\to\infty} (\Sigma_{1\leq i\leq n} |e_i-\bar{e}| \, / \, n)$. For non-random sequences refined object-inductivistic prediction strategies exist, whose success dominates OI's success (see §xx).

I identify prediction games with *possible worlds*. Apart from the definition of a prediction game, I make no assumptions about these possible worlds. The stream of events (e) can be any sequence you like. Should (e) be non-random, then more refined object-inductivistic strategies may exist (as explained), but nothing which concerns the behaviour of xMI hangs on that question. I also do not assume a fixed list of players − the list of players may vary from world to world, except that it always contains xMI and the (virtual) OI. I make the realistic assumption that xMI has *finite* computational means, whence I restrict my investigation to prediction games with *finitely* many players.

According to my knowledge, the use of prediction games for epistemological purposes is new in the philosophical literature. There are, however, three related approaches in related fields:

1) In *formal learning theory* (cf. Kelly 1996) only *one* player, an object-inductive scientist, plays against a stream of events, and it is investigated which cognitive tasks can reliably be achieved under which conditions on the stream of events. For inductive prediction tasks the general result is negative, because of the possibility of 'demonic' streams of events which at every time n produce the opposite of the object-inductivist's prediction. This insight goes back to Putnam (1963), and it is a variant of Hume's lesson.[1] In contrast, my prediction games consist of several prediction me-

---

[1] Cf. cf. Friend et al. 2007, ix). In formal learning theory one considers especially hypotheses evaluation tasks which are not considered here. Kelly's major result about prediction tasks is this: an infinite stream of events (e) is correctly predictable by a scientific method after some finite time iff (e) is among a recursively enumerable set of possible data streams (1996, 260ff).

thods playing against each other, and my investigation does not focus on the question of the *reliability* but of the *optimality* of methods. Even if for every meta-inductive prediction method there exist suitably chosen 'demonic' streams of events for which its predictive success is zero, such a method may still be optimal, provided one can prove that in all 'demonic' cases also al other accessible methods must have zero success.

2) A second field which comes very close to our approach, although it has not been related to the problem of induction, is the *non-probabilistic variant* of the theory of *universal prediction* (cf. Merhav/Feder 1998), which has been developed in the fields of decision and learning theory (for an overview Bianchi and Lugosi 2006). In this approach one considers *online predictions based on expert advice*: a forecaster (who corresponds to our meta-inductivist) predicts an arbitrary event sequence based on the predictions of a set of experts (who correspond to our 'non-MI-players'). One speaks of 'universal' prediction theory because the event-sequence may be arbitrary; and the setting is called 'online learning' because the players have simultaneously to learn from past events to make new predictions. In §xx we make use of a central theorem achieved in this field.

3) A third related field is the comparative investigation of the efficiency of prediction methods by Gigerenzer and the *ABC-research group* (for 'Adaptive Behavior and Cognition') by real experiments and computer simulations. Although this approach focuses on object-inductive prediction methods, some of my results bear tight relations to this approach.

## 3. Simple Meta-Induction,  Take-the-Best, and Its Limitations

For one-favorite meta-inductivists, binary prediction can be obtained as a *subclass* of real-valued prediction games, by restricting to event sequences containing only binary events, and to non-MI-players predicting only binary values. This implies that also the meta-inductivist delivers a binary prediction, because she predicts what her

favorite predicts. Therefore, our theorems about one-favorite meta-inductivists apply to real-valued as well as to binary prediction games.

I have investigate the prediction game with help *mathematical analysis* as well as *computer simulations.* The performance of a type of meta-inductivist has always two sides: (i) its *long-run* behavior, which is of central significance, and (ii) its *short-run* performance, which is also important: although one should be willing to buy *some* short-run losses of a prediction method for sake of its long-term optimality, these short-run losses should not be too large, and they should be under rational control.

In this section I investigate the performance of MI and its relative, Gigerenzer's prediction rule *TTB* (for *Take-the-Best*). From now on, $\text{maxsuc}_n$ denotes the maximal success rate of the non-MI-players at time n. The set of non-MI-players is said to contain a (unique) *best* player $B \in \{P_1,\ldots,P_m\}$ iff there exists a time point $n_B$ such that for all later times B's success rate is greater than the success rate of all other non-MI-players. $n_B$ is called B's *winning* time. The central result about MI is theorem (1.1), which tells us that MI predicts long-run optimal whenever there exists a best non-MI-player.

*Theorem 1:* For each prediction game $((e), \{P_1,\ldots,P_m, MI\}$ whose player-set contains a best player B, the following holds:

*(1.1) (Long-run:)* MI's success rate approximates the maximal success of the non-MI-players (from below): $\lim_{n\to\infty}(\text{maxsuc}_n - \text{suc}_n(MI)) = 0$.

*(1.2) (Short-run:)* $(\forall n \geq 1:)$ $\text{suc}_n(MI) \geq \text{maxsuc}_n - (n_B/n)$, where $n_B$ is B's winning time.

The proof of theorem 1 is obvious and just explained informally: after the time point $n_B$ MI's favorite will be B forever, but until time $n_B$ MI's success may be zero in the worst case, due to switching favorites (see below). This yields theorem (1.2), and (1.1) follows.

In determining her favorites the meta-inductivist must buy some losses, compared

to the best non-MI-method. These losses result from the fact that in order to predict for time n+1, the meta-inductivist can only take into account the non-MI-players' success rates until time n. Whenever MI recognizes that her present favorite $P_1$ has lost one point compared to some new best player $P_2$, then MI has also lost this one point compared to $P_2$, before MI decides to switch to $P_2$. So for each switch of favorites MI looses a score of 1 in the binary prediction game, and a non-zero score $\leq 1$ in the real-valued game, compared to the best non-MI-player. These losses may accumulate. The assumption of theorem 1 excludes that MI can have more than finitely many losses due to switching favorites; so these losses must vanish in the limit. Theorem 1.2 informs us about the maximal short-run loss of MI. Since the time point $n_B$ may come arbitrary late, MI's cumulative short run loss may be arbitrarily high. Nevertheless, the result of theorem (1.2) is at least something, because it shows that a high short-run loss of MI is caused by a late arrival of B's winning time. In conclusion, MI's optimality is restricted to prediction games which contain a best player whose winning time doesn't occur to late.

Note that the condition of theorem 1 is rather general; for example, it does not imply that the success rates of the non-MI-players have to converge to a limit. If this is the case, then the condition of theorem 1 is satisfied if there exists one player with maximal limit success. Illustrations of the behavior of MI by computer simulations can be found in Schurz (2008).

