

## Chapter 4

### Regression Analysis

In Chapter 3 we used some basic statistical methods - data transformations and one-way and two-way analysis of variance - for analysing the Olympic performances in swimming, running, jumping and throwing events over the period from 1928 to 2000. These methods identified trends and showed how the performances in these sports improved over the years.

In this chapter we develop statistical models based on linear regression analysis. We also summarise the results obtained from these models and thus develop an overall view that will facilitate comparison of the performances in the four sports. This approach also facilitates the detection of outstandingly good and bad performances.

In fact, the two-way analysis of variance method for comparison is the special case of multiple linear regression in which there are two sets of binary predictor variables corresponding to the determinant and the covariate. So it is instructive to repeat the analyses given in Chapter 3 using multiple regression. We do this for each sport in turn, and then combine the results.

#### 4.1 Swimming Performances

Figure 4.1 shows the results from fitting the multiple linear regression model to the men's swimming speeds in each of the 13 events for the 17 Olympics from 1928 to 2000. This table gives the regression coefficients for years and events, with 1928 as the baseline for year and the 50 meters freestyle as the baseline event. The table also gives the corresponding standard errors, from which confidence intervals may be computed, and p-values for testing the null hypotheses that the population coefficients are zero.

The bottom line of the table gives the r-squared and the adjusted r-squared (adjusted for the number of predictors used in the regression model), together with the residual sum of squares (rss), the residual standard deviation (s), and the p-value for testing the overall statistical significance of the model.

linear regression analysis: response = speed

predictor	coeff	St.Error	p-value
constant	1.8849	0.013819	0
year	( 0 )		0
1928	0.026254	0.012579	0.038932
1932	0.049603	0.012579	0.00013359
1936	0.056506	0.012579	1.5949e-005
1948	0.076478	0.013389	7.8156e-008
1952	0.124	0.012089	0
1956	0.15185	0.012089	0
1960	0.20992	0.01182	0
1964	0.2284	0.010807	0
1968	0.29366	0.010807	0
1972	0.3392	0.010925	0
1976	0.32986	0.010925	0
1980	0.35794	0.010807	0
1984	0.37103	0.010718	0
1988	0.3833	0.010718	0
1992	0.37995	0.010718	0
1996	0.40085	0.010718	0
2000			
event	( 0 )		0
50f	-0.21206	0.011424	0
100f	-0.42584	0.011486	0
100ba	-0.63388	0.012146	0
100br	-0.36369	0.012146	0
100bu	-0.39404	0.012146	0
200f	-0.56907	0.01199	0
200ba	-0.73635	0.011472	0
200br	-0.54321	0.011757	0
200bu	-0.60069	0.012607	0
200md	-0.50629	0.011424	0
400f	-0.69895	0.01199	0
400md	-0.60134	0.011424	0
1500f			

r-sq: 0.99302(0.99144) rss: 0.049054 df: 124 sd: 0.01989 p-value: 0

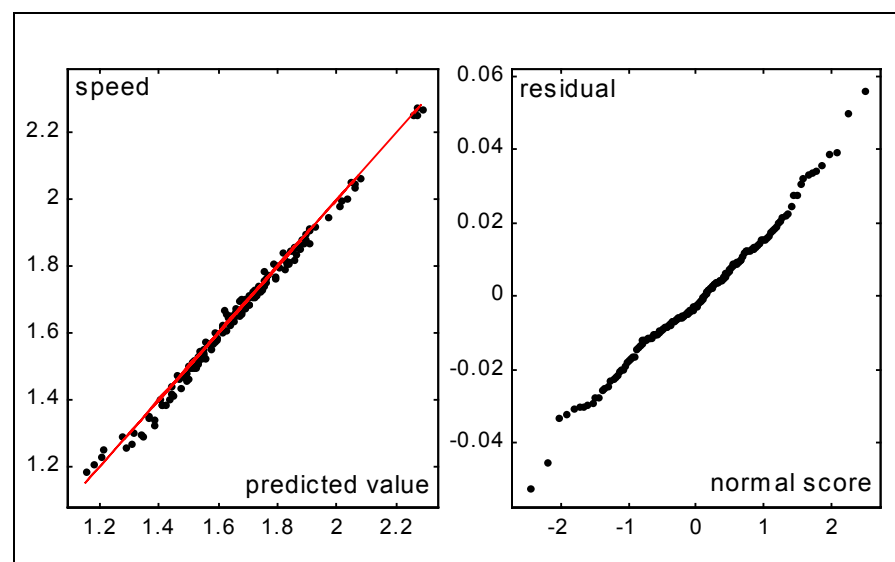


Figure 4.1: The results from fitting the multiple linear regression model to men swimming

The graphs below this table show (a) the relation between the speeds and their predicted values based on the model, and (b) the normal scores plot of residuals. These plots show (a) how well the model fits the data and (b) the plausibility of the normality assumption.

The coefficients for year increase steadily from 0 (the baseline value) in year 1928 to 0.40 in 2000. These coefficients show how the performances of the men swimmers improved over the whole period.

Since our objective is to compare the performances in the various sports (swimming, running, jumping and throwing), we need to make sure that these improvements are comparable.

