

A Brief Investigation into Tit-for-Tat's Performance
Against Human Opponents, with Short Game Lengths

Alexander Mieczyslaw Kasprzyk

1997

Introduction

The original form of the Prisoner's Dilemma can be best described with a little example. Two prisoners are being held in separate rooms and the police are interrogating them about a crime they jointly committed. Each prisoner has a choice – should he inform on his partner's actions ("defect" from his partner's perspective) or remain silent ("cooperate" with his partner). If both of the prisoners remain silent (cooperate) they will both be released. The police, however, are aware of the situation and offer an incentive: if one of them defects then the prisoner will be granted immunity and be given a reward for his troubles. His partner, however, will be assigned a sentence twice as harsh as that which would be given had they both cooperated. Of course, if both the prisoners defect they will both serve the maximum sentence and neither gets a reward.

What should the prisoners do? Obviously the best tactic is to defect, since you stand to lose little and gain everything. However, when played over several turns the game becomes interesting. Tactics begin to play an important role. How can you maximise your winnings?

This is the basic form of the Iterated Prisoner's Dilemma ("iterated" because it is played over several turns). In fact it comes in many more flavours. The version I am going to concentrate on is the "selfish" version, where the aim is to minimise your opponent's score. At first glance this appears easy, since if you constantly defect your opponent will never receive any rewards, however this is against the rules of the "selfish" game. A selfish player must respond in some way to the other player's previous actions – a simple "always defect" strategy is not allowed because it does not react in any way to the opponent's actions.

What is the best tactic?

In the 1970s a computer tournament was organised at Michigan by Robert Axelrod to find the answer. In this tournament people were invited to submit programs to play a game of Prisoner's Dilemma which would last for 200 moves. Fourteen programs were submitted for the first round, and much to the amazement of Axelrod the crown went to a very simple strategy which has become known as the "Tit-for-Tat" strategy. Submitted by psychologist Anatol Rapoport, tit-for-tat would start out with a purely random first move (though in some versions the first move is always to cooperate), and from then onwards would simply imitate the opponent's previous move.

A typical game against tit-for-tat would look something like that shown in figure 1. Here the pay-off matrix reflects the scores for the player only. In a non-selfish game we would also be required to consider the scores obtained by tit-for-tat. The scores in the pay-off matrix can be varied depending on choice, but traditionally they are values similar to those shown in the figure. Throughout this project I will use the same scores as those shown in the example.

Player's Pay-off Matrix		
		Tit-for-Tat
		C D
Player	C	1 0
	D	10 0

Tit-for-Tat	D	C	D	D	C	D	C	C	C	D
Player	C	D	D	C	D	C	C	C	D	D
Score	0	10	0	0	10	0	1	1	10	0

Player's Final Score = 32

Fig.1 A typical Tit-for-Tat game

Aims

My first aim in this project is see whether tit-for-tat performs equally exceptionally when playing against human opponents rather than computer opponents. Since a human would soon grow tired if required to play a game for 200 moves as in the Axelrod tournament, we will not only be playing tit-for-tat against humans but also playing it for shorter length games than previously.

I will test for a difference in means of the final scores a selection of five different algorithms (tit-for-tat included) when played against twelve human players. If an analysis of variance test shows any difference between the means I will then go on to ascertain which algorithm performed the best. My prediction based on the Axelrod tournaments is that tit-for-tat should have the highest mean score. It may also be useful to consider whether the variance (ie. consistency of score) varies between tit-for-tat and other players.

The second aim is to investigate the correlation between game length and total score. Is the relationship linear and, if so, can the final score be predicted if the length of the game is known? This will be achieved using linear regression, and a test of whether the regression coefficient is different from the theoretical coefficient for a purely random player will also be performed.

To test whether tit-for-tat performs exceptionally when playing against humans over short periods of time it will be necessary to perform an experiment. We will first need something to gauge the performance of tit-for-tat against. I have selected three other strategies apart from the tit-for-tat method from the original Axelrod tournament. These three strategies can be loosely described as:

Overall Tit-for-Tat – if the player has defected more times than they have cooperated then the computer will defect, otherwise it will cooperate.

Other turn Tit-for-Tat – if the player has defected more than once in the past three moves then the computer will defect, otherwise it will cooperate.

Advantage Player - the computer attempts to spot patterns in the style of play by looking at consecutive moves. For example, if the player tends to defect after cooperating then the computer will attempt to predict this.

I will also include a purely random algorithm which, although strictly speaking invalidates the rules for selfish Prisoner's Dilemma because it does not respond to the opponent's previous moves, will provide an invaluable benchmark against which to gauge tit-for-tat's performance.

Data Collection and Experimental Design

I have readily available twelve people who are willing to take part in my trials.

What factors other than the algorithm they face could effect their final score? The length of the game played is an obvious factor. Longer games will have larger total scores than shorter games. This can be removed by keeping the length of the games fixed at, say, twenty moves.

As people play they improve and develop tactics. This is a factor which may well bias the final scores in favour of the last few games played. The effect of this can be minimised by blocking the players in such a way as to vary the order in which the opponents are played. Table 1 shows the table I created to try and minimise this problem. The numbers in the table reflect the order in which the human will face the various opponents. For example, human player number five will first play tit-for-tat and work their way through the other opponents, finishing with the random player.

One of the problems I faced designing this table was the fact that the number of opponents and the number of human players are not equal. If they were I could draw up a latin square which would effectively remove any possible bias. Because the numbers are not equal I could either choose to cut down my number of human opponents to five (which I felt would be a mistake since it drastically reduces the sample size), increase the number of computer opponents to twelve (which was impractical due to time restrictions and the patience of the people taking part in the experiment), or assign values myself in such a way as to try and remove any bias.

		Order of Play				
		1	2	3	4	5
Human Player	1	i	ii	iii	iv	v
	2	v	i	ii	iii	iv
	3	iv	v	i	ii	iii
	4	iii	iv	v	i	ii
	5	ii	iii	iv	v	i
	6	i	ii	iii	iv	v
	7	v	i	ii	iii	iv
	8	iv	v	i	ii	iii
	9	iii	iv	v	i	ii
	10	ii	iii	iv	v	i
	11	i	v	ii	iii	iv
	12	ii	iv	iii	i	v

i - Random player
 ii - Tit-for-Tat player
 iii - Overall Tit-for-Tat player
 iv - Other turn Tit-for-Tat player
 v - Advantage player

Table 1

The pattern up to player 10 will remove any bias by varying the order of play in a regular nature. Players number 11 and 12 did not easily fit into this pattern, so I decided to assign these players opponents in a random order taken from the first row of my random number tables. This was done as shown in figure 2.

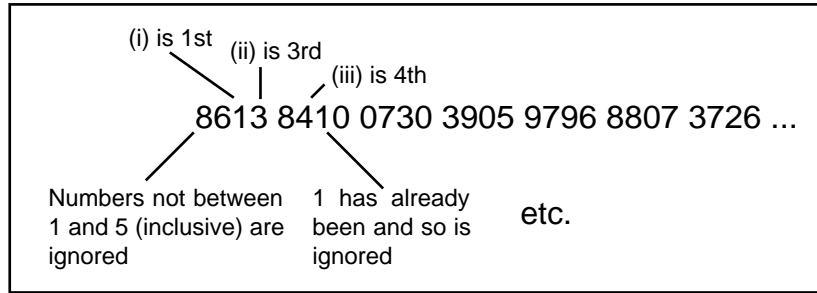


Fig. 2

By varying the order in which the opponents were played in this manner I hoped to remove any bias due to the human player's ability improving with time.

Conversely, the human player may grow bored as time progresses. By keeping the length of each game down to twenty moves and restricting the number of opponents to five I hoped to keep people interested. Also, any variation due to tiredness will be accounted for in the same way I accounted for improvements in playing tactics with time.