The assumption of theorem 1 is violated whenever the success rates of two or more leading non-MI-players oscillate endlessly around each other. There exist two sorts of success-oscillations: *convergent* oscillations and *non-convergent* oscillations. Convergent oscillations are given when two (or more) leading players oscillate in their success-rate around each other with constant period and diminishing amplitude; i.e. their success-difference converges against zero. Here MI looses one success point in every half oscillation period. In the *worst* case, two alternative players A and B oscillate around each other with the smallest possible period of 4 time units. In this worst case, MI gets systematically *deceived* by the alternative players, because the

alternative players predict incorrectly exactly when they are in the position of being MI's favorite. In the result, the success rates of the two alternative players converges against 1/2, while the meta-inductivist's success remains zero for all time. A computer simulation of this scenario is shown in fig. 1.[2]

Theorem 1 as well as the negative result of fig. 1 can be generalized to the prediction rule *Take-the-Best* (TTB) of the ABC-research group (Gigerenzer et al. 1999, chs. 2-4). Although TTB is usually treated as an object-inductive (rather than a meta-inductive) prediction method, this difference is just one of *interpretation*, and not of substantial content. The predictions of the non-MI-players correspond to the *cues* in Gigerenzer's setting. There are the following two differences between MI and TTB as it is used in the Gigerenzer setting:
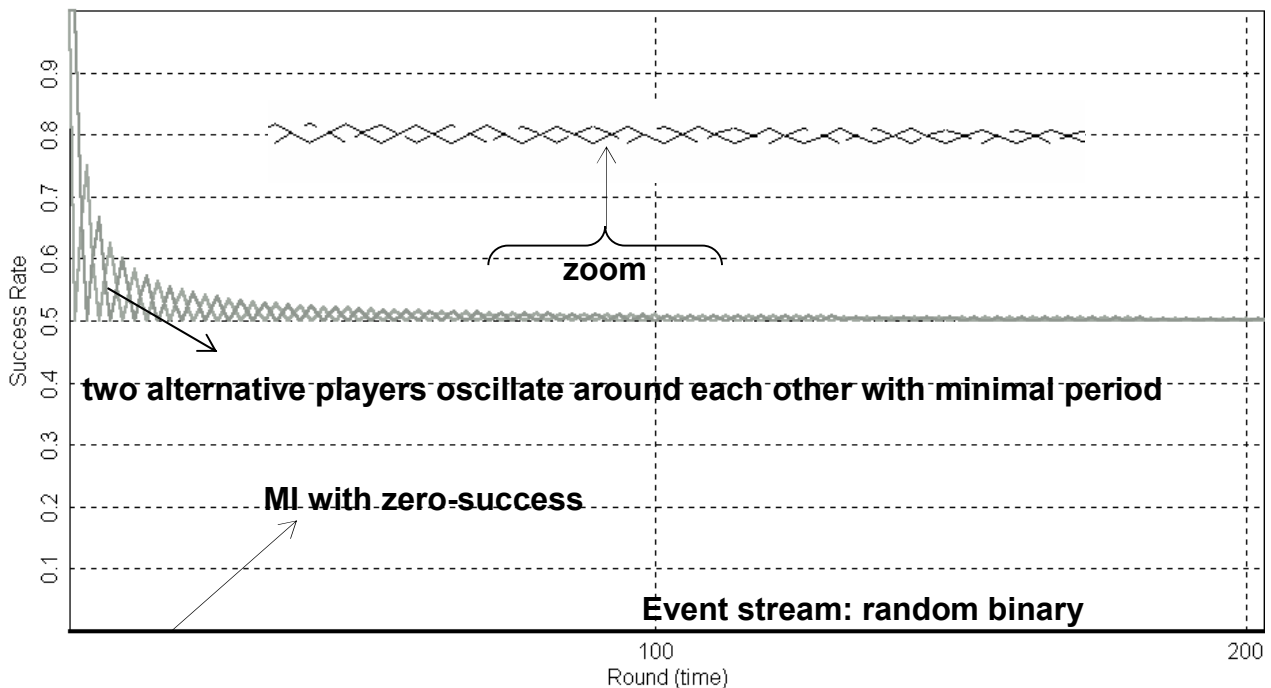


*Fig. 1: MI against two best alternative players in convergent oscillation*

(1.) The TTB strategy works like MI except that it is assumed that the cues need not make a prediction at *every* time. Thus, TTB chooses that non-MI-player as her

---

[2]    OI has been omitted in this scenario, but its addition cannot avoid MI's breakdown. If we add an OI with limit success 1/2 and assume that the limit success of the two oscillating players is slightly greater than 1/2, then MI's limit success will be slightly greater than zero.

favorite for time n who delivers a prediction for time n and has the first-best success rate among those non-MI-players who have delivered a prediction for time n.

(2.) Gigerenzer assumes that all frequencies converge to limiting frequencies, i.e. probabilities, and moreover, that the success probabilities of the cues (the so-called 'ecological validities', see p. 130), are estimated by repeated *random samplings* from finite domains. These estimations are 'inductively safe' modulo random errors. This downplays the very problem of induction. In the situation of online learning one inductively infers from past events to future events. This is not random sampling, because you cannot sample from future events, but only from past events. (In terms of sampling theory, you sample from a finite subset of an possibly infinite domain.) If the future is different from the past, inductive inference leads to systematic failure. If the rule TTB is applied to the situation of online learning in our prediction games, then theorem 1 can be generalized as follows: if there exists a time point n* after which the (strict) success ordering of the non-MI-players (or cues) remains constant, then TTB's success rate converges to a weighted average of the success rates of the non-MI-players (or cues) conditional to times when they delivered a prediction, weighted by their frequencies of delivering a prediction. A detailed outline of this generalization is left to another paper.