How can we compare the performances in the different sports?

To answer this question, note that the performances in the swimming and running events are measured in terms of speeds, whereas the logarithms of the distances are used to compare the jumping and throwing performances.

A more valid way of comparing performances is based on percentage improvements. Taking logarithms is equivalent to using percentages, because an improvement of  $x$  percent from  $y$  to  $y(1+x/100)$  corresponds to an improvement from  $\log(y)$  to  $\log(y)+\log(1+x/100)$  in the logarithms, that is,  $\log(1+x/100)$ . Note that if natural logarithms are used, this is close to  $x/100$ , or  $x$  percent. (This is because the Taylor expansion of  $\ln(1+x/100)$  is  $(x/100)+\frac{1}{2}(x/100)^2+\dots$ ). If the logarithms are taken to some other base, it doesn't change the result, because the same constant multiplies everything.

In transforming the distances we used logarithms to base 10. So to provide a valid basis for comparing the four sports, we should repeat the analysis for the swimming and running events after transforming the speeds by taking base 10 logarithms. Doing this enables us to use percentage improvements to compare all performances.

Figure 4.2 shows the corresponding result for the men's swimming performances, after transforming the speeds using base 10 logarithms.

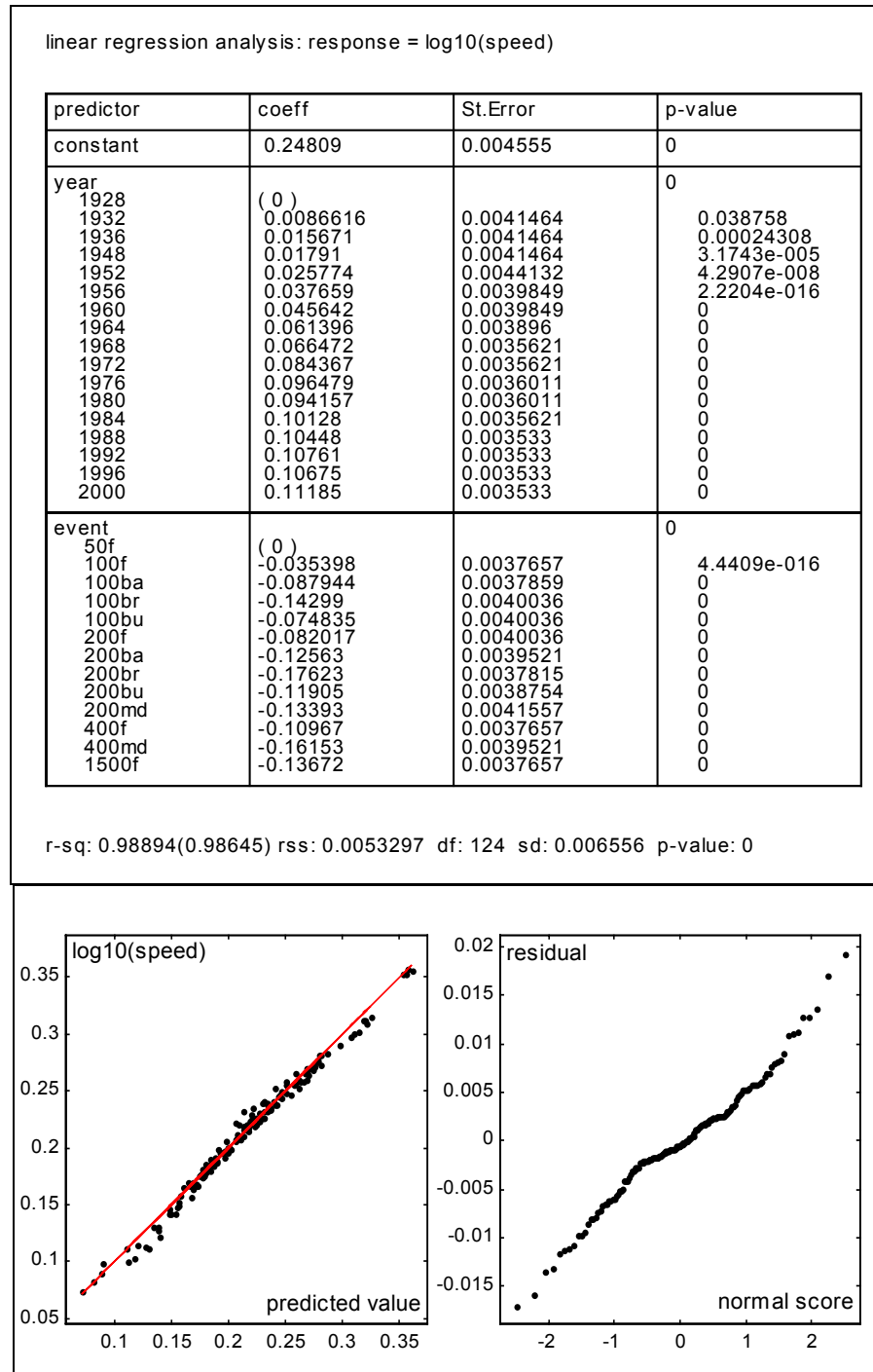


Figure 4.2: Men swimming using logarithm of speed to measure performance

The results show that, although the adjusted r-squared has decreased to 98.64%, the model still fits the data extremely well. The normal scores plot shows some slight curvature, but the normality assumption is still acceptable.