To prevent any bias I assigned each human a number at random by selecting their names out of a bag. The first person who's name emerged was human player number one, the second was number two, etc. These numbers are used through out the project (eg. table 1).

The people taking part in this experiment were all friends of mine and, as such, knew of my interest in the Prisoner's Dilemma and many of them knew about the tit-for-tat algorithm. This posed a serious problem when collecting results. To remove any bias these individuals could consciously or unconsciously introduce into my results I needed to make sure people were unaware as to which opponent they were facing. This blind play was made possible by programming a computer with the order in which each player should face each opponent and withholding this information from the player. Figure 3 shows a snap-shot of a game. Notice how there is no indication of what opponent the human is facing.

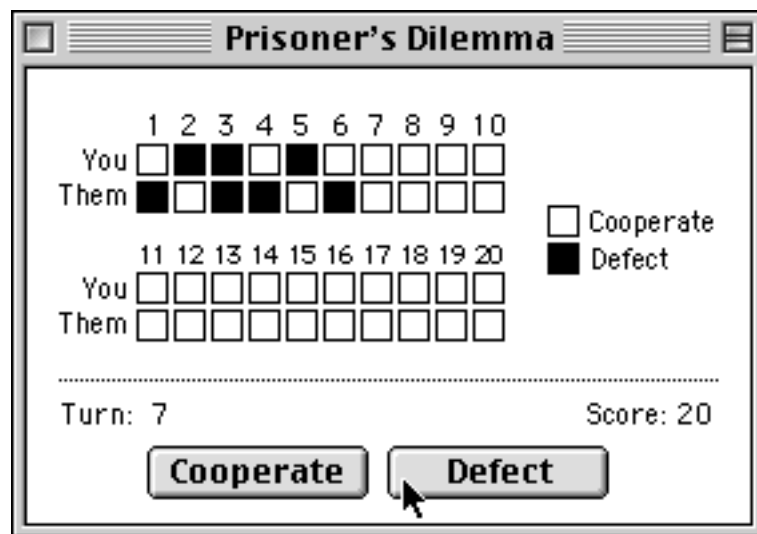


Fig. 3

The response variable I will be measuring will, of course, be the total score of each game.

Obtaining the experimental data for the second project aim will require more work. It will be necessary to vary the length of the games so that a range of readings are recorded. I feel that a

range of game lengths between 10 and 29 moves in length (inclusive) should be acceptable. Very few people would be willing to play games which last much longer than that.

How should I collect this data? Should I give each person a game to play against tit-for-tat which lasts a different length, or should the same person play all the games? Before I can answer this question it will be necessary to see whether there is any significant difference between the various human player's scores. I will be able to use the within sample variance from an analysis of variance test on the first set of data collected as a guide. If this variance level is noticeably large, I will need to use the same person for the entire set of games, otherwise I can divide up the various games at random amongst the 12 human players.

How Does Tit-For-Tat Perform?

The data in table 2 was collected as described in the section above on experimental design.

Table 2

	i	ii	iii	iv	v
1	53	46	76	4	121
2	76	41	11	11	15
3	62	75	25	12	135
4	22	44	96	32	63
5	65	57	117	44	83
6	49	65	12	31	57
7	43	46	33	10	32
8	37	49	68	32	12
9	47	53	3	13	58
10	54	67	12	4	72
11	47	65	24	44	45
12	66	45	23	22	25

i - Random player
 ii - Tit-for-Tat player
 iii - Overall Tit-for-Tat player
 iv - Other turn Tit-for-Tat player
 v - Advantage player

Human Player

Final score of 12 human players against five different styles of play

Before any statistical tests can be performed on the data it is necessary to determine whether the scores obtained by each human player against a specific computer opponent follow a normal distribution or not.

Normally it is possible to carry out a χ^2 goodness of fit test to see whether the values follow a normal distribution. In this case, however, this is not really possible due to the small sample size. Such a test would require the combining of so many classes that we can effectively rule out even attempting to perform such a test on the data.

Fortunately there exists a rather crude but practical alternative. By plotting the values on normal graph paper we can say whether or not the data follows a normal distribution. First, each sample must be sorted into order with the smallest first, the largest last. Each of these values is plotted against a corresponding probability calculated by:

$$\frac{n_i}{n_{total} + 1} \quad \text{where } n_i \text{ is the position the value comes in the ordered list}$$

$$n_{total} \text{ is the total number of items in the list}$$

For example the 5th item from a sample of 12 would have a corresponding probability of $5/13 = 0.385$ (to 3 d.p.). Table 3 shows the values I plotted on my normal graph paper.

Table 3

p	i	ii	iii	iv	v	p	v
0.08	22	41	3	4	12	0.09	12
0.15	37	44	11	4	15	0.18	15
0.23	43	45	12	10	25	0.27	25
0.31	47	46	12	11	32	0.36	32
0.38	47	46	23	12	45	0.45	45
0.46	49	49	24	13	57	0.55	57
0.54	53	53	25	22	58	0.64	58
0.62	54	57	33	31	63	0.73	63
0.69	62	65	68	32	72	0.82	72
0.77	65	65	76	32	83	0.91	83
0.85	66	67	96	44	121		
0.92	76	75	117	44	135		

Scores converted for plotting on normal graph paper

When I plotted the values for opponent (v) I could see that the sample would follow a normal distribution with the exception of last two values in the list. These values certainly appear to be outliers, however I can find no explanation for their unusually high values. For the sake of normality I decided to remove them and recalculate the probabilities for (v) to test for normality. This new smaller sample is normal.

As a result of plotting the points on normal graph paper normality can be assumed for (i), (ii), (iv) and the smaller sample for (v). We can now proceed with performing an analysis of variance test on these samples to see whether there is any difference between the mean scores of these four different styles of play.

If tit-for-tat truly performs better than the other algorithms would expect some difference in the means to exist. Should such a difference be confirmed by the analysis of variance test we can then go on to discover for which styles of play this difference exists. If the analysis of variance test concludes that there is no difference between the means we can say that tit-for-tat does not perform significantly differently against human opponents over short game lengths.

Table 4

	i	ii	iv	v	
1	53	46	4	15	
2	76	41	11	63	
3	62	75	12	83	
4	22	44	32	57	
5	65	57	44	32	
6	49	65	31	12	
7	43	46	10	58	
8	37	49	32	72	
9	47	53	13	45	
10	54	67	4	25	
11	47	65	44		
12	66	45	22		
Total	621	653	259	462	= 1995

$$SS_b = \frac{T_1^2}{n_1} + \frac{T_2^2}{n_2} + \dots + \frac{T_k^2}{n_k} - \frac{G^2}{n} = \frac{621^2}{12} + \frac{653^2}{12} + \frac{259^2}{12} + \frac{462^2}{10} - \frac{1995^2}{46} = 8083.03 \quad \text{d.f.}_b = 3$$

(to 2 d.p.)

$$SS_t = \sum x^2 - \frac{G^2}{n} = 105973 - \frac{1995^2}{46} = 19450.72 \quad \text{d.f.}_t = 45$$

$$SS_w = SS_t - SS_b = 11367.68 \quad \text{d.f.}_w = 42$$

$$\hat{\sigma}_b^2 = \frac{8083.03}{3} = 2694.34 \quad \text{(to 2 d.p.)}$$

$$\hat{\sigma}_w^2 = \frac{11367.68}{42} = 270.66 \quad \text{(to 2 d.p.)}$$

H_0 : there is no difference between the means of the different styles of play

H_1 : the mean of one or more of the styles of play is different

One tailed test at the 5% level
Critical value from $F_{3,42}$ is 2.83

$$F_{\text{test}} = \frac{2694.34}{270.66} = 9.95 \quad \text{(to 2 d.p.)}$$

Since $9.95 > 2.83$ reject H_0 and conclude that there is enough evidence at the 5% level to say that there is a significant difference between the mean score of one or more of the styles of play.