The ABC-group is especially interested in comparing the success of TTB with the success of *refined object-inductive* strategies such as linear regression or Bayes rule. In this paper I consider such refined rules only in the margin because they do not affect my results on meta-induction. Gigerenzer has argued repeatedly that in spite of its simplicity, TTB is almost always as good as these refined prediction strategies. Hogarth and Karelaia (2006) have shown that in scenarios with highly compensatory cues, TTB is inferior to refined object-inductive strategies, but they point out that these scenarios are rare. My results, however, reveal another restriction of TTB in scenarios of online learning: TTB will only perform well if the success rates of the cues converge sufficiently fast either towards a limit or at least to a unique success-ordering among the cues. This is assumption is implicitly granted by the explained

random sampling methods of the ABC-group. However, in scenarios of online learning with oscillating event frequencies and success rates, as for example in *predictions of the stock market*, 'inductive safety' cannot be assumed. In such a case it would be a bad recommendation to put all of one's money always on the presently most successful stock, instead of leaving it on one stock for some time (which corresponds to ε-meta-induction in §4), or distributing it over several stocks in form of a stock portfolio (which corresponds to weighted average meta-induction in §5). Fig. 2 illustrates breakdown of TTB when playing against non-MI-players (cues) with convergently oscillating success rates of the described worst-case sort.
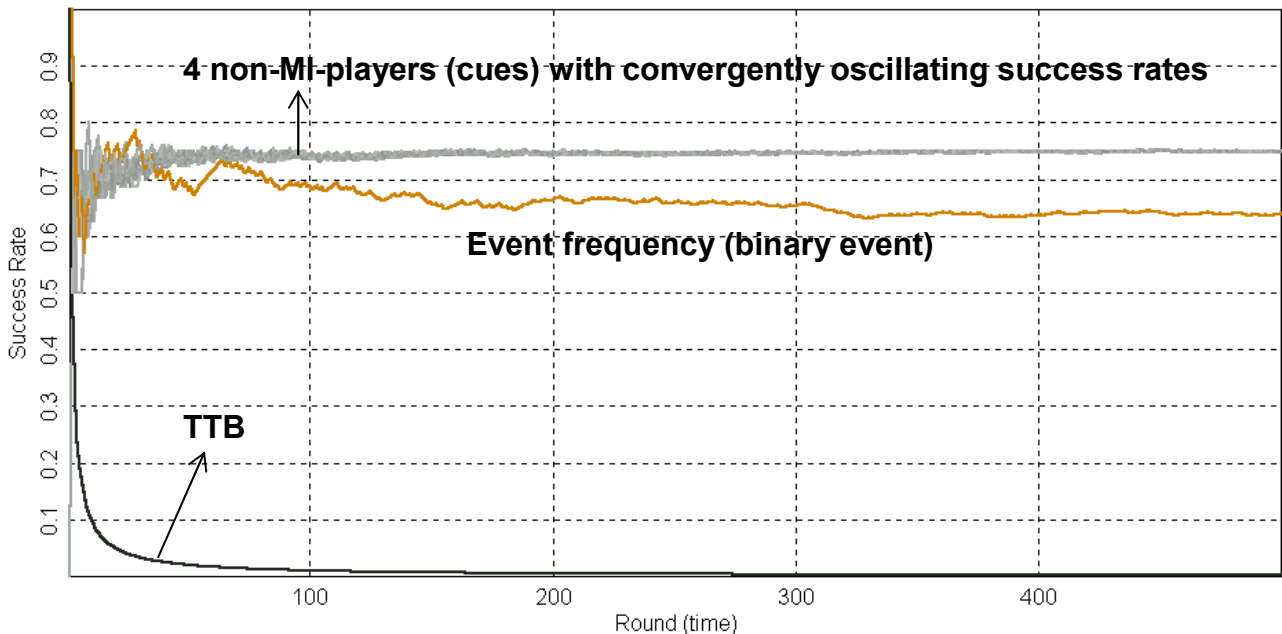


*Fig. 2: Breakdown of TTB with convergently success-oscillating non-MI-players*

## 4. Improvements of One-Favorite Meta-Inductivists

The meta-inductivist has a simple and robust defense strategy to permanent losses causes by convergent oscillations: don't switch favorites if their success difference is *practically insignificant*. I call this new type of meta-inductivist the *ε-meta-inductivist* εMI: εMI switches his favorite only if the success difference between his present favorite and a new better favorite exceeds a small threshold ε which is consid-

ered as practically insignificant. The performance of εMI is illustrated by the computer simulation in fig. 3, in which εMI plays against the two alternative players of the convergent oscillation scenario of fig. 1: when the success difference between the alternative players has become smaller  than ε, εMI stops to switch but sticks to one player, with the result that εMI's success rate recovers and ε-approximates the maximum success of the two alternative players.
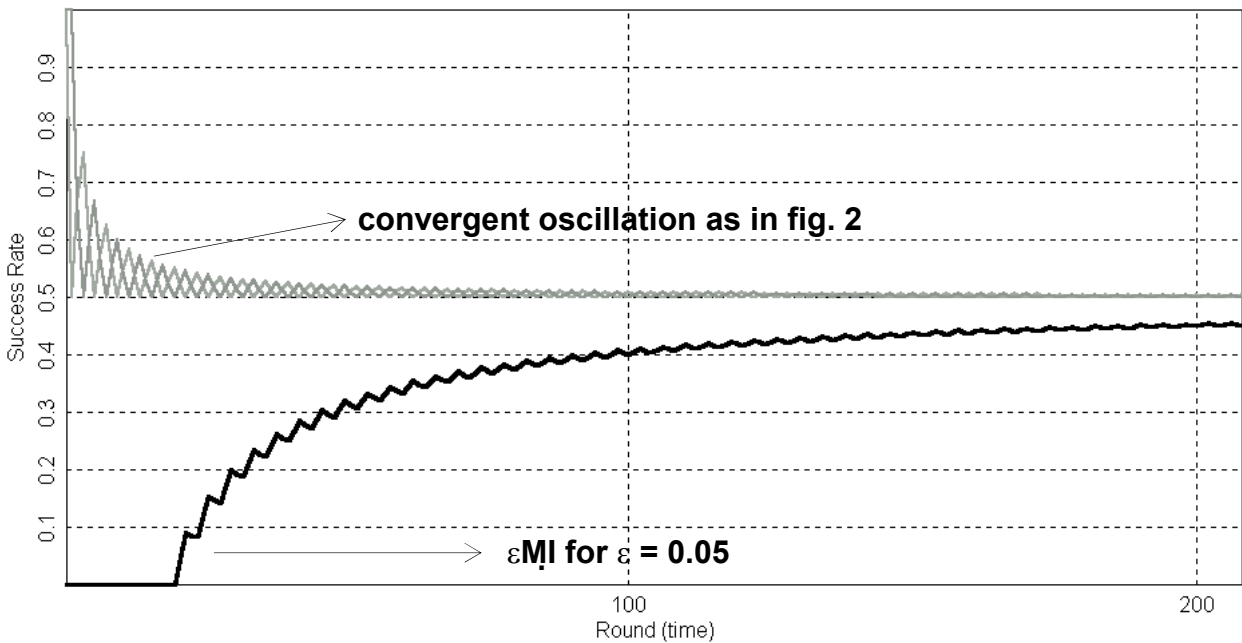


*Fig. 3: εMI  in the convergent oscillation scenario of fig.2*

The move from MI to εMI gives rise to stronger theorem than (1.1), namely theorem (2.1). We say that a prediction game contains a subset BP $\subseteq$ {P$_1$,…,P$_m$} of *ε-best* non-MI-players iff there exists a time n$_{BP}$, the *winning* time of BP, such that for all times later than n$_{BP}$, (a) each player in B is more successful than each non-MI-player outside BP, and (b) the successes of all BP-players are ε-close to each other. Theorem (2.1) establishes that εMI ε-approximates the maximal success rate if there exists a subset of ε-best non-MI-players. Again note the generality of this condition: it does not imply, but it is implied by the convergence of the non-MI-players' success rates towards a limit.