Figure 4.3 shows the corresponding result for the women's swimming performances.

linear regression analysis: response = log10(speed)

predictor	coeff	St.Error	p-value
constant	0.17514	0.0050049	0
year	( 0 )		0
1928	0.018356	0.0048257	0.00022808
1932	0.022796	0.0048257	6.4887e-006
1936	0.035266	0.0048257	3.6136e-011
1948	0.039985	0.0048257	2.2204e-013
1952	0.053553	0.0045992	0
1956	0.063823	0.0045992	0
1960	0.076603	0.0044447	0
1964	0.081536	0.0040394	0
1968	0.096085	0.0040394	0
1972	0.11605	0.0040766	0
1976	0.12144	0.0040766	0
1980	0.1214	0.0040394	0
1984	0.126	0.0040116	0
1988	0.12774	0.0040116	0
1992	0.12647	0.0040116	0
1996	0.13283	0.0040116	0
2000			
event	( 0 )		0
50f	-0.032624	0.003932	2.0917e-013
100f	-0.086428	0.003932	0
100ba	-0.13746	0.0041676	0
100br	-0.072044	0.0040394	0
100bu	-0.070815	0.0041676	0
200f	-0.11308	0.0041676	0
200ba	-0.16632	0.003932	0
200br	-0.10831	0.0041676	0
200bu	-0.12208	0.0043258	0
200md	-0.096865	0.003932	0
400f	-0.14382	0.0041155	0
400md	-0.10082	0.0041676	0
800f			

r-sq: 0.98775(0.98482) rss: 0.0054492 df: 117 sd: 0.0068245 p-value: 0

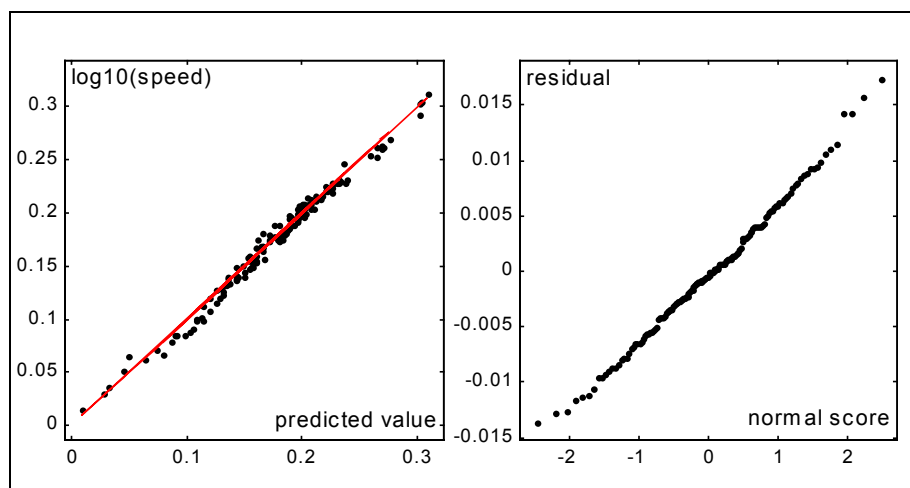


Figure 4.3: Women swimming using logarithm of speed to measure performance

In this case we see that the adjusted r-squared is 98.48%, again indicating that the model fits the data well. And again the plot of the speeds against their predicted values shows a close linear relation. The normal scores plot indicates that the normality assumption is reasonable.

## 4.2 Running Performances

Figure 4.4 shows the corresponding result for the men's running performances. Again the speeds are transformed using base 10 logs before fitting the model.