Several interesting results can be gained from this analysis of variance test. Firstly, and perhaps most importantly with regards to the design of the experiment and the second aim of the project, the within sample variation is low. This is a good sign since a high within sample variance would suggest a poorly designed experiment because a factor other than those accounted for in the design is playing a key role. The low within sample variance also suggests that there is no real difference between the scores of the various human players and so the task of playing games with a range of different lengths can be divided up between them.

The analysis of variance test confirmed that there was some significant difference between the mean score of one of more of the various styles of play. This would be a good sign if it were not for the fact that from simply looking at the total scores for the four different methods we can

suggest that tit-for-tat actually has one of the total highest scores. Since the scores we are looking at are the human player's total score for each game, if tit-for-tat was performing better than the other algorithms we would be expecting a lower total score from the rest, not a higher total score. More disturbingly, it would appear that tit-for-tat is no more effective than a simple random algorithm since there is very little difference between their total scores.

We can test whether tit-for-tat performs any differently (whether it be better or worse) than a random player by performing a paired sample difference in means test. The test is for paired samples because each human played against the two algorithms. Before we perform this test, however, we will test for a difference between tit-for-tat's scores and algorithm (iii) scores. Because algorithm (iii) does not follow a normal distribution (and hence did not feature in the analysis of variance test) we will have to make use of non-parametric tests.

A sign test for paired data is a quick and simple test to see whether tit-for-tat's median score is less than the median of (iii). This is the result we would expect if tit-for-tat is a more effective player than algorithm (iii).

Table 5

	Tit-for-Tat (X)	Type iii (Y)	Sign of X - Y
1	46	76	-
2	41	11	+
3	75	25	+
4	44	96	-
5	57	117	-
6	65	12	+
7	46	33	+
8	49	68	-
9	53	3	+
10	67	12	+
11	65	24	+
12	45	23	+

let R = number of +'s

H_0 : no difference between the players (median of differences is 0)

H_1 : Tit-for-Tat's median is lower than Type iii (median is less than 0)

One tailed test at the 5% level

Under H_0 , $R \sim \text{Bin}(12, 0.5)$ $r = 8$

$$P(R \geq 8) = 1 - 0.8062 = 0.1938$$

Since $0.1938 > 0.05$ accept H_0 and conclude that there is enough evidence at the 5% level to say that there is no difference between the two players (Tit-for-Tat and Type iii).

The conclusion drawn from this rather crude Sign Test is not the result I was expecting. I feel that it is therefore worth investigating this further by performing a more accurate Wilcoxon Signed Rank Test. However, since the sign test concluded that there was no difference between the medians of the two samples, and because the evidence from the analysis of variance test is suggesting that tit-for-tat is in fact performing at a worse than expected standard

when compared with the other algorithms, I will now test to see whether there is any difference between the two algorithms. If the test shows that there is a difference, then we will go on to consider what form this difference takes.

Table 6

	Tit-for-Tat	Type iii	Difference	Rank
1	46	76	-30	-4.5
2	41	11	30	4.5
3	75	25	50	7.5
4	44	96	-52	-9
5	57	117	-60	-12
6	65	12	53	10
7	46	33	13	1
8	49	68	-19	-2
9	53	3	50	7.5
10	67	12	55	11
11	65	24	41	6
12	45	23	22	3

$$S_- = 27.5$$

$$S_+ = 50.5 \quad \therefore S = 27.5$$

H_0 : no difference between the two types of play

H_1 : some difference between the two types of play

Two tailed test at the 10% level

Critical region is $S \leq 17$

Since $27.5 > 17$ accept H_0 and conclude that there is enough evidence at the 10% level to say that there is no difference between the two styles of play (Tit-for-Tat and Type iii).

The results from these two tests confirm quite clearly that tit-for-tat is not performing any better than algorithm (iii). There is no difference between the two styles of play.

I will now go on to perform the difference in means test we discussed above. Because the sample size is small and σ^2 is estimated we will have to perform a t-test. The result from the analysis of variance test, combined with the result from the Wilcoxon Signed Rank Test, suggest that tit-for-tat is not performing as spectacularly as we had hoped. In fact it is looking as though tit-for-tat is actually quite a weak algorithm when it comes to playing short games against human opponents.

Table 7

	Random	Tit-for-Tat	d	d ²
1	53	46	7	49
2	76	41	35	1225
3	62	75	-13	169
4	22	44	-22	484
5	65	57	8	64
6	49	65	-16	256
7	43	46	-3	9
8	37	49	-12	144
9	47	53	-6	36
10	54	67	-13	169
11	47	65	-18	324
12	66	45	21	441
			-32	3370

$$\bar{x} = \frac{-32}{12} = -2.667 \text{ (to 3 d.p.)}$$

$$s^2 = \frac{3370}{12} - \frac{(-32)^2}{12^2} = 273.722 \text{ (to 3 d.p.)}$$

$$\hat{\sigma}^2 = \frac{12 \times 273.722}{11} = 298.606 \text{ (to 3 d.p.)}$$

H₀: μ_x = μ_y (the means of the two populations are the same)

H₁: μ_x ≠ μ_y (the means of the two populations are different)

Two tailed test the the 5% level

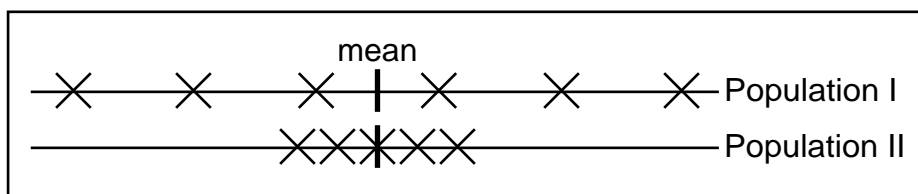
Critical values are ±2.201 (v = 11)

$$t_{\text{test}} = \frac{-2.667}{\sqrt{\frac{269.606}{12}}} = -0.535 \text{ (to 3 d.p.)}$$

Since -0.535 > -2.201 accept H₀ and conclude that there is enough evidence at the 5% level to say that the means of the two populations (Random player and Tit-for-Tat player) are the same.

The results of this test are even more disturbing. It would appear that not only does tit-for-tat perform exceptionally poorly when playing against human opponents rather than computer opponents, but that in fact the effectiveness of tit-for-tat is no better than that of a purely random player.

Mean scores are not the only factor that could set tit-for-tat apart from the rest. It is possible that, although on average tit-for-tat performs poorly, it performs at a consistent standard. It is clear that a smaller variance (ie. a more consistent player) is preferable. Figure 4 illustrates this.



An illustration of why a smaller variance is preferable

Fig. 4

To test whether tit-for-tat has indeed got a smaller variance than a purely random player we will perform a Fisher test. One of the criteria for an F-test is that the two samples are independent, but we have already asserted that they are in fact paired. I will therefore perform the F-test on a random sample of six of the tit-for-tat results and the remaining six results from the random player.

Before we perform the test it is worth looking at the range of scores on a simple diagram. Perhaps there will be an obvious difference in variances and we will not need to perform the test. Figure 5 illustrates the scores in this way. From looking at the figure it is tempting to conclude that the random player does have a greater range of scores, however I still feel that it is worth performing an F-test.