*Theorem 2:* For every prediction game ((e), {$P_1$,…,$P_m$, εMI}) whose non-MI-players set contains a subset BP of ε-best players, the following holds:

*(2.1) (Long run:)* εMI's success *ε-approximates* the maximal success of the non-MI-players (from below): $\lim_{n\to\infty}(\text{maxsuc}_n - \text{suc}_n(\varepsilon\text{MI})) \leq \varepsilon$.

*(2.2) (Short run:)* ($\forall n \geq 1$:) $\text{suc}_n(\text{MI}) \geq \text{maxsuc}_n - (n_{BP}+1/n) - 2\cdot\varepsilon$, where $n_{BP}$ is BP's winning time.

A proof of theorem 2 is found in Schurz (2008, §4). The worst-case bound of εMI's short-run loss provided by theorem (2.2) (namely $\frac{n_{BP}+1}{n} - 2\cdot\varepsilon$) is not especially good. At least, theorem (2.2) tells us that εMI's short run loss decreases with an early arrival of the winning time $n_{BP}$ of the ε-best players.

The ε-meta-inductivist is long-run optimal in a broader class of possible worlds than MI, on the cost that its optimality is not strict but *approximate*. I think that approximate optimality is still good enough to count as a justification, because the loss of an approximately optimal strategy, compared to the best strategy, is *always small* – smaller than ε. Moreover, for almost all practical purposes there exists a choice of ε which is small enough to count as *practically insignificant*, and theorem (2.1) holds for *all* choices of ε. However, there exists a *trade-off* in respect to the short-run performance, since a small ε goes usually hand in hand with a large $n_{BP}$ in theorem (2.2). So, the freedom to make ε very small is limited by the interest in keeping the short-run loss small.

The assumption of a subset of ε-best players is violated in prediction games with *non-convergent* success-oscillations. Here we find the *worst cases* for one-favorite meta-induction. If the success rates of two or more leading alternative players oscillate around each other in a *nonconvergent* manner with a nondiminishing amplitude of δ > ε, then εMI will be deceived. The minimal periods of such nonconvergent oscillations must grow exponentially in time. The worst case are so-called *systematic*

*deceivers*: they are assumed to *know* (e.g., by clairvoyance) whether the meta-inductivist will choose them as favorite for the next time, and they use this information to deceive the meta-inductivist by delivering a *worst* (i.e. minimal-score) prediction whenever the meta-inductivist chooses them as their favorite, while they deliver a correct prediction whenever they are not chosen as a favorite. For natural loss functions, the worst prediction for time n is 0 if $e_n \geq 1$, and is 1 otherwise. Hence the score of the worst prediction is 0 in the binary prediction and a value between 0 and 0.5 in the real-valued prediction game; while the score of a correct prediction is always 1.

If εMI is playing against k deceivers, then at each time there will be $k-1$ deceivers which predict correctly because they are not εMI's favorite. The computer simulation in fig. 4 shows a binary prediction game in which four alternative players deceive the ε-meta-inductivist. As long as a deceiver $D_1$ is εMI's favorite, $D_1$ predicts the wrong result until his success is more than ε below some deceiver $D_2$. At this time εMI switches his favorite from $D_1$ to $D_2$, $D_1$ starts to predict correctly and $D_2$ starts to predict wrong results, until the next switch of εMI occurs, etc. In this way, εMI's success rate is turned down to zero, while the mean success of the four deceivers per oscillation is 3/4.
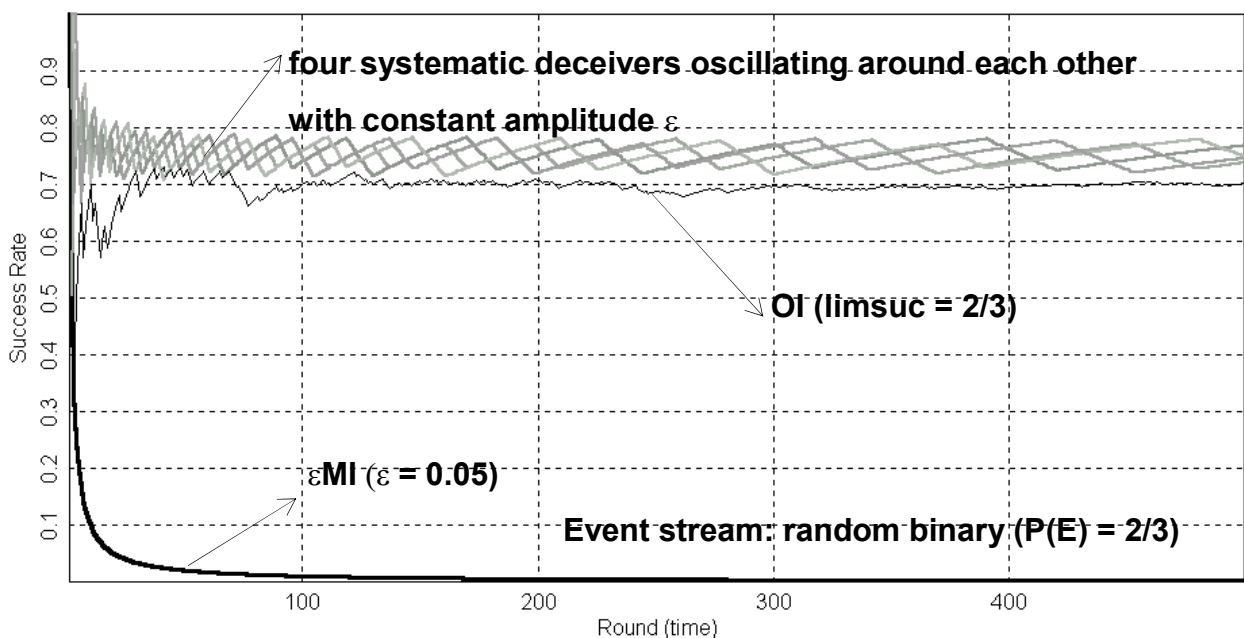


*Fig. 4: Systematic deception of εMI by non-convergent oscillations*

The negative result of fig. 4 generalizes to all kinds of *one-favorite* meta-inductivists: they must fail to be optimal whenever they play against k ≥ 2 systematic deceivers (and use them as favorites), because in that case they have zero-success, while the deceivers will have a limit success of $\frac{k-1}{k}$, because in the long run each deceiver is εMI's favorite in 1 out of k times.