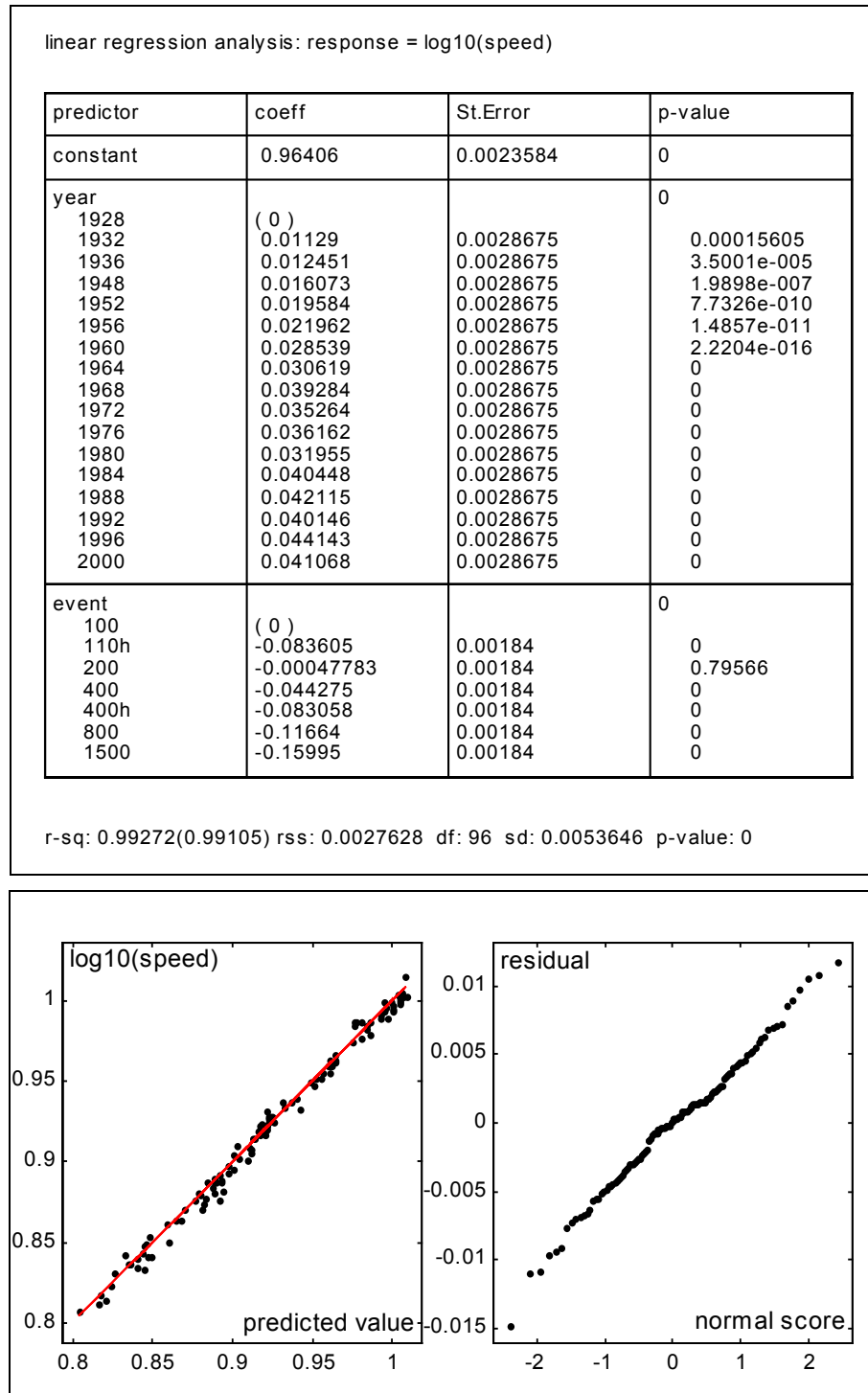


Figure 4.4: Men running using logarithm of speed to measure performance

In this case we see that the adjusted r-squared is 99.11%, indicating that the model fits the data very well. The plot of the speeds against their predicted values shows a close linear relation, and the normal scores plot indicates that the normality assumption is reasonable.

Figure 4.5 shows the corresponding result for the women's running performances.