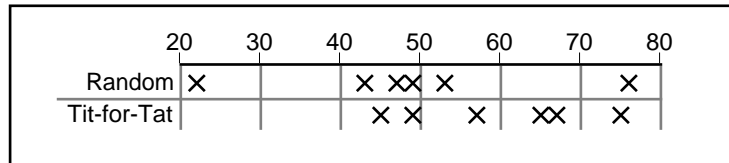


Fig. 5

Table 8

Random (x)	Tit-for-Tat (y)
53	75
76	57
22	49
49	67
43	65
47	45
290	358

$$\bar{x} = \frac{290}{6} = 48.333 \text{ (to 3 d.p.)}$$

$$s_x^2 = \frac{15528}{6} - \frac{290^2}{6^2} = 251.889 \text{ (to 3 d.p.)}$$

$$\hat{\sigma}_x^2 = \frac{6 \times 251.889}{5} = 302.267 \text{ (to 3 d.p.)}$$

$$\bar{y} = \frac{358}{6} = 59.667 \text{ (to 3 d.p.)}$$

$$s_y^2 = \frac{22014}{6} - \frac{358^2}{6^2} = 108.889 \text{ (to 3 d.p.)}$$

$$\hat{\sigma}_y^2 = \frac{6 \times 108.889}{5} = 130.667 \text{ (to 3 d.p.)}$$

$H_0: \sigma_x = \sigma_y$ (the variances of the two populations are equal)

$H_1: \sigma_x > \sigma_y$

One tailed test at the 5% level

Upper critical value from $F_{5,5}$ is 5.05

$$F_{\text{test}} = \frac{302.267}{130.667} = 2.313 \text{ (to 3 s.f.)}$$

Since $2.313 < 5.05$ accept H_0 and conclude that there is enough evidence at the 5% level to say that the variances of the two populations (Random player and Tit-for-Tat player) are equal.

This test has yet again demonstrated that tit-for-tat is not the exceptional player Axelrod's tournaments suggested it would be. We can now quite confidently say that tit-for-tat performs very poorly against a human opponent over a short game length. In fact the evidence suggests that tit-for-tat performs no differently from a purely random player, having a similar mean and variance.

What is the Relationship between Game Length and Total Score?

To investigate the relationship between game length and total score data must be collected where games varying in length from 10 moves to 29 moves are played against tit-for-tat. Because of the low within sample variance obtained by the analysis of variance test we can safely use different human players to help collect the data. Twenty games needed to be played. These were assigned to the twelve human players (they were assigned by simply giving player 1 a game of length 10, player 2 a game of length 11, etc. When we ran out of players we started with player one again. Ideally twenty human players would have been preferable, but this was not possible). The data collected from these games is recorded in table 9.

Game Length (x)	Score (y)
10	50
11	23
12	25
13	25
14	15
15	41
16	48
17	51
18	33
19	32
20	45
21	59
22	67
23	60
24	62
25	67
26	78
27	87
28	81
29	75

Table 9

Before we perform any regression on this data, let us look at the theory behind a tit-for-tat game and the relationship we would expect to see emerge if a human player decided to cooperate or defect purely at random.

If we select a turn at random from a game of Prisoner's Dilemma only the four arrangements shown in figure 6 are possible.

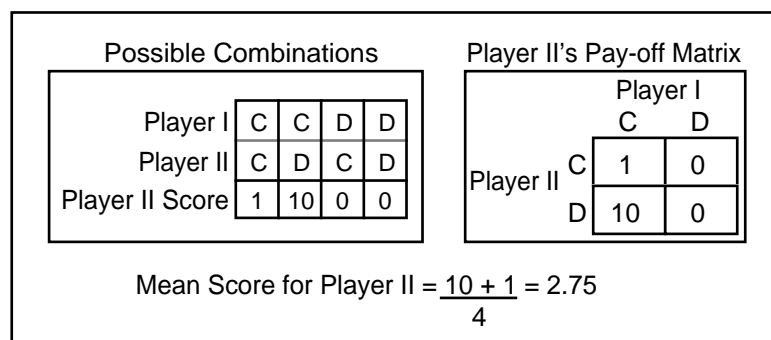


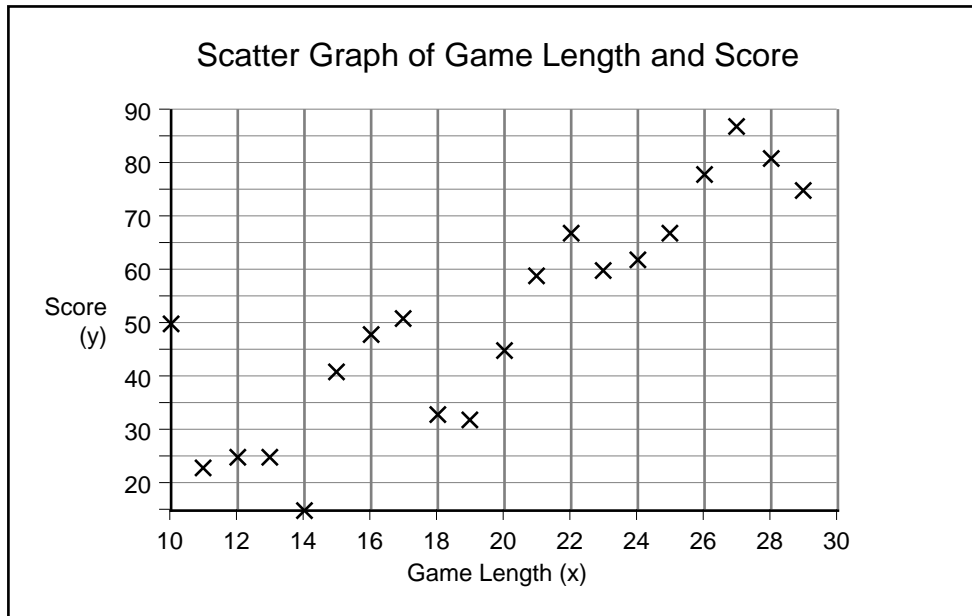
Fig. 6

This means that the mean score increase we would expect to see per move would be 2.75. Furthermore, we would expect the graph to pass through the origin, allowing us to predict a

linear regression equation of:

$$y = 2.75x$$

The first point to notice will be whether the graph appears to be linear. Graph 1 shows the data collected plotted on a scatter diagram. It does indeed appear to follow a linear relationship.



Graph 1

We will now go on to calculate the line of best fit for this data so we can compare it with the theoretical line of best fit.

If the gradient of the line of best fit is significantly less than the theoretical gradient, this would imply that human players perform poorly against tit-for-tat's algorithm. If, however, the gradient is greater, this would indicate that humans achieve (on average) higher than expected scores when playing against tit-for-tat. Based on our conclusions from the previous tests we would expect the latter case to occur with this data.

Before we calculate the line of best fit for the data in table 9 it is worth finding the Spearman's rank correlation coefficient (r_s) and testing whether a straight line graph is appropriate. Since the theory and the graph both indicate positive correlation we will test to see whether r_s is significantly close to 1. It is not appropriate to perform a test on the product moment correlation coefficient (r) since the game length is certainly not normally distributed.

This test indeed demonstrated positive correlation at the 5% level (see table 10 and the associated hypothesis test). We can now go on to calculate a line of best fit for this data. The working for this is summarised in table 11 with the final regression line plotted on graph 2.

Table 10

Game Length in Rank Order	Score in Rank Order	d	d ²
1	10	9	81
2	2	0	0
3	3.5	0.5	0.25
4	3.5	0.5	0.25
5	1	4	16
6	7	1	1
7	9	2	4
8	11	3	9
9	6	3	9
10	5	5	25
11	8	3	9
12	12	0	0
13	15	2	4
14	13	1	1
15	14	1	1
16	16	0	0
17	18	1	1
18	20	2	4
19	19	0	0
20	17	3	9
			174.5

$H_0: r_s = 0$

$H_1: r_s > 0$

One tailed test at the 5% level
Critical value is 0.3805

$$r_s = 1 - \frac{6 \sum d^2}{n(n^2 - 1)} = 0.8688 \text{ (to 4 d.p.)}$$

Since $0.8688 > 0.3805$ reject H_0 and conclude that there is enough evidence at the 5% level to say that the data is positively correlated.

