

Evaluation of the Nietvoorbij Wine Score Card and Experimental Wine Panelists Utilizing Pattern Recognition Techniques

P. C. VAN ROOYEN, Oenological and Viticultural Research Institute (O V R I), Private Bag X5026, Stellenbosch 7600, South Africa.

The author wishes to acknowledge the valuable contribution of Mr. J. J. Theart in supervising the winemaking and organising the tasting sessions, as well as the faithful attendance and diligent efforts of the 18 panelists who participated in the 1979 wine evaluation programme of the O V R I.

Accepted for publication: April 1982

Principal component analysis (PCA) was applied to four cultivar wines, submitted repetitively amongst more than 450 experimental wines, to a sensory evaluation panel of 18 members who used a scoring system comprising overall wine quality and 11 wine descriptors. Inconsistent judges could be eliminated by the evaluation of scatter diagrams. After eliminating 11 judges, the scores of the remaining 7 were used to evaluate the score card in terms of weights placed on individual descriptors in a multiple regression equation, which related the 11 parameters to overall quality scores. Deviations from actual score card weights are discussed in terms of previous PCA analyses, and it is argued that both cultivar and the composition of a mixed data set with respect to these factors, could affect the relative importance of certain parameters. Fitting a similar equation to a large data set consisting of about 480 wines, comprising 23 different cultivars, confirmed the need for further investigations concerning relative score card weights, as well as a critical evaluation of score card parameters for evaluating widely diverging experimental wines.

Sensory evaluation of wines is central to research in oenology and viticulture. Many researchers base inferences from experimental results on quality scores obtained from a taste panel. Various score cards are used internationally, and have doubtless been subjected to evaluations at various times. The Nietvoorbij experimental score card, in the form reported by Tromp & Conradie (1979), has been in use at O V R I for many years. Wine evaluation at the O V R I has to a large extent been based on the A side of the card, where overall quality is evaluated. The B side, where 11 descriptors related to wine quality are scored, has been used mainly as an informative aid to determine possible reasons for overall score deviations. The elimination of inconsistent panelists before calculating final scores has also been done in the past by using routine statistical methods based on overall quality scores. A long-term policy of panel selection and adequate training is, however, not served adequately by this low resolution selection procedure, and does not supply information which pinpoints specific shortcomings regarding both individual judges and the scoring system. Wine descriptors on the B side of the score card were not used for this purpose because of limited time, the complexity and general unavailability of good statistical methods, as well as the general notion that the latter data have no bearing on the overall score of the wine (Tromp & Conradie, 1979). However, recent studies have shown the usefulness in this respect of principal component analysis (PCA) in conjunction with linear multiple regression analysis (Wu, Bargmann & Powers, 1977; Kwan & Kowalski, 1980 a), and provided guidelines for the effectivity of a score card using wine descriptors as basis. The most promising aspect of the recent work by Kwan & Kowalski (1980 a, 1980 b) was the use of these methods to sift panelists and, most important, to relate the chemical composition of wines to their sensory scores.

The aim of this study was to evaluate the 1979 Nietvoorbij panelists, as well as the current wine score card,

by using the methods outlined by Kwan & Kowalski (1980 a) and applying them to reference wines of three different cultivars originating from the Nietvoorbij Experimental farm, and made in the experimental cellar under controlled conditions. A final analysis was also made using the 1979 sensory evaluations of a comprehensive range of more than 480 experimental wines, comprising 23 different cultivars.

MATERIALS AND METHODS

Sensory evaluation: During regular annual sensory evaluations conducted as described by Tromp & Conradie (1979), pure cultivar wines of Cabernet Sauvignon, Colombard and Chenin blanc were submitted to a panel of 18 tasters. All wines were fermented to dryness, and each wine was submitted in its own class, i.e. the two whites amongst the other dry white wines. The wines were submitted at least four times during the tasting period of several weeks, and were tasted on a random basis amongst the approximately 800 wines that were evaluated in the 1979 season. The first two wines were of the 1979 vintage, while a Cabernet Sauvignon and Chenin blanc of the 1978 vintage were also used in order to evaluate a possible effect of bottle ageing for one year. The score card based on that reported by Tromp & Conradie (1979) was modified slightly and adapted for direct keypunching (Fig. 1).

Attendance at the tasting sessions was good, so that the numbers of original data vectors for the four reference wines were: Cabernet Sauvignon 1979 (70), Colombard 1979 (137), Cabernet Sauvignon 1978 (80) and Chenin blanc 1978 (62). Most panelists had more than 5 years' experience as experimental wine panelists, as well as 10 to 30 years' experience as general wine tasters. With one exception all were males.