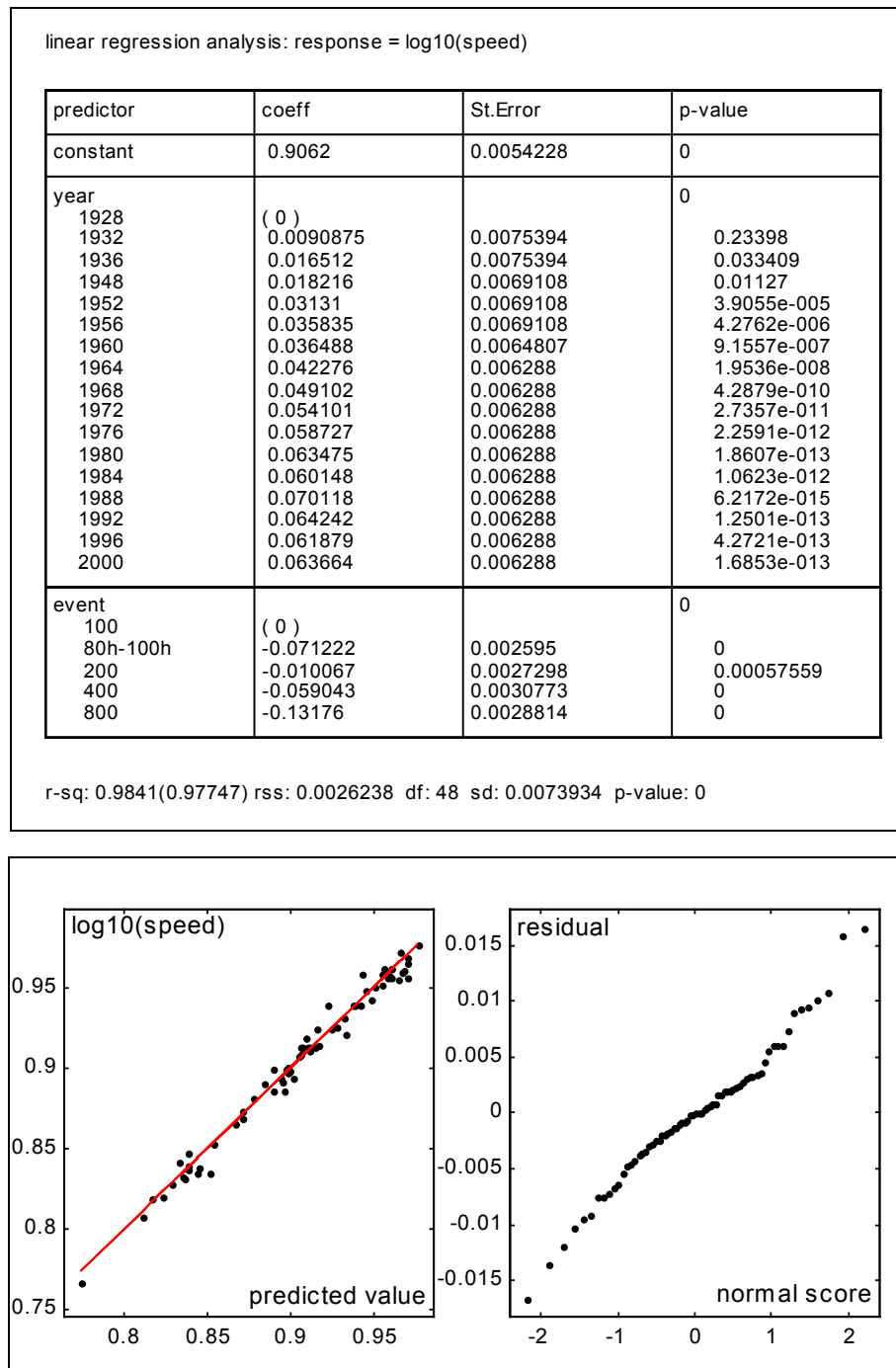


Figure 4.5: Women running using logarithm of speed to measure performance

The adjusted r-squared is now 97.75%, which is substantially less than the goodness-of-fit for the men's running events. This is not due to any particularly poorly fitting results, but rather to an overall larger variance in the speeds for the women's events. However, the plot of the log-transformed speeds against their predicted values does not show any curvature, and the normal scores plot indicates that the normality assumption is very reasonable.

### **4.3 Jumping and Throwing Performances**

Figure 4.6 shows the results after fitting the multiple regression model for the men's and women's jumping performances.

The adjusted r-squared is now 99.91%, showing that this model provides an almost perfect fit to the jumping performances. It should be noted, however, that the r-squared statistic is high because the events form three distinct clusters, as can be seen from the plot of the logs of the distances against their predicted values. In fact there is some evidence that one of the performances in the middle cluster (corresponding to the long jump) is unusually low. This low outlier shows up more clearly in the normal scores plot.

Figure 4.7 shows the results for the men's and women's throwing performances.

The adjusted r-squared is 99.28%, which again indicates that the model fits the data very well. But as for the jumping performances, the plot of the logs of the distances against their predicted values is clustered, this time into two groups. The lower group contains the men's and women's shot put events. The normal scores plot highlights four low outliers, one of which is extremely low.

The extremely low outlier corresponds to the performance by Tapio Korjus from Finland, who won the men's javelin event at the Seoul Olympics in 1988 with a throw of 84.28 meters. To see how poor this throw was, note that this distance had been exceeded 32 years earlier at the 1956 Melbourne Olympics by Egil Danielsen from Norway with a throw of 85.71 meters. It is interesting to note that the overall best performance for its time among all Olympic throwing events, according to our model,



occurred way back in 1932. This outstanding performance was achieved by Matti Jarvinen, also from Finland, who threw the javelin a distance of 72.71 meters.

linear regression analysis: response = log distance

predictor	coeff	St.Error	p-value
constant	0.29335	0.0061069	0
year	( 0 )		0
1928	0.0079971	0.0079624	0.31911
1932	0.015641	0.0079624	0.053976
1948	0.005144	0.0075837	0.50011
1952	0.016917	0.0075837	0.02933
1956	0.029962	0.0075837	0.00020203
1960	0.041759	0.0075837	7.4595e-007
1964	0.049707	0.0075837	1.2615e-008
1968	0.06034	0.0075837	4.7075e-011
1972	0.057194	0.0075837	2.4721e-010
1976	0.058499	0.0075837	1.242e-010
1980	0.070969	0.0075837	1.8163e-013
1984	0.071087	0.0075837	1.7097e-013
1988	0.081497	0.0075837	8.8818e-016
1992	0.078709	0.0075837	3.5527e-015
1996	0.079278	0.007443	1.3323e-015
2000	0.073567	0.007443	2.3315e-014
event	( 0 )		0
M.High jump	0.57439	0.0038623	0
M.Long jump	0.88619	0.0038623	0
M.Triple jump	-0.074486	0.0038623	0
W.High jump	0.47923	0.0041128	0
W.Long jump	0.81392	0.0090623	0
W.Triple jump			