Table 11

Game Length (x)	Score (y)	$y = \alpha + \beta x$	ϵ_i
10	50	21.943	28.057
11	23	25.023	-2.023
12	25	28.102	-3.102
13	25	31.182	-6.182
14	15	34.262	-19.262
15	41	37.341	3.659
16	48	40.421	7.579
17	51	43.501	7.499
18	33	46.580	-13.580
19	32	49.660	-17.660
20	45	52.740	-7.740
21	59	55.820	3.180
22	67	58.899	8.101
23	60	61.979	-1.979
24	62	65.059	-3.059
25	67	68.138	-1.138
26	78	71.218	6.782
27	87	74.298	12.702
28	81	77.377	3.623
29	75	80.457	-5.457
			-0.000

$n = 20$

$\sum x = 390$

$\sum y = 1024$

$\sum xy = 22016$

$\sum x^2 = 8270$

$\sum y^2 = 60970$

$$S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} = 665$$

$$S_{yy} = \sum y^2 - \frac{(\sum y)^2}{n} = 8541.2$$

$$S_{xy} = \sum xy - \frac{(\sum x)(\sum y)}{n} = 2048$$

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}} = 3.0797 \text{ (to 4 d.p.)}$$

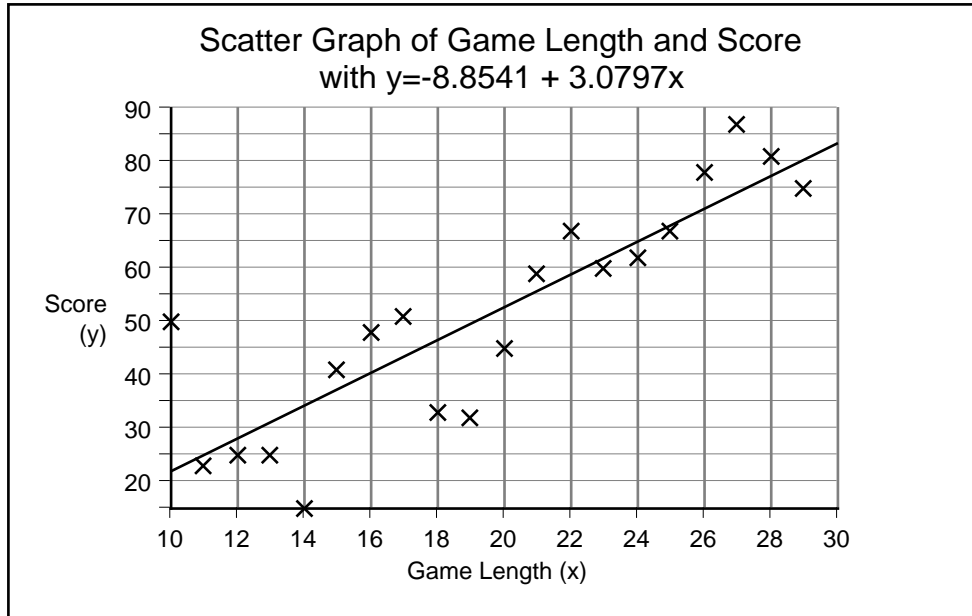
$$\hat{\alpha} = \bar{y} - \bar{x}\hat{\beta} = -8.8541 \text{ (to 4 d.p.)}$$

$y_i = \alpha + \beta x_i + \epsilon_i$

$\therefore y_i = -8.8541 + 3.0797x_i + \epsilon_i$

$$r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = 0.859 \text{ (to 3 d.p.)}$$

$$\hat{\sigma}^2 = \frac{S_{yy} (1 - r^2)}{n - 2} = 124.11 \text{ (to 2 d.p.)}$$



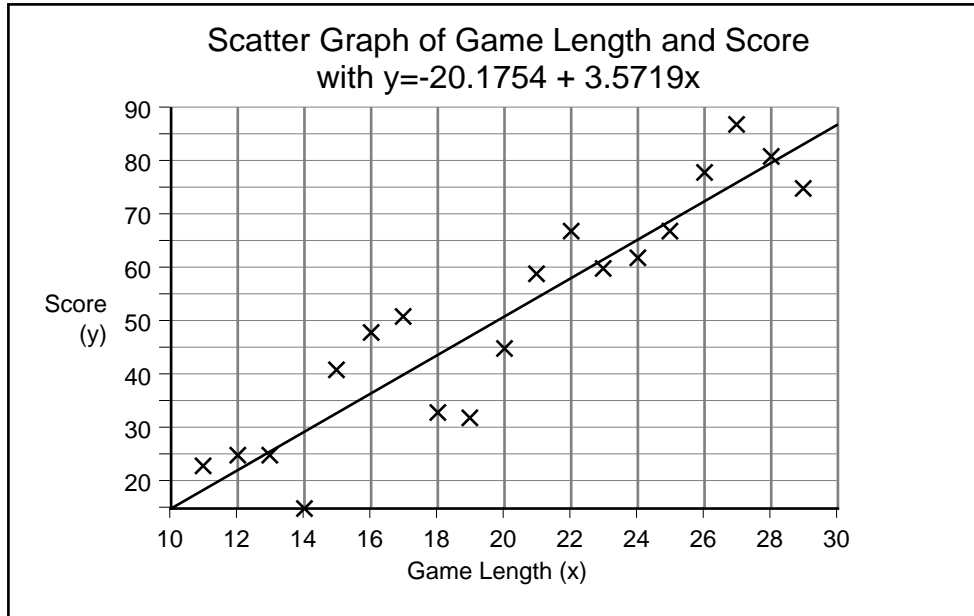
Graph 2

The sum of the residuals is exceptionally close to zero, suggesting that this regression model is appropriate. One residual value, however, stands out clearly – 28.057 for a game length of 10. By excluding this point it should be possible to calculate a more accurate line of best fit.

Table 12

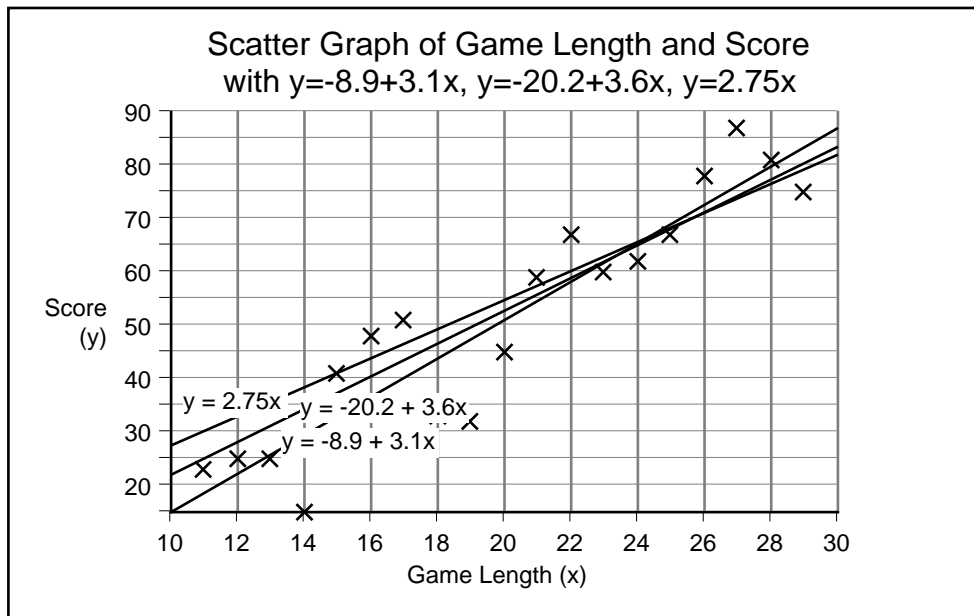
Game Length (x)	Score (y)	$y = \alpha + \beta x$	ϵ_i
11	23	19.116	3.884
12	25	22.688	2.312
13	25	26.260	-1.260
14	15	29.832	-14.832
15	41	33.404	7.596
16	48	36.975	11.025
17	51	40.547	10.453
18	33	44.119	-11.119
19	32	47.691	-15.691
20	45	51.263	-6.263
21	59	54.835	4.165
22	67	58.407	8.593
23	60	61.979	-1.979
24	62	65.551	-3.551
25	67	69.123	-2.123
26	78	72.695	5.305
27	87	76.267	10.733
28	81	79.839	1.161
29	75	83.411	-8.411
			-0.000

