# Detecting match-fixing in tennis

**Tim Paulden** explains how spotting anomalous movements in the betting markets
can help shine a light into the murky world of tennis match-fixing

*This is a heavily reduced version of the article appearing in the June 2016 issue of Significance Magazine.*

The research that ATASS Sports have conducted into tennis match-fixing since early 2015 represents our first foray into the fledgling field of *forensic sports analytics* – the application of statistical modelling to help identify and eliminate corruption within the sports sector.

To set the scene, consider the following question: How might an unscrupulous tennis player try to make money through match-fixing? Perhaps the most obvious scheme would be for him to agree in advance to deliberately lose a certain match, and have a complicit third party place large bets on *his opponent* winning. This simple strategy is understood to be one of the most likely modes of tennis match-fixing – particularly at lower echelons of the sport, where the size of the betting market for a single first-round match may dwarf the prize money for an entire tournament.

However, the perpetrators of such a fix might inadvertently leave behind a crucial fragment of evidence: the perturbations to the odds caused by their bets. To explain briefly, whenever a betting market experiences a large influx of money in support of one player, that player's odds will tend to shorten (become less generous) – just as a share price will rise if demand exceeds supply. Since our shady third-party bettor is already virtually certain of the match result, he may be willing to put down large bets at odds that would normally be seen as unfavourable – leading to the odds becoming distorted in an anomalous way.

In our research, we decided to focus on spotting anomalies in the "in-play" betting markets rather than the "pre-match" markets, since in-play volumes now often exceed pre-match volumes by a factor of 10 or more – making it highly likely that the majority of suspicious third-party betting would occur in-play. (It is worth noting that in early 2016, a large pre-match odds data set was analysed by Buzzfeed News, who concluded that there were 15 players who regularly lost matches in which there had been lopsided pre-match betting activity. However, the methodological shortcomings of this analysis received significant criticism from the well-respected DW on Sport blog, and later from the Guardian.)

Of course, identifying anomalous in-play odds movements is tricky because the odds for each player will naturally fluctuate during the match due to the changing scoreline. To spot when the odds are shifting anomalously, we must get a handle on how each player's win probability *should* have evolved during the match, given the sequence of observed points.

As explained in the unabridged version of this article (published in the June 2016 edition of Significance Magazine), our experiments led us to investigate an extremely parsimonious point-by-point model in which the two players – A and B,  say – each have a fixed probability of winning a point when serving, given respectively by
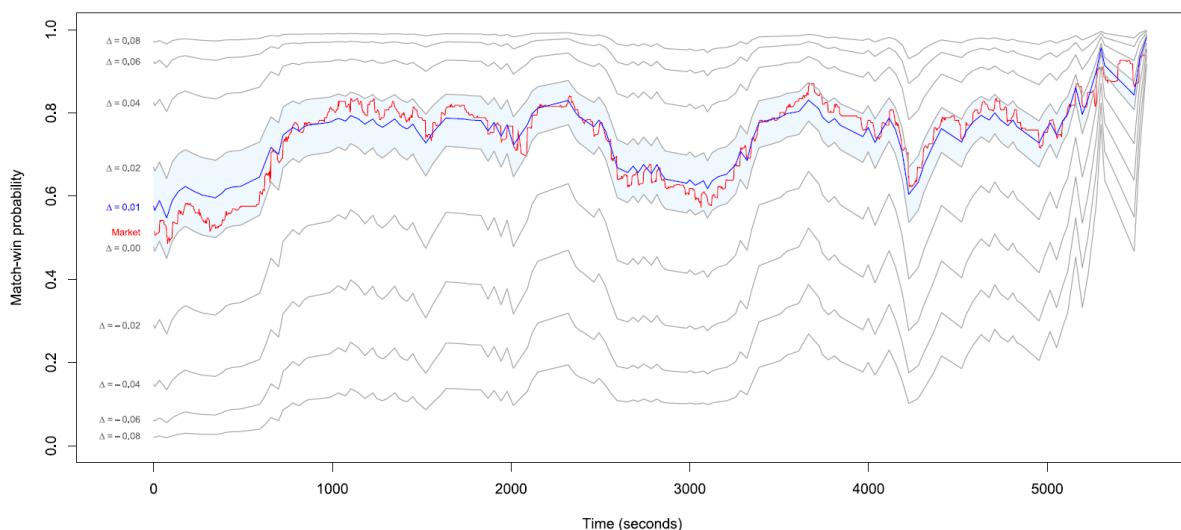
$$p_A = 0.645 + \Delta$$
$$p_B = 0.645 - \Delta$$

Here, $\Delta$ is a "dominance parameter" that encodes how much better one player is than the other, while 0.645 represents an average on-serve point win probability for men's tennis.