r-sq: 0.99932(0.99909) rss: 0.0078615 df: 62 sd: 0.01126 p-value: 0

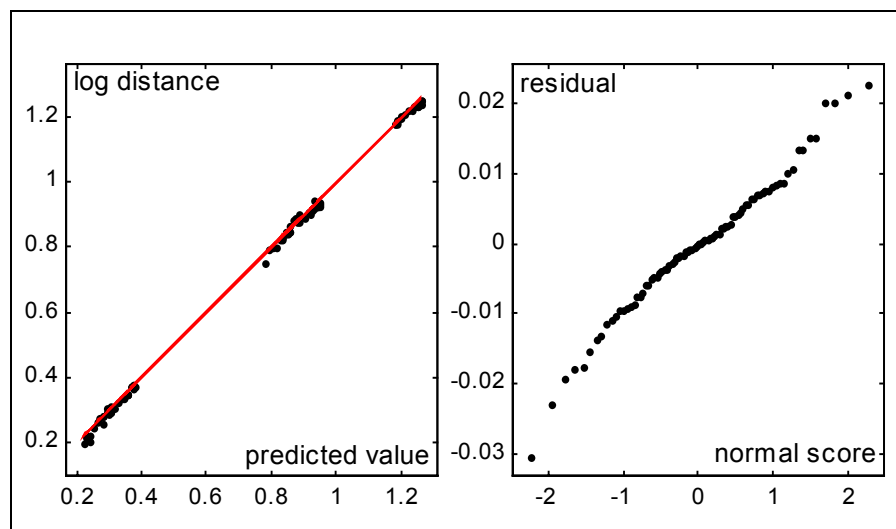


Figure 4.6: Jumping using logarithm of distance to measure performance

linear regression analysis: response = log distance

predictor	coeff	St.Error	p-value
constant	1.1737	0.010535	0
year	( 0 )		0
1928	0.013726	0.012871	0.28901
1932	0.032597	0.012871	0.013016
1936	0.023709	0.012479	0.060581
1948	0.060917	0.012479	4.4109e-006
1952	0.090379	0.012479	1.3118e-010
1956	0.10663	0.012479	2.6223e-013
1960	0.12168	0.012479	8.8818e-016
1964	0.1405	0.012479	0
1968	0.1599	0.012479	0
1972	0.17135	0.012479	0
1976	0.17891	0.012479	0
1980	0.16415	0.012479	0
1984	0.18873	0.012479	0
1988	0.17446	0.012479	0
1992	0.17394	0.012479	0
1996	0.1731	0.012479	0
2000			
event	( 0 )		0
M.shot put	0.49102	0.00727	0
M.discus	0.62751	0.00727	0
M.javelin	0.55094	0.00727	0
M.hammer	-0.030087	0.0077146	0.00018274
W.shot put	0.47001	0.00727	0
W.discus	0.47665	0.007405	0
W.javelin			

r-sq: 0.99418(0.99279) rrs: 0.041331 df: 92 sd: 0.021195 p-value: 0

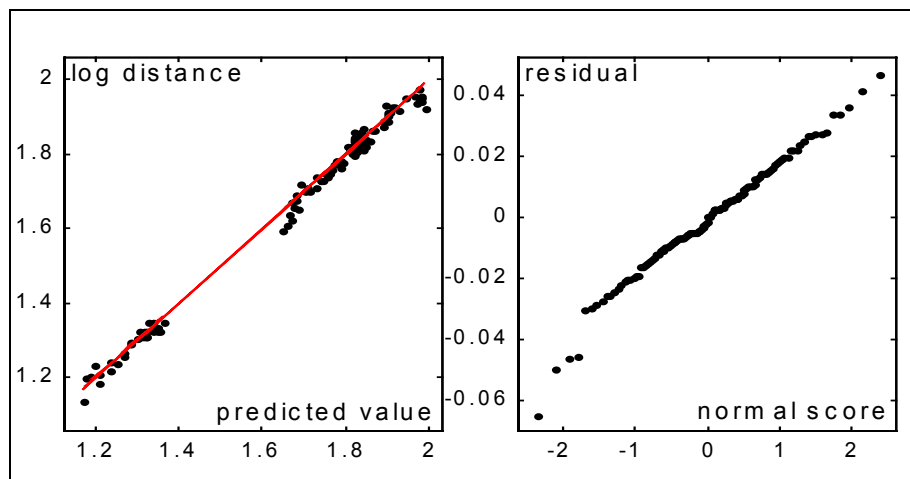


Figure 4.7: Throwing using logarithm of distance to measure performance

The other low outliers correspond to the winners of the women's discus events in 1928, 1932, and 1948. Given that these all occurred in earlier years, they provide some evidence that women's discus performances improved at a faster rate than other throwing events.

#### 4.4 Linear Regression Model

Following are models of each events.

- ◆ The model for men's swimming is

$$\log_{10}(\text{speed}) = 0.2481 + \text{year effect} + \text{event effect}$$

*year effect:*

$$[ \quad 0(^{\circ}28) \quad 0.0087(^{\circ}32) \quad 0.0158(^{\circ}36) \quad 0.0179(^{\circ}48) \quad 0.0258(^{\circ}52) \quad 0.0377(^{\circ}56) \\ 0.0456(^{\circ}60) \quad 0.0614(^{\circ}64) \quad 0.0665(^{\circ}68) \quad 0.0844(^{\circ}72) \quad 0.0965(^{\circ}76) \quad 0.0942(^{\circ}80) \\ 0.0103(^{\circ}84) \quad 0.1045(^{\circ}88) \quad 0.1076(^{\circ}92) \quad 0.1068(^{\circ}96) \quad 0.1119(2000) \quad ]$$