$$\begin{aligned}
 n &= 19 \\
 \sum x &= 380 \\
 \sum y &= 974 \\
 \sum xy &= 21516 \\
 \sum x^2 &= 8170 \\
 \sum y^2 &= 58470 \\
 \hline
 S_{xx} &= \sum x^2 - \frac{(\sum x)^2}{n} = 570 \\
 S_{yy} &= \sum y^2 - \frac{(\sum y)^2}{n} = 8539.6841 \quad (\text{to 4 d.p.}) \\
 S_{xy} &= \sum xy - \frac{(\sum x)(\sum y)}{n} = 2036 \\
 \hline
 \hat{\beta} &= \frac{S_{xy}}{S_{xx}} = 3.5719 \quad (\text{to 4 d.p.}) \\
 \hat{\alpha} &= \bar{y} - \bar{x}\hat{\beta} = -20.1754 \quad (\text{to 4 d.p.}) \\
 y_i &= \alpha + \beta x_i + \epsilon_i \\
 \therefore y_i &= -20.1754 + 3.5719x_i + \epsilon_i \\
 \hline
 r &= \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = 0.923 \quad (\text{to 3 d.p.}) \\
 \hat{\sigma}^2 &= \frac{S_{yy} (1 - r^2)}{n - 2} = 74.54 \quad (\text{to 2 d.p.})
 \end{aligned}$$



Graph 3

The two calculated lines and the theoretical line have been plotted together on graph 4 below.



Graph 4

From this we can see that all three equations are very close estimates for the data. The question remains, is the gradient significantly different from the theoretical gradient of 2.75? If not, we can assume that human players perform no better or worse against tit-for-tat than a random opponent would. I will use the more refined model for the test.

$$H_0: \beta = 2.75$$

$$H_1: \beta \neq 2.75$$

Two tailed test the the 5% level

Critical values obtained from t_{17} are ± 2.110

$$t_{\text{test}} = \frac{3.5719 - 2.75}{\sqrt{\frac{74.54}{19}}} = 0.415 \text{ (to 3 d.p.)}$$

Since $0.415 < 2.110$ accept H_0 and conclude that there is enough evidence at the 5% level to say that $\beta = 2.75$.

We can conclude that the line of best fit fits the theoretical line of fit.

The theory predicted that the graph would pass through the origin, however neither of the linear regression equations does. It is worth checking to see whether the origin lies within 95% confidence intervals for 'a' (the point of y-axis intersection). Once again I will make use of the more refined model.

$$a \sim N\left(\alpha, \frac{\sigma^2 \sum x^2}{n S_{xx}}\right)$$

$$\hat{\alpha} = -20.1754 \text{ (to 4 d.p.)}$$

$$\hat{\sigma}^2 = 74.54 \text{ (to 2 d.p.)}$$

$$\sum x^2 = 8170$$

$$n = 19$$

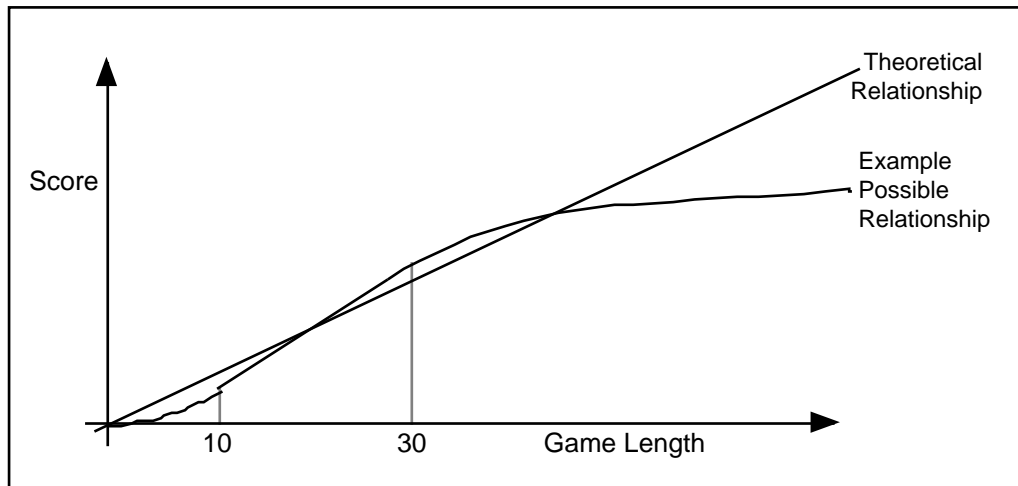
$$S_{xx} = 570$$

$$\therefore a \sim N(-20.1754, 56.232)$$

95% confidence intervals for 'a' are:

$$-20.1754 \pm 2.110 \sqrt{56.232} = (-36.00, -4.35)$$

It is interesting to note that the origin (0,0) does not lie within the 95% confidence intervals calculated for 'a' for the more refined linear regression model. This suggests to me that the linear model for total human player score may not hold for very short games (ie. game lengths less than 10) and so may well not hold for game lengths greater than 30 either. If this is the case we may well find that human players under perform when playing against tit-for-tat when the number of moves is as high as in Axelrod's tournaments. The true equation linking final human score with game length may well just be approximately linear between the game lengths I have chosen to investigate. Graph 5 illustrates one possible relationship between game length and score.



Graph 5

Conclusion

This project has certainly led to some interesting conclusions.

My first aim in this project was to see whether tit-for-tat performed as exceptionally as it did in the original Axelrod tournament when playing against human opponents rather than computer opponents and over shorter games than Axelrod used. The stark conclusion is that it does not.

Not only did tit-for-tat constantly underperform when compared to the other algorithms, but there was no significant difference between the mean of a purely random player and that of tit-for-tat. A test on the variance, which could have restored some of the reputation of tit-for-tat by concluding that it was more consistent than the random player, concluded that, once again, there was no significant difference between the algorithms. We can conclude from this that tit-for-tat performs poorly against humans over short game lengths, but must somehow improve over longer (200) length games (as the evidence from the Axelrod tournament suggests).

The second aim was to investigate the correlation between game length and total score. Is the relationship linear and, if so, can the final score be predicted if the length of the game is known?

A theoretical correlation was calculated and the experimental data, on the whole, supports this prediction. It would appear that for every additional move played the score increases, on average, by 2.75. An interesting conclusion from the linear regression is that there is no significant difference between the gradient of the theoretical line and the line of best fit calculated from regression. This means that human players do not perform differently when playing against tit-for-tat than a random player for game lengths between 10 and 29 moves.

This project has yielded many surprising results. It appears that conclusions drawn from games of length 200 moves against computer opponents can not be extended to cover shorter games against human opponents. In fact quite the opposite conclusion emerges and we must say that tit-for-tat is, at best, a rather poor algorithm.