SECTION A

Unacceptable	0	1
Average quality	—	2
with faults	0	3
	—	4.
Average quality	0	5
	+	6
Above average quality with some	0	7
outstanding characteristics	+	8
Superior	0	9

SECTION B

		Superior Outstanding	Very Good Above average	Good Average	Acceptable	Unacceptable	Motivation
Eye	Clarity ²⁶	5	4	3	2	1	
	Colour ²⁷	5	4	3	2	1	
Nose	Typicality ²⁸	5	4	3	2	1	
	Maturation bouquet ²⁹	5	4	3	2	1	
	Purity ³⁰			3	2	1	
Mouth	Acidity ³¹			3	4 + 2 -	5 + 1 -	
	Astringency ³²			3	4 + 2 -	5 + 1 -	
	Bitterness ³³			3	2	1	
Overall impression	Fullness ³⁴	5	4	3	2	1	
	Flavour ³⁵	5	4	3	2	1	
Overall impression	Harmony ³⁶	5	4	3	2	1	

FIG. 1
The Nietvoorbij experimental score card.

Data processing: The 12 sensory evaluation scores (11 descriptors + overall quality) obtained at each tasting from each panelist were regarded as features or data vectors, each giving the total evaluation of the wine by a particular panelist in terms of the 12 properties. The data were analysed using a batch-process version of the pattern recognition system "ARTHUR" (Harper *et al* 1977), executed on a UNIVAC 110 computer of the

University of Stellenbosch. The principal component factor analysis subprogram KAPRIN/KATRAN was used but was preceded by AUTOSCALE to normalise the data, and CORREL to examine feature/feature and feature/property covariances and correlations where applicable. VARVAR was used to generate line printer plots of three factors against each other in order to identify panelists deviating from the cluster of uniform

sensory evaluations of their colleagues. After the elimination of inconsistent panelists, the programme LEAST (least square multi-linear regression analysis) was used to determine the weights these judges placed on the individual sensory parameters (Kwan & Kowalski, 1980 a).

After treating each of the four wines separately, LEAST was again used, employing as input the combined, selected data of the three cultivars to establish parameter weights when the effect of cultivar was disregarded. A similar run was also made using the same information for more than 400 wines, spanning a range of 23 cultivars of the varieties Colombard, Chenin blanc, Cape Riesling, Bukettraube, Sauvignon blanc, Kerner, Chardonnay, Weisser Riesling, Fernao Pires, Vital, Muscat d'Alexandrie, Sémillon, Palomino, Clairette blanche, Raisin blanc, Ugni blanc, Pinotage, Cinsaut, Cabernet Sauvignon, Tinta Barocca, Chenel and two local *Vitis vinifera* crosses.

RESULTS AND DISCUSSION

Consistency amongst judges: Cabernet Sauvignon, 1979 vintage: Principal component analysis, extracting factors which the judges used to evaluate this wine, re-

vealed that three factors, explaining 68,6 per cent of the variation, could be used to evaluate deviations from uniformity for the individual panelists. The factor loadings for the first three principal components are given in Table 1. The first principal factor can be regarded as the single best summary of individual attributes exhibited in the data (Kwan & Kowalski, 1980 a), and has as major contributors overall quality, typicality, flavour and harmony. In the second factor colour, acidity and maturation play the major parts, and in factor three clarity, purity, bitterness and acidity are important. Two dimensional plots of the three factors against one another are given in Figures 2, 3 and 4. Numbers on the figures are unique codes for the 18 different judges. Missing code numbers are ones overprinted by the line printer, but are not considered important, as only outliers are evaluated. Judges, will, moreover, be penalised over all four wines, so that the chances of being an "underlier" outside the cluster for all the wines are very small. The darkened areas on the three figures were arbitrarily chosen as the cluster of uniform evaluation. Judges appearing more than once outside or near the edge of the cluster (Fig. 2: numbers 9, 11, 13, 17, 18, 23, 24; Fig. 3: numbers 9, 11, 13, 17, 18, 23, 24; Fig. 4:

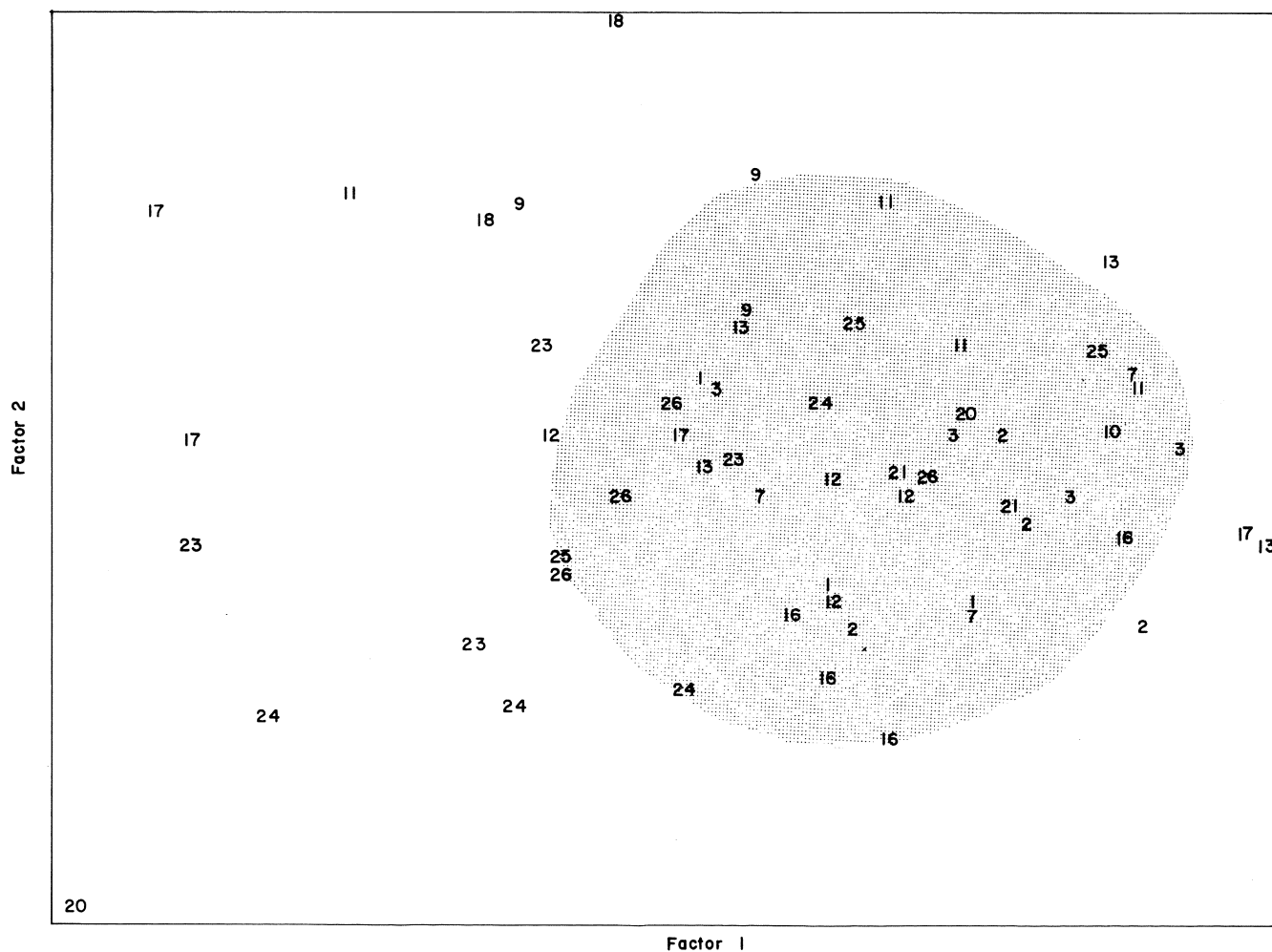
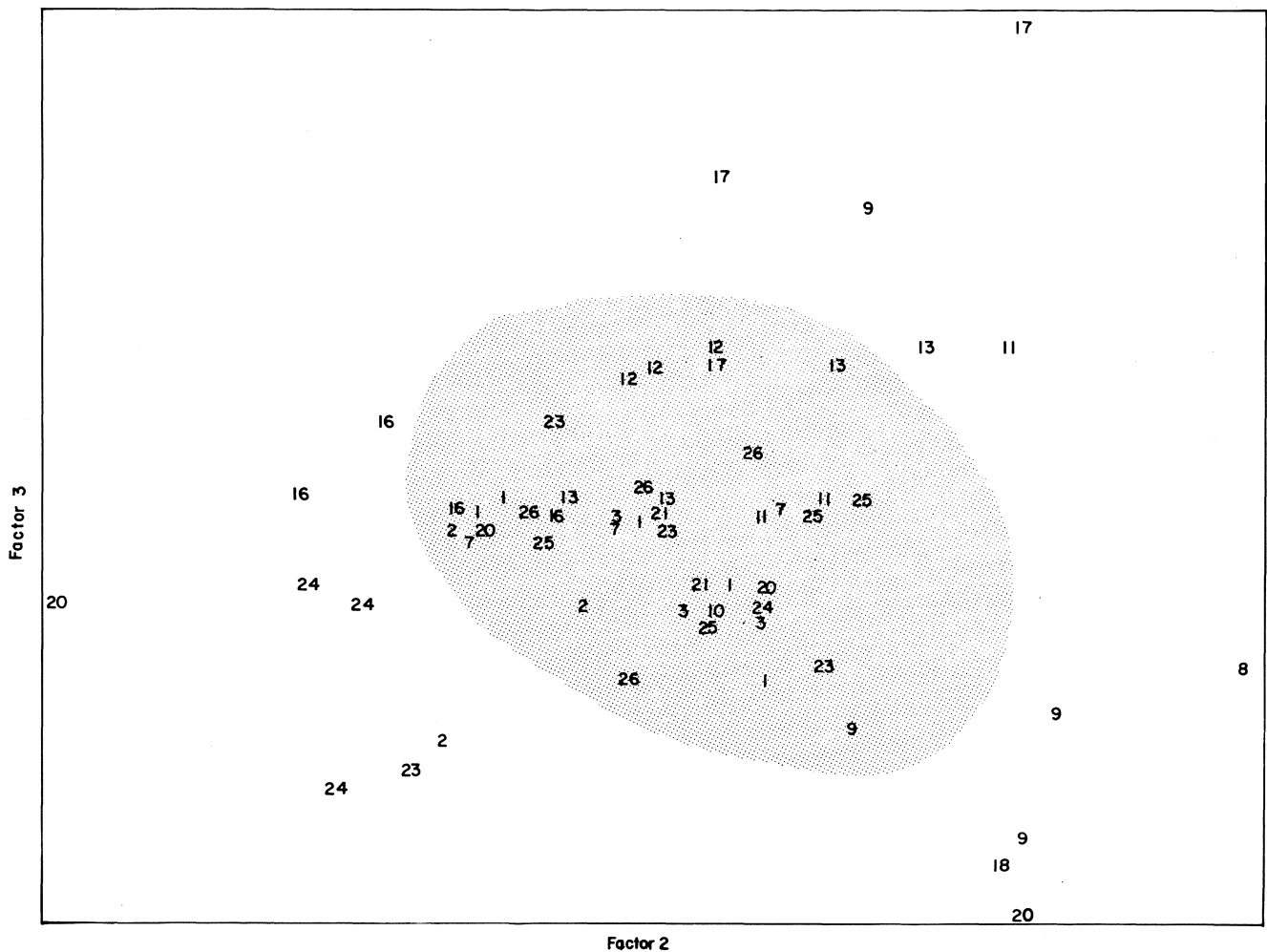