In spite of this drawback of εMI I wish to emphasize that εMI is an important improvement over MI. Convergent oscillations of success rates (or more generally, relative frequencies) are nut 'unrealistic' in our real world: they occur under conditions which involve periodic developments, e.g. as in trajectories of predator-prey-systems. In contrast, non-convergent oscillations of relative (success) frequencies with exponentially growing periods are rather strange and presumably extremely rare in our real world.

i have tried to improve the performance of εMI further by assuming that εMI does never favor non-MI-players who deceive to him a certain degree, which means that their success-rate conditional on times at which their where εMI's favorites is more than a certain threshold below their unconditional success rate. I have called this kind of meta-inductivist the *avoidance* meta-inductivist aMI, and it can be shown that aMI predicts optimal with respect to *non-deceiving* non-MI-players, *even if* deceiving non-MI-players are present (cf. Schurz 2008, § 5). However, the fact remains that neither aMI nor any other one-favorite meta-inductive method can predict optimally in regard to deceivers.

Do meta-inductivists strategies exist which are indeed universally optimal? Our negative result about one-favorite meta-inductivists entails that if they exist, they must be found in the class of *multiple-favorite* meta-inductivists. In the next section we investigate their most important variant: weighted average meta-inductivists.

*5. Weighted Average Meta-Induction*

A weighted average meta-inductivist predicts a weighted average of the predictions of the non-MI-players. Since the weighted average of several predictions of zeros and ones is a real value between 0 and 1, this method cannot be applied to binary prediction games, in which all predictions must be either '0' or '1'. It can only be applied to real-valued prediction games. In this form, the method of weighted average prediction has been studied in the theory of (non-probabilistic) *universal prediction* which was mentioned in §2. The results in this literature have not at all been related to the problem of induction, but the problem setting is similar to my prediction games.

The weighted average meta-inductivist is abbreviated as *wMI* and defined as follows. For every non-MI-player P we define $at_n(P) := suc_n(P) - suc_n(wMI)$ as P's *attractiveness* (as a favorite) at time n. Let PP(n) be the set of all non-MI-players with *positive* attractiveness at time n. Then wMI's prediction for time 1 is set to 1/2, and for all times >1 with non-empty PP(n) $\neq \varnothing$ it is defined as follows:

*Definition of weighted average prediction:*

$$\forall n \geq 1: \ p_{n+1}(wMI) = \frac{\sum_{P \in PP(n)} at_n(P) \cdot p_n(P)}{\sum_{P \in PP(n)} at_n(P)}$$

*In words:* wMI's prediction for the next round is the attractiveness-weighted average of the attractive players' predictions for the next round. Should it happen that PP(n) gets empty, $p_{n+1}(wMI)$ is reset to 1/2.

Informally explained, the reason why wMI cannot be deceived is the following: a non-MI-player who tries to deceive wMI would be one who starts to predict incorrectly as soon as his attractiveness for wMI is higher than a certain threshold. The success rates of such wMI-adversaries must oscillate around each other. But wMI does not favor just one of them (who predicts incorrectly in turn), but wMI predicts according to an attractiveness-weighted average of correctly and incorrectly predict-

ing adversaries, and therefore wMI's long-run success must approximate the maximal long-run success of his adversaries.

The next theorem (theorem 3) does not hold for arbitrary but only for those loss functions $l(p_n,e_n)$ which are *convex* in $p_n$. By definition, $l(p_n,e_n)$ is convex in its argument $p_n$ iff (for fixed weights) the loss of the weighted average of two predictions is smaller-equal than the weighted average of the losses of the two predictions. It is easy to see that the natural loss-function $l(p_n,e_n) := |p_n-e_n|$ is convex. There exist many other convex loss-functions, e.g. $|p_n-e_n|^q$ for $q\geq1$, and theorem 3 applies to all of them. Theorem 3.1 establishes that wMI is indeed a *universally* long-run optimal prediction strategy, even in the strict (and not approximate) sense. Also wMI's short-run performance is good, as theorem 3.2 reveals. The number m of non-MI-players (or strategies) is under complete control and the worst-case short-run loss $\sqrt{m/n}$ quickly vanishes for times $n \gg m$.

*Theorem 3:* For every prediction game $((e), \{P_1,\ldots,P_m,wMI\})$ whose loss-function $l(p_n,e_n)$ is *convex* in the argument $p_n$, the following holds:

*(3.1) (Long-run:)* $suc_n(wMI)$ (strictly) approximates the non-MI-players' maximal success: $\lim_{n\to\infty} (maxsuc_n - suc_n(wMI)) = 0$.

*(3.2) (Short run:)* $(\forall n\geq1:)$ $suc_n(wMI) \geq maxsuc_n - \sqrt{m/n}$.

The proof of theorem 3 is found in Schurz (2008, §5); it rests on corollary 2.1 of Cesa-Bianchi and Lugosi (2006, 12f). In prediction games satisfying the conditions of theorem 1, the strategy wMI will soon converge to the simple MI-strategy, since after some time, only the best player will have positive attractiveness, whence wMI will predict as if she would favor this best player forever.

Theorem 3 does not apply to binary prediction games, because even under the assumption that the events and the non-MI-player's predictions are binary, wMI's predictions are not binary. The failure of theorem 3 for binary prediction games can be

recognized from the following example by Cesa-Bianchi and Lugosi (2006, 67): assume a meta-inductivist playing against two non-MI-players, one of them constantly predicting 1 and the other constantly predicting 0, and a 'demonic' event-sequence produces constantly the opposite of the meta-inductivist's predictions. Then whatever the meta-inductivist predicts, his success rate will be constantly zero, while the maximal success rate of the two non-MI-players must always be $\geq 0.5$. Thus, we have obtained a second general negative result: a meta-inductive strategy which predicts long-run optimal for an *individual* player in arbitrary *binary* prediction games does not exist. In combination with theorem 3 this is a deep result, insofar it shows that a continuous nature is more friendly to the inductivist than a discrete nature.