More Theory Behind the Tit-for-Tat Opponent

Mathematically, how can we describe the scores obtained when playing against a tit-for-tat opponent? Is it possible to calculate the probability of a particular score occurring? The following work does certainly not amount to a mathematical proof, but it does provide some answers to these questions. Because of the lack of rigour in the following arguments I have refrained from including this work in my main project and have relegated it to the end.

It is clear that the possible scores the player (as opposed to tit-for-tat) can obtain on any move other than the first move depends on the previous move the player made. This is because tit-for-tat merely copies the previous move the player made.

let X = "tit-for-tat's opponent's score after 'n' moves" where n is a constant

I originally approached this from a purely abstract direction. If we consider the sequence of moves we can say that the possible score on any move other than the initial move (which is purely random) depends on the player's previous move. Figure 7 shows the "transition rules".

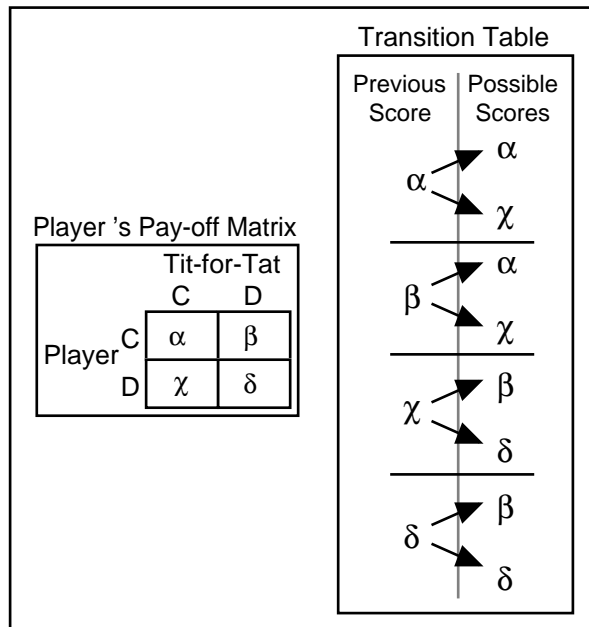


Fig. 7

From this I attempted to abstract a typical game of length 'n' moves. It is apparent that only certain arrangements of scores (ie. moves) are allowed. The typical payoff matrix for a Prisoner's Dilemma game scores 0 for ' β ' and ' δ ', so I can express any score (x) as:

$$x = r\chi + c\alpha$$

where r = number of times χ was scored

c = number of times α was scored

The moves for a game seem to break down into two basic building blocks:

$$(\beta, \alpha, \dots, \alpha, \chi) \text{ and } (\delta, \dots, \delta)$$

These blocks can finish abruptly at the end of a game (ie. after 'n' moves) and a game can start anywhere in the middle of a block (though obviously this would make little difference to the second block). The first block can contain any number of α scores, including none [ie. the set would become (β, χ)] and the second block can contain one or more δ scores.

The task of modelling a purely abstract game is immense and beyond my abilities at the moment. Instead I chose to focus on modelling a game where $n = 10$ and hope that a pattern presents itself which can be abstracted for any game.

I had the computer play all the possible games of length 10 (2^{11} possibilities) and charted the results in table 13.

		c										
		0	1	2	3	4	5	6	7	8	9	10
0	x	0	α	2α	3α	4α	5α	6α	7α	8α	9α	10α
	f	2	1	1	1	1	1	1	1	1	1	1
1	x	χ	$\chi+\alpha$	$\chi+2\alpha$	$\chi+3\alpha$	$\chi+4\alpha$	$\chi+5\alpha$	$\chi+6\alpha$	$\chi+7\alpha$	$\chi+8\alpha$	$\chi+9\alpha$	
	f	19	25	29	31	31	29	25	19	11	1	
2	x	2χ	$2\chi+\alpha$	$2\chi+2\alpha$	$2\chi+3\alpha$	$2\chi+4\alpha$	$2\chi+5\alpha$	$2\chi+6\alpha$	$2\chi+7\alpha$			
	f	64	119	153	160	140	99	49	8			
3	x	3χ	$3\chi+\alpha$	$3\chi+2\alpha$	$3\chi+3\alpha$	$3\chi+4\alpha$	$3\chi+5\alpha$					
	f	91	185	220	180	95	21					
4	x	4χ	$4\chi+\alpha$	$4\chi+2\alpha$	$4\chi+3\alpha$							
	f	50	85	65	20							
5	x	5χ	$5\chi+\alpha$									
	f	7	5									

where f = frequency of score

Table 13

As you can see from the table, the frequency of the scores when $r = 0$ (ie. when no move scored χ) follow a very clear pattern. It is a simple matter (based on the building blocks described above) to see that the frequency when $c = 0$ and $r = 0$ will always be two:

$$\delta, \delta, \dots, \delta \quad \text{or} \quad \delta, \delta, \dots, \delta, \beta$$

And the frequency will always be one when $r = 0$ and $c \geq 1$:

$$\delta, \delta, \dots, \delta, \beta, \alpha, \alpha, \dots$$

since the above sequence can never stop repeating a score of α without scoring a χ , but since $r = 0$ this can never happen and thus the score on the 'n'th move must be a α . No other possibilities exist.

Armed with the frequencies for when $c = 0$ and $r \geq 1$ I proceeded to discover how these frequencies were obtained. By using the transition rules and the building blocks I managed to draw up table 14.

r	f	
1	19	$9C_1x^2+9C_0x^1$
2	64	$8C_2x^2+8C_1x^1$
3	91	$7C_3x^2+7C_2x^1$
4	50	$6C_4x^2+6C_3x^1$
5	7	$5C_5x^2+5C_4x^1$

where ${}^nC_r = \binom{n}{r} = \frac{n!}{r!(n-r)!}$

Table 14

The above table did make certain assumptions about the pattern which might evolve. For example, I chose to represent 9 by ${}^9C_0 \times 1$. The advantage in representing numbers in this form will soon become apparent.

I then proceeded to work out the combinations for the frequencies for when $c = 1$ and $r \geq 1$. These were exceedingly difficult to work out and became too complicated after $r = 2$. The combinations I worked out were added to a table (see table 15). I also managed to calculate a few other values with are indicated on the table.

		c										Table 15
		0	1	2	3	4	5	6	7	8	9	
r	1	9C0x1+ 9C1x2	8C0x1+ 8C1x3	7C0x1+ 7C1x4	6C0x1+ 6C1x5	5C0x1+ 5C1x6	4C0x1+ 4C1x7	3C0x1+ 3C1x8	2C0x1+ 2C1x9	1C0x1+ 1C1x10	0C0x1+ 0	
	2	8C1x1+ 8C2x2	7C1x2 + 7C2x5	6C1x3+ 6C2x9	5C1x4+ 5C2x14	4C1x5+ 4C2x20	3C1x6+ 3C2x27	2C1x7+ 2C2x35	1C1x8+ 0			
	3	7C2x1+ 7C3x2	6C2x3+ 6C3x7	5C2x6+ 5C3x16	4C2x10+ 4C3x30	3C2x15+ 3C3x50	2C2x21+ 0					
	4	6C3x1+ 6C4x2	5C3x4+ 5C4x9	4C3x10+ 4C4x25	3C3x20+ 0	N.B. only bold values were calculated						
	5	5C4x1+ 5C5x2	4C4x5+0									

The other values in table 15 were predicted from observing the pattern which seemed to be emerging and then checked to verify that they gave the required frequency. This is certainly not the preferred way of approaching this problem, but I did not have the required skill to approach this from a more rigorous perspective.