*event effect:*

$$[ \quad 0(50 \text{ free}) \quad -0.0354(100 \text{ free}) \quad -0.0879(100 \text{ back}) \quad -0.1430(100 \text{ breast}) \\ -0.0748(100 \text{ butterfly}) \quad -0.0820(200 \text{ free}) \quad -0.1256(200 \text{ back}) \quad -0.1762(200 \text{ breast}) \\ -0.1191(200 \text{ butterfly}) \quad -0.1339(200 \text{ medley}) \quad -0.1097(400 \text{ free}) \\ -0.1615(400 \text{ medley}) \quad -0.1367(1500 \text{ free}) \quad ]$$

- ◆ The model for women's swimming is

$$\log_{10}(\text{speed}) = 0.1751 + \text{year effect} + \text{event effect}$$

*year effect:*

$$[ \quad 0(^{\circ}28) \quad 0.0184(^{\circ}32) \quad 0.0228(^{\circ}36) \quad 0.0353(^{\circ}48) \quad 0.0400(^{\circ}52) \quad 0.0536(^{\circ}56) \\ 0.0638(^{\circ}60) \quad 0.0766(^{\circ}64) \quad 0.0815(^{\circ}68) \quad 0.0961(^{\circ}72) \quad 0.1161(^{\circ}76) \quad 0.1214(^{\circ}80) \\ 0.1214(^{\circ}84) \quad 0.1260(^{\circ}88) \quad 0.1277(^{\circ}92) \quad 0.1265(^{\circ}96) \quad 0.1328(2000) \quad ]$$

*event effect:*

$$[ \quad 0(50 \text{ free}) \quad -0.0326(100 \text{ free}) \quad -0.0864(100 \text{ back}) \quad -0.1375(100 \text{ breast}) \\ -0.0720(100 \text{ butterfly}) \quad -0.0708(200 \text{ free}) \quad -0.11310(200 \text{ back}) \quad -0.1663(200 \text{ breast}) \\ -0.1083(200 \text{ butterfly}) \quad -0.1221(200 \text{ medley}) \quad -0.0969(400 \text{ free}) \\ -0.1438(400 \text{ medley}) \quad -0.1008(800 \text{ free}) \quad ]$$

- ◆ The model for men's running is

$$\log_{10}(\text{speed}) = 0.9641 + \text{year effect} + \text{event effect}$$

*year effect:*

$$[ \quad 0(^{\circ}28) \quad 0.0113(^{\circ}32) \quad 0.0125(^{\circ}36) \quad 0.0161(^{\circ}48) \quad 0.0196(^{\circ}52) \quad 0.0220(^{\circ}56) \\ 0.0285(^{\circ}60) \quad 0.0306(^{\circ}64) \quad 0.0393(^{\circ}68) \quad 0.0353(^{\circ}72) \quad 0.0362(^{\circ}76) \quad 0.0320(^{\circ}80) \\ 0.0404(^{\circ}84) \quad 0.0421(^{\circ}88) \quad 0.0401(^{\circ}92) \quad 0.0441(^{\circ}96) \quad 0.0405(2000) \quad ]$$

*event effect:*

$$[ \quad 0(100 \text{ sprint}) \quad -0.0836(100 \text{ hurdles}) \quad -0.0007(200 \text{ sprint}) \quad -0.0443(400 \text{ sprint}) \\ -0.0831(400 \text{ hurdles}) \quad -0.1166(800 \text{ sprint}) \quad -0.156(1500 \text{ sprint}) \quad ]$$

- ◆ The model for women's running is

$$\log_{10}(\text{speed}) = 0.9062 + \text{year effect} + \text{event effect}$$

*year effect:*

$$\begin{bmatrix} 0(^{\circ}28) & 0.0091(^{\circ}32) & 0.0165(^{\circ}36) & 0.0182(^{\circ}48) & 0.0313(^{\circ}52) & 0.0358(^{\circ}56) \\ 0.0365(^{\circ}60) & 0.0423(^{\circ}64) & 0.0491(^{\circ}68) & 0.0541(^{\circ}72) & 0.0587(^{\circ}76) & 0.0635(^{\circ}80) \\ 0.0602(^{\circ}84) & 0.0701(^{\circ}88) & 0.0642(^{\circ}92) & 0.0619(^{\circ}96) & 0.0637(2000) & \end{bmatrix}$$

*event effect:*

$$\begin{bmatrix} 0(100 \text{ sprint}) & -0.0712(80 \text{ hurdles}, 100 \text{ hurdles}) & -0.0101(200 \text{ sprint}) \\ -0.0590(400 \text{ sprint}) & -0.1318(800 \text{ sprint}) & \end{bmatrix}$$

- ◆ The model for jumping is

$$\log_{10}(\text{distance}) = 0.2934 + \text{year effect} + \text{event effect}$$