Nevertheless I have found a way to apply theorem 3 indirectly also to the prediction of binary events, namely by means of assuming a *collective* of k meta-inductivists, abbreviated as $cwMI_1,\ldots,cwMI_k$, and by considering their *mean success rate* ('cwMI$_i$' stands short for 'collective weighted-average meta-inductivist no. i'). I regard wMI's real-valued prediction as an *ideal* (though non-admissible) prediction, which is approximated by the mean value of the k binary predictions of the collective of cwMI-meta-inductivists as follows: $[p_n \cdot k]$ cwMI's predict 1, and $k - [p_n \cdot k]$ cwMI's predict 0. In this way, one obtains a universal optimality result for the *mean success rate* of collective of cwMI's, abbreviated as $\overline{suc}_n(cwMI)$, which is formulated in theorem 4 (proof in Schurz 2008, §8). The additional worst-case loss term $\frac{1}{2 \cdot k}$ reflects the maximal loss due to approximation of the ideal prediction by k binary predictions; this loss can be made arbitrarily small by increasing the number of meta-inductivists.

*Theorem 4:* For every binary prediction game $((e), \{P_1,\ldots,P_m, cwMI_1,\ldots,cwMI_k\})$:

*(4.1) (Long run:)* $\overline{suc}_n(cwMI)$ $\frac{1}{2 \cdot k}$-approximates the non-MI-players' maximal success: $\lim_{n\to\infty} (\overline{suc}_n(cwMI) - maxsuc_n) \leq \frac{1}{2 \cdot k}$.

*(4.2) (Short run:)* $(\forall n \geq 1:)$ $\overline{suc}_n(cwMI) \geq maxsuc_n - \sqrt{m/n} - \frac{1}{2 \cdot k}$.

Figure 5 shows a computer simulation of a collective of ten cwMI's playing against 4 specially designed cwMI-adversaries, who predict incorrectly as soon as their attractiveness gets higher than a variable threshold.
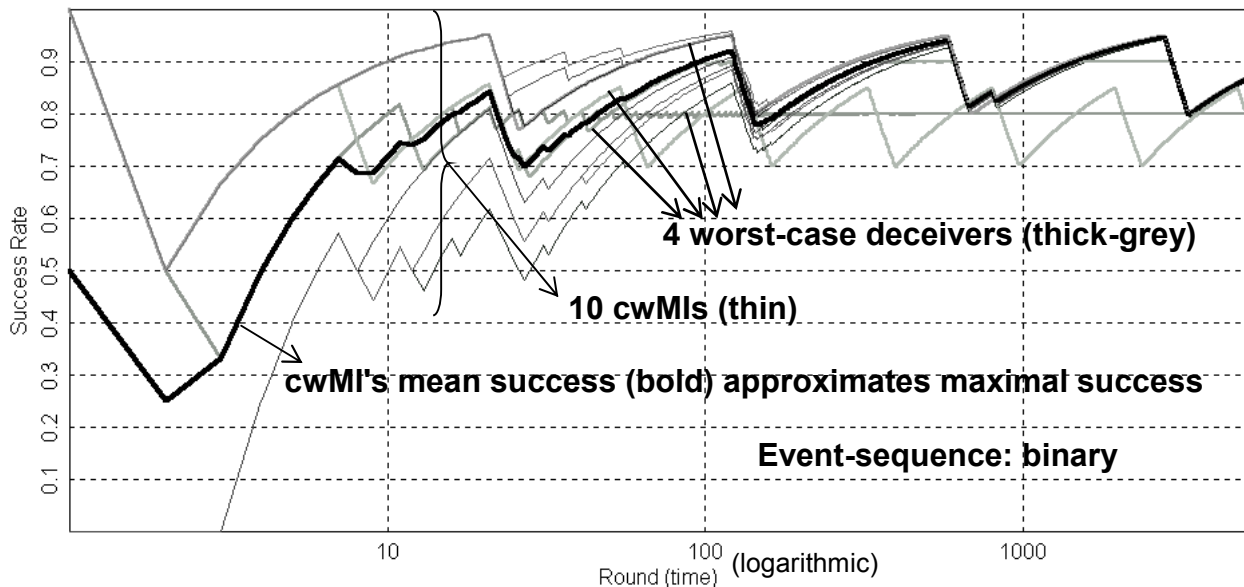


*Fig. 5: Ten cwMI's against four cwMI-adversaries.*

The relation of the cwMI-strategy to the situation of an *individual* meta-inductivist in binary prediction games is the following. The cwMI-adversaries may conspire against a particular individual, say against $cwMI_3$, and constantly deceive $cwMI_3$ (alternatively, a 'demonic' event sequence may constantly deceive $cwMI_3$). But the cwMI-adversaries cannot deceive the other cwMI's at the same time, and their anti-$cwMI_3$-conspiration will not affect the optimality result for the cwMI's mean success.

The collective actions of the cwMI's may be based on *rational agreement*. But it is sufficient to assume that the frequency of cwMI's who favor a non-MI-player P is proportional to P's attractiveness (provided this is positive). In the latter perspective, theorem 4 describes the predictive success of a population of cwMI's in a typical evolutionary setting as described in §7.

## 6. Conclusions for Epistemology

Table 1 summarizes the results on the optimality of meta-induction. Thereby, I relativize the notion of approximate optimality to a class of strategies Σ w.r.t. which xMI is optimal, and to a class of worlds **W** in which the strategies in Σ are played.

| xMI-strategy: | kind of optimality: | strategies in Σ are xMI-accessible and: | worlds in **W** contain finitely many xMI-accessible players satisfying: |
|---|---|---|---|
| MI (th.1) | strict opt. | no condition | ∃ best non-MI-player |
| εMI (th.2) | ε-opt. | no condition | ∃ set of ε-best non-MI-players |
| wMI (th.3) | strict opt. | no condition | real-valued game |
| cwMI (th.4) | $\frac{1}{2 \cdot k}$ -opt. | no condition | binary game |

*Table 1: Optimality of xMI w.r.t. Σ in **W***

I think the achieved results on the optimality of meta-induction are strong enough to show that a non-circular justification of meta-induction is possible. As I have explained in §2, this analytic justification of meta-induction implies an a posteriori justification of object-induction in our real word, because so far object-induction has turned out to be the most successful prediction strategy.