The values by which the combinations are multiplied follow a very definite pattern. They are the sum of the corresponding values in the row above and the adjacent column. For example, the combination in cell $c=4, r=3$ has the formula ${}^3C_2 \times 15 + {}^3C_3 \times 50$. The values 15 and 50 are gain from $5+10$ and $20+30$ respectively (from cell $c=4, r=2$ and cell $c=3, r=3$). Is there a more general way to calculate these values?

		c										Table 16									
		0	1	2	3	4	5	6	7	8	9										
r	1	1	2	1	3	1	4	1	5	1	6	1	7	1	8	1	9	1	10	1	-
	2	1	2	2	5	3	9	4	14	5	20	6	27	7	35	8	-				
	3	1	2	3	7	6	16	10	30	15	50	21	-								
	4	1	2	4	9	10	25	20	-												
	5	1	2	5	-																

It is obvious that the first set of values for each column are Pascal's Triangle and as such can be calculated using the formula:

$${}^{r+c-1}C_c \quad \text{where } r = \text{number of times } \chi \text{ was scored}$$

$$c = \text{number of times } \alpha \text{ was scored}$$

With a little more thought you can see that the second set of values in each column can be calculated as shown in figure 8. Essentially they are the sum of two offset Pascal's Triangles.

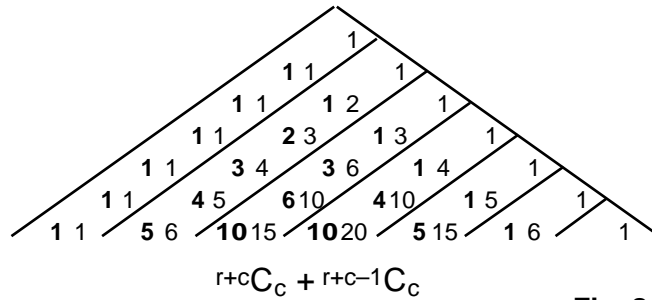


Fig. 8

Putting all this information together and assuming that the formulas can be abstracted we get the following formula:

let $X =$ "tit-for-tat's opponent's score after 'n' moves"
 where n is a constant

for a final score represented in the form $x = r\chi + c\alpha$
 where r = number of times χ was scored
 c = number of times α was scored

$$P(X=x) = \frac{1}{2^{n+1}} \left\{ \binom{r+c-1}{c} \left[\binom{n-r-c}{r-1} + \binom{n-r-c}{r} \right] + \binom{n-r-c}{r} \binom{r+c}{c} \right\}$$

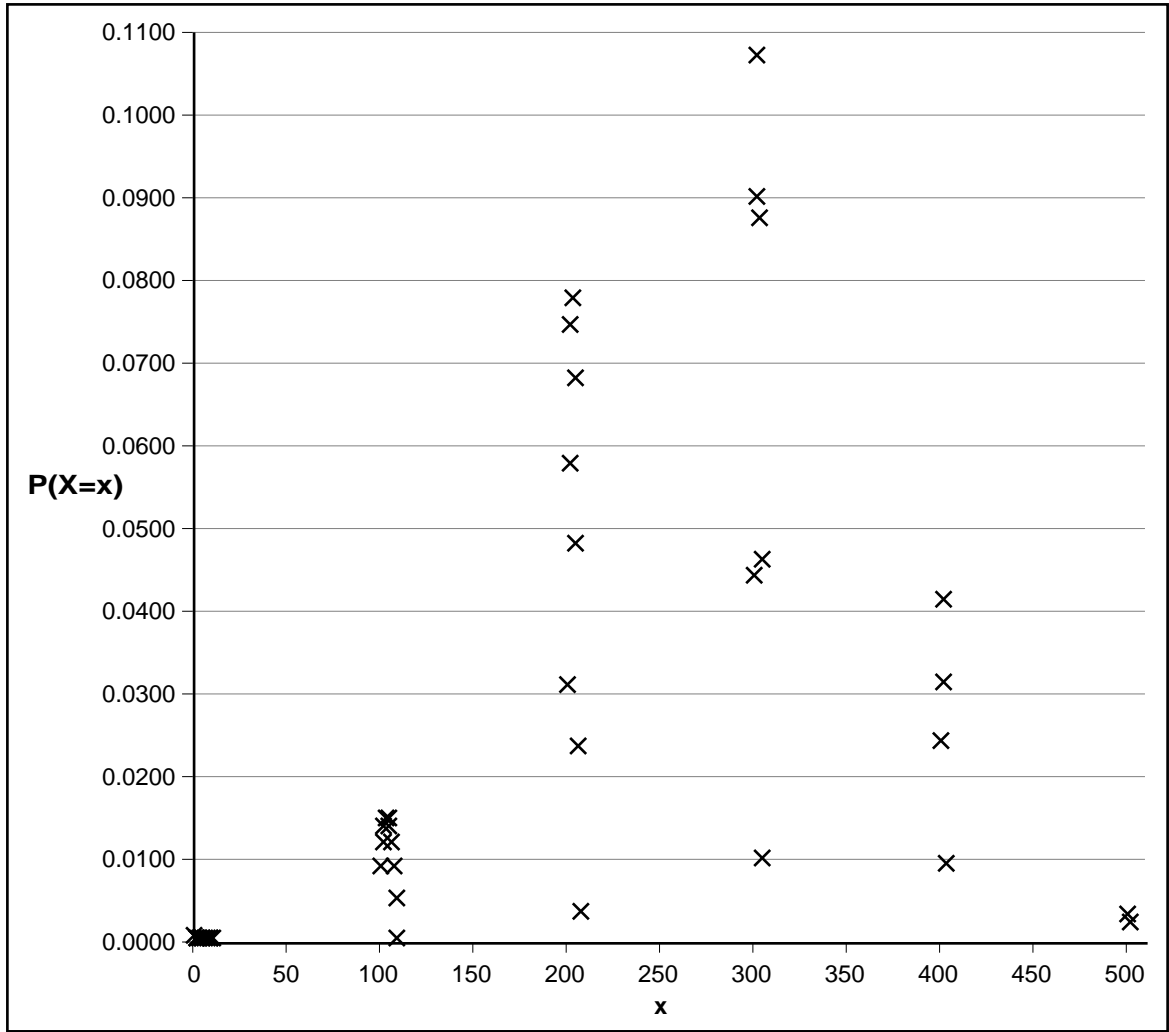
This formula relies on you being able to express any score in the form $x = r\chi + c\alpha$, but not any value of 'r' and 'c' is valid. In fact this is a simple matter to see that only certain combinations are possible. From looking at the building blocks described above you can see that for $r=0$, all values of 'c' up to and including 'n' are possible. After that the pattern is the same as we saw for in figure 12. For any value of 'r', the maximum possible value for 'c' is $n-2r+1$. Using this information and the formula above I drew up a table predicting the frequencies of scores for a game where $n=5$. This is shown in table 17. I then had the computer play all possible games when $n=5$. The predicted frequencies matches the true frequencies without exception indicating (but certainly not proving) that the formula is correct.

		c						
		0	1	2	3	4	5	
r	0	x	0	α	2α	3α	4α	5α
		f	2	1	1	1	1	1
	1	x	χ	$\chi+\alpha$	$\chi+2\alpha$	$\chi+3\alpha$	$\chi+4\alpha$	
		f	9	10	9	6	1	
	2	x	2χ	$2\chi+\alpha$	$2\chi+2\alpha$			
		f	9	9	3			
	3	x	3χ					
		f	1					

where f = frequency of score

Table 17

What does this distribution look like graphically? Because only certain discrete values are possible and the graph depends on the values of α and χ it depends on the situation. However, graph 5 shows the distribution when $n=10$, $\alpha=1$ and $\chi=100$.



Graph 6