This simple assumption turns out to be surprisingly powerful, because given the values of $p_A$ and $p_B$ for a particular match, we can apply the tennis formulae of O'Malley (2008) to determine the players' respective match win probabilities – not only at the beginning of the match, but in every possible match scenario.
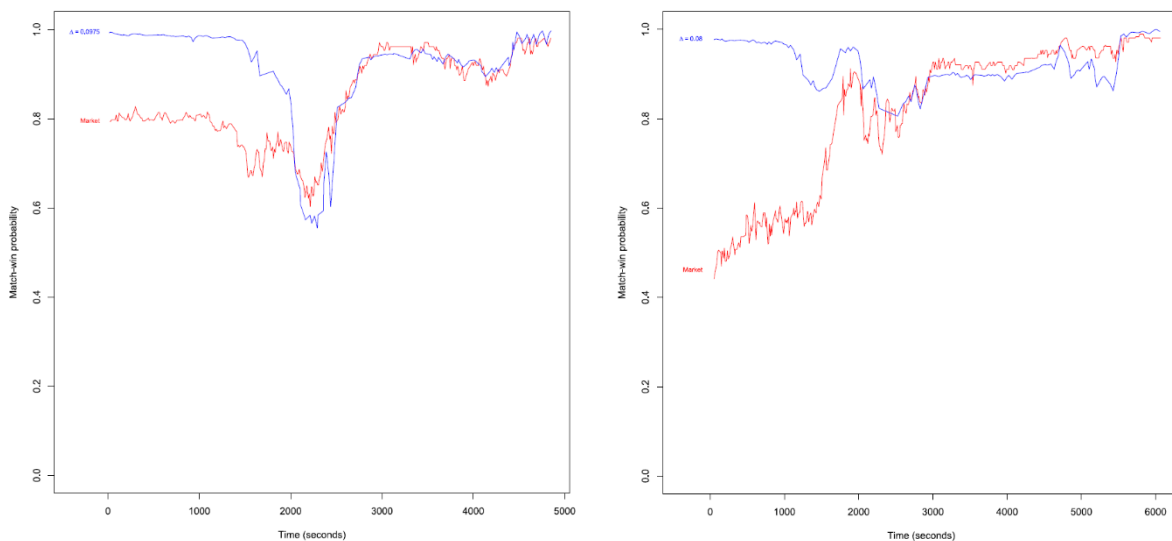
How well can such a simple one-parameter model capture the evolution of a real betting market? The figure below depicts a typical match from July 2014, with the red line showing the evolution of the market's win probability for player A at each point of the match, and the nine grey lines showing how the model win probability for player A varies for different values of $\Delta$ (from –0.08 at the bottom to 0.08 at the top). It is natural to define the "best" $\Delta$ to be that minimising the *discrepancy score* – i.e. the mean absolute difference between the market probability and model probability over the course of the match. In this case, the minimum discrepancy score of 0.0232 arises for $\Delta = 0.01$, and is indicated by the blue line.



Following the above framework, we investigated an archive of around 5000 ATP matches spanning the period from 2013 to mid-2015, and all levels down to the Challenger Tour, and found the best $\Delta$ for each one by minimising the discrepancy score. The results we obtained were extremely encouraging, with the optimal discrepancy score being small for almost all matches examined (80% of matches had a score less than 0.03; 99% of matches had a score less than 0.06). These results provide powerful empirical evidence that even an extremely

parsimonious model can provide a reasonable approximation to the evolution of a betting market across the vast majority of tennis matches.

Naturally, not all matches with large discrepancy scores represent instances of match-fixing: the anomaly could be due to an injury, or any number of extraneous factors. However, these are the matches that merit closer scrutiny. When we examined the matches in the rightmost 1% of the distribution (with a discrepancy score exceeding 0.06), we found there were matches in which the odds evolved in a highly irregular fashion that we could not rationally explain. The graphs from two of these matches (one played in August 2014, and one played in February 2015) are shown below.



Under our model, there is no value of Δ that is remotely consistent with the evolution of the betting market depicted in these graphs. In fact, for both matches, the eventual winner would have needed to be a "dead cert" at the start of the match (as shown by the starting point of the blue model line) for the market probabilities observed later on to make sense. In the second of these matches, the market jolted so irregularly that the eventual winner was deemed significantly more likely to win when trailing by a set than they had been at the start of the match – a clearly absurd situation.

Moreover, when we investigated the background to these particular matches in detail, we found that in both cases, the evolution of the betting market had been identified as being highly irregular by numerous tennis blogs – including DW on Sport – and other websites, such as Slate.com. In other words, our system had been able to successfully filter down a collection of several thousand matches to a small subset that included those specifically flagged as suspicious by tennis experts.

*Since our initial work was undertaken in summer 2015, this stream of research has blossomed into a three-year collaborative PhD project with Lancaster University on the topic of "forensic sports analytics".*