As I have remarked in §1, the meta-inductivists of table 1 are not generally dominant (w.r.t the respective Σ and **W**). The main reason for this fact is the existence of *refined* (meta-) inductive strategies, which I call their *conditionalized* versions. They exploit correlations in non-random worlds which obtain between the events $e_n$ and prior events, internal or external to the sequence (e), with help of Reichenbach's principle of the *narrowest reference class* (1949, § 72). Assume $\{R_1,…,R_r\}$ is a partition of the events prior to the given time, such that the given person has reliable information about which cell of $R_i$ was realized before the given time, and the cells are statistically relevant for the events of the sequence E (i.e., $\bar{e}_n |R_i \neq \bar{e}_n |R_j$ for j≠i, where

$\bar{e}_n | R_i$ is the $R_i$-conditionalized mean value of e up to time n. Let $R: |N \rightarrow \{R_1, \ldots, R_r\}$ be the function which assigns to each time n the cell R(n) of the partition which was realized before time n. Then the *conditionalized OI-strategy* transfers the conditionalized mean value $\bar{e}_n | R(n)$, and in the binary case its integer-rounding $[\bar{e}_n | R(n)]$, to the next time (note that $[\bar{e}_n | R(n)]$ coincides with the conditional frequency $f_n(E|R(n))$. Provided that all involved mean values and frequencies converge to a limit, one can prove that conditionalizing to reference partitions may only *improve* the success, compared to the simple OI (cf. Good 1983, ch. 17).

Also the conditionalized meta-inductivist works with a reference partition, but she conditionalizes the success frequencies of the other players to the cells of this partition. For example, the conditionalized εMI favors the first-best player P in the list of non-MI-players whose conditional success rate $suc_n(P|R(n))$ is maximal modulo ε, or the conditionalized wMI computes the attractivenesses of the non-MI-players in terms of their conditional success rates $suc_n(P|R(n))$. While a simple xMI ε-approximates the maximal success always from below, the success rate of a conditionalized xMI may be even strictly greater than the success rates of all other players. This fact does not affect the approximate optimality of the simple xMI, because we assume that refined meta-inductivistic techniques, if they are accessible, are among the methods of the alternative players. Hence with a "non-xMI-player" we mean a "non-simple-xMI-player". However, the fact shows that a simple meta-inductivist can *improve* his results by getting access to refined meta-inductivist (or object-inductivist) techniques. This will become important in the next section in which we consider evolution-theoretic applications of prediction games.

## 6. *Applications to the Evolution of Cognition*

In order for prediction games to make sense in an evolutionary setting, I change two interpretations of them and add one additional restriction.

*(Interpretation 1:)* I consider inductive strategies are strategies of *learning* within the individuals lifetime, while *non-inductive* strategies correspond to *genetically determined* strategies which cannot be modified by individual learning.

*(Interpretation 2:) Meta-inductivistic* strategies are strategies of learning from the performance of other successful individuals of one's population.

Under interpretation 2, the success of meta-inductivistic techniques reflects exactly the advantage of populations which possess the capability of *cultural* evolution, i.e. evolution by imitation and learning (in the sense of Richard Dawkins concept of 'memes'; cf. 1989, ch. 11). Several evolutionary theorists have discussed the question under which conditions *generation-wise* cultural evolution is superior to individual learning or to genetically determined strategies (cf., e.g., Boyd und Richerson 1985, 127). One such condition is, for example, that environmental condition do not completely change from one generation to the next one. However, my results on the optimality of meta-induction provide a *general* argument why the imitation of the successful members of one's population – which is the general basis of cultural evolution – brings an advantage to the mean success of the members of a population. In other words, meta-induction and cultural evolution are two sides of the same coin, and the optimality of meta-induction provides a general argument for the advantage of populations being capable of cultural evolution.

*(Additional restriction:)* Perfect clairvoyants, which have to be considered for the sake of the epistemological argument, do not play any realistic role in the evolution-theoretic setting. Evolutionary organisms are never perfect. I therefore assume the *constrain of imperfection* which says that for each non-MI-strategy there exist some environmental conditions in which its success is very low, and/or some environments in which their success is very high. Under this condition, the *conditionalized* meta-inductive strategy which I have explained before turns out to be not only weakly but even *strongly dominant*. In what follows, condMI denotes the conditionalized version of MI; it is strongly dominant (given the constraint of imperfection) under the condition of theorem 1; otherwise one has to use the conditionalized version of wMI). Fig.

6 illustrates conMI at hand of a prediction game with 5 different environments which change in the average after 50 rounds, but in an unpredictable way. conMI's success rates climbs high above the success rates of the non-MI-players, because in each given environment condMI takes advantage of exactly that strategy which performs best in *this* environment. For sake of comparison, fig. 6 informs also about success rate of the unconditional MI under the *hypothetical* assumption that condMI's predictions are *not accessible* to MI − otherwise MI would of course predict equally good as conMI (apart from a small short-run loss).
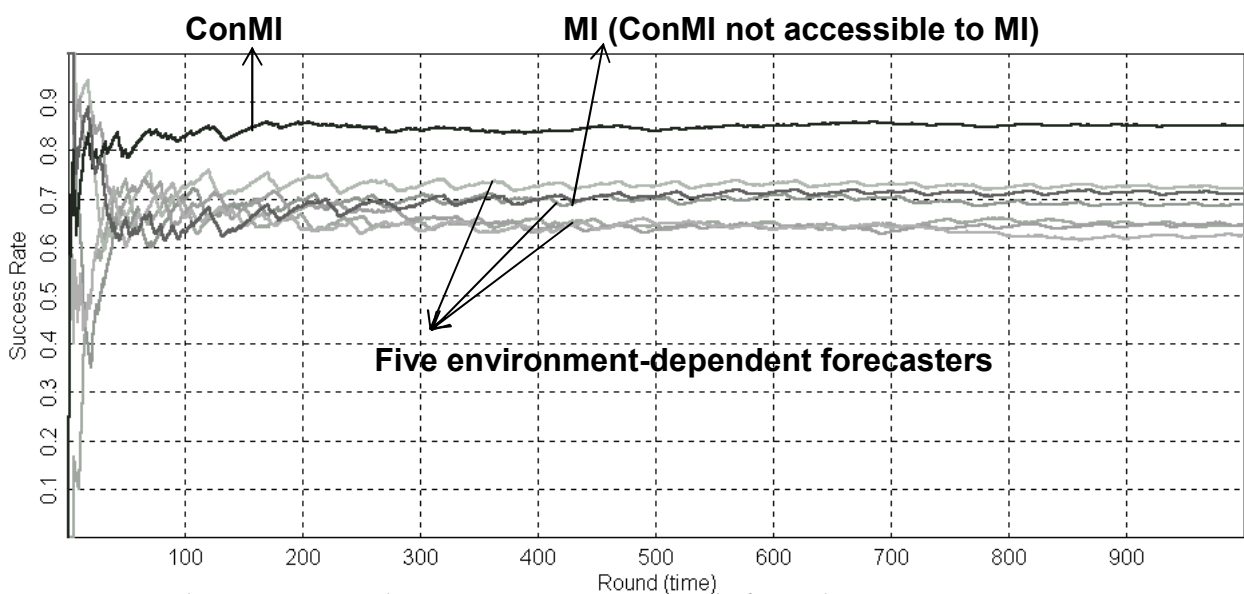