*year effect:*

$$\begin{bmatrix} 0(^{\circ}28) & 0.0080(^{\circ}32) & 0.0156(^{\circ}36) & 0.0051(^{\circ}48) & 0.0169(^{\circ}52) & 0.0300(^{\circ}56) \\ 0.0418(^{\circ}60) & 0.0497(^{\circ}64) & 0.0603(^{\circ}68) & 0.0572(^{\circ}72) & 0.0585(^{\circ}76) & 0.0710(^{\circ}80) \\ 0.0711(^{\circ}84) & 0.0815(^{\circ}88) & 0.0787(^{\circ}92) & 0.0791(^{\circ}96) & 0.0736(2000) & \end{bmatrix}$$

*event effect:*

$$\begin{bmatrix} 0(\text{M.high jump}) & 0.5744(\text{M.long jump}) & 0.8862(\text{M.triple jump}) \\ -0.0745(\text{W.high jump}) & 0.4792(\text{W.long jump}) & 0.8139(\text{W.triple jump}) \end{bmatrix}$$

- ◆ The model for throwing is

$$\log_{10}(\text{distance}) = 1.1737 + \text{year effect} + \text{event effect}$$

*year effect:*

$$\begin{bmatrix} 0(^{\circ}28) & 0.0137(^{\circ}32) & 0.0326(^{\circ}36) & 0.0237(^{\circ}48) & 0.0609(^{\circ}52) & 0.0904(^{\circ}56) \\ 0.1066(^{\circ}60) & 0.1217(^{\circ}64) & 0.1405(^{\circ}68) & 0.1599(^{\circ}72) & 0.1714(^{\circ}76) & 0.1789(^{\circ}80) \\ 0.1642(^{\circ}84) & 0.1887(^{\circ}88) & 0.1745(^{\circ}92) & 0.1739(^{\circ}96) & 0.1731(2000) & \end{bmatrix}$$

*event effect:*

$$\begin{bmatrix} 0(\text{M.shot put}) & 0.4910(\text{M.discus}) & 0.6275(\text{M.javelin}) & 0.5509(\text{M.hammer}) \\ -0.0301(\text{W.shot put}) & 0.4700(\text{W.discus}) & 0.4767(\text{W.javelin}) & \end{bmatrix}$$

In the next section we combine the results from these models.

#### 4.5 Overall trends

Figure 4.8 graphs the coefficients for year from the models for the six sports (separating men's and women's swimming and running events but combining the men's and women's jumping and throwing events) shown in Figures 4.2 to 4.7.

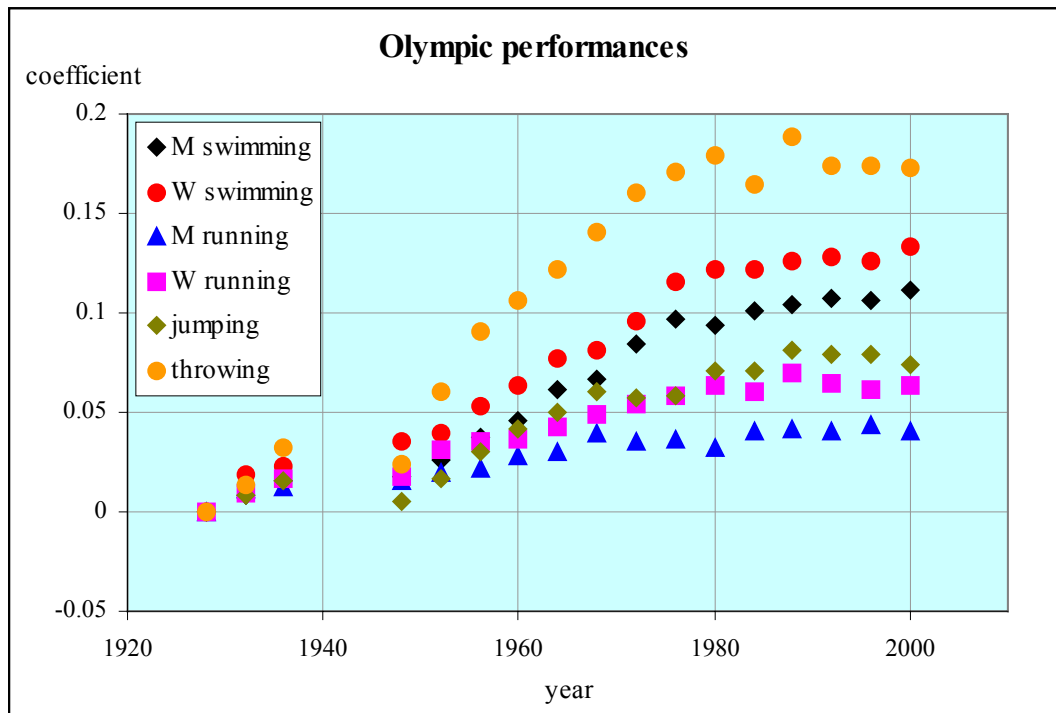


Figure 4.8: Comparison of performances in different sports

You can see from this graph that the performances increase with time, but more so in some events than others. The largest increase is in the throwing events, followed by the women's and the men's swimming events. The improvements in jumping and in the women's running over the years are similar. However the least improvements occur in the men's running events.

For each sport there is a steady improvement from 1928 to 1972, followed by a period in which there is little or no improvement. The main exception is in swimming, where both men's and women's performances improved noticeably at the Sydney 2000 Olympics.