*Fig. 7: CondMI in an evolutionary scenario with five changing environments*

A further application to cognitive psychology arises when we compare the success of *one-favorite* meta-induction with the *weighted average* meta-induction. In fact, we find both strategies in humans as members of cultural evolution: some people tend to choose *one* authority whom they believe everything; others compare the opinions of several authorities and weigh them. What is the better strategy? Our results tell us that MI and TTB are not robust in situations of oscillating success rates. Indeed, humans who tend to believe one authority do *not easily* change the authority in which they believe − they behave more like εMI, which is a much more robust one-favorite strategy than MI. In the end we have seen that the superior though more complicated meta-induction method is weighted average meta-induction, although for binary pre-

diction games, this method does not work at the individual but only at the collective level.

*8. Local, General and Universal Prediction Strategies: Revising the Paradigm of Bounded Rationality*

The *bounded rationality* paradigm of the ABC-group maintains that *all* prediction methods, whether complex or simple, are *ecological* in the sense that their success is restricted to certain environmental conditions. We have seen that this claim is not generally valid. I want to suggest a more refined picture, by distinguishing between three kinds of prediction strategies with respect to the range of worlds in which they are optimal.

(1.) *Local* strategies are those whose optimality is restricted to *particular* kinds of situations, or worlds. They presuppose that certain correlations hold and fail whenever these correlations don't hold. Local strategies *cannot learn*, they are genetically determined or elsewhere dogmatically fixed strategies. For example, the strive of an insect towards light is a local strategy; the insect cannot change this strategy in an environment of electric light bulbs in which this strategy kills it. (Note that I speak of learning at the level of the individual, not at the global level of evolution itself). Another example of a local strategy is Gigerenzer's *recognition heuristics* (1999, p. 37ff.), insofar this strategy presupposes that the layman's recognitions of events are statistically correlated with certain comparative properties of these events; e.g. with the size of recognized cities. For example, it would be a bad recommendation to apply the recognition heuristics to the mathematical intelligence of recognized persons or the attractiveness of recognized tourist places in nature.

(2.) Object-inductive strategies *can learn*: they can change their beliefs about correlations and cues in the light of new evidence. Hence, the optimality of these strategies is certainly not local, but they are more-or-less *general* (with 'general' I do not mean 'strictly universal' but general in the 'more-or-less' sense). The existence of cog-

nitive *general-purpose* mechanisms has also been emphasized by cognitive psychologists (e.g. Over 2003). Nevertheless, all object-inductive strategies fail in certain worlds, whence they are not universally optimal.

(3.) Universally optimal would be a prediction strategy if it predicts optimal in *every* possible environmental condition. According to the *bounded rationality* paradigm of the ABC-research group, universally optimal strategies don't exist. We have seen that although this is true for object-inductive strategies, it is false for meta-inductive strategies, since weighted average meta-induction is indeed universally optimal among all *accessible* prediction methods. The restriction to "accessibility" is, of course, crucial, and without this restriction a universal optimality result can impossibly be achieved.

## *References*

Boyd, R. and Richerson, P. J. (1985), *Culture and the Evolutionary Process*. Chicago: Univ. of Chicago Press.

Cesa-Bianchi, N., and Lugosi, G. (2006): *Prediction, Learning, and Games*, Cambridge Univ. Press, Cambridge.

Dawkins, R. (1989), *The Selfish Gene*, 2nd edition. Oxford: Oxford Univ. Press.

Friend, M., Goethe, N.B., and Harizanov, V.S. (2007, eds.), *Induction, Algorithmic*

Gigerenzer, G., et al. (1999): *Simple Heuristics That Make Us Smart*, Oxford Univ. Press, Oxford.

Good, I. J. (1983): *Good Thinking*, Univ. of Minnesota Press, Minneapolis.

Hogarth, R.M., and Karelaia, N. (2006): " 'Take-The-Best' And Other Simple Strategies", *Theory and Decision* 61, 205-249.

Kelly, K.T. (1996): *The Logic of Reliable Inquiry*, Oxford Univ. Press, New York.

Merhav, N., and Feder, M. (1998): "Universal Prediction", *IEEE Transactions on Information Theory* 44(6), 2124-2147.

Over, D.E. (2003): "From massive modularity to meta-representation: the evolution of higher cognition", in: D.E. Over (ed.), *Evolution and the psychology of thinking: The debate,* Psychology Press, . Hove, UK, 122-140.

Putnam, H. (1963): "Probability and Confirmation", in: *The Voice of America. Forum Philosophy of Science,* reprinted in: Putnam, H., *Mathematics, Matter, and*

*Method, Cambridge Univ. Press, Cambridge 1975.*

Reichenbach, H. (1938): *Experience and Prediction*, University of Chicago Press.

Reichenbach, H. (1949): *The Theory of Probability*, University of California Press.

Rescher, N. (1980): *Induction*, University of Pittsburgh Press.

Salmon, W. C. (1957): "Should We Attempt to Justify Induction?", *Philosophical Studies* 8, No. 3, 45-47.

Salmon, W. (1974). "The Pragmatic Justification of Induction", in: R. Swinburne, *The Justification of Induction*, Oxford University Press, Oxford, 85 – 97.

Schurz, G. (2008): "A Game-Theoretical Approach to the Problem of Induction", to appear in: C. Glymour et al. (eds.), *Proceedings from the 13th International Congress of Logic, Methodology and Philosophy of Science*, King's College Publications, London.

Skyrms, B. (1975): *Choice and Chance*, Dickenson, Encinco (4th ed. Wadsworth 2000).

Weibull, J. (1995). *Evolutionary Game Theory*, MIT Press, Cambridge/Mass.