FIG. 2
Plot of first against second principal factor, 1979 Cabernet Sauvignon.



Factor 2

FIG. 4

Plot of second against third principal factor, 1979 Cabernet Sauvignon.

Cabernet Sauvignon, 1978 vintage: Application of the same procedure outlined above resulted in 67,7 per cent of the variation being explained by the first three principal components. The factor loadings are given in Table 3. The relative factor loadings for the first principal component are similar to those of the previous wine, but clarity, purity and acidity for factor two, and clarity, maturation and bitterness for factor three are important in this case.

Colombard, 1979 vintage: The first three principal components extracted from the Colombard data explain 67,6 per cent of the total variation. The factor loadings for this wine are given in Table 4. The first factor is mainly concerned with overall quality, flavour and harmony, whereas the second is primarily loaded with respect to clarity, colour and astringency. Maturation bouquet, acidity and astringency figure prominently in the third principal factor.

TABLE 3
Factor loadings for the first three principal components:
1978 Cabernet Sauvignon

Factor	I	II	III
Variation explained	38,8 per cent	17,4 per cent	11,4 per cent
Overall quality	-0,42	-0,02	-0,03
Clarity	-0,11	+0,41	+0,61
Colour	-0,23	+0,17	+0,05
Typicality	-0,37	-0,05	-0,14
Maturation bouquet	-0,25	-0,20	-0,41
Purity	-0,17	-0,45	-0,09
Acidity	-0,07	+0,53	-0,30
Astringency	+0,07	+0,35	-0,17
Bitterness	-0,22	-0,33	+0,51
Fullness	-0,35	+0,17	+0,09
Flavour	-0,42	+0,04	+0,10
Harmony	-0,41	+0,14	-0,17

TABLE 4
Factor loadings for the first three principal components:
1979 Colombard

Factor	I	II	III
Variation explained	40,3 per cent	15,3 per cent	12,0 per cent
Overall quality	+0,42	+0,22	+0,13
Clarity	+0,05	-0,48	+0,02
Colour	+0,14	-0,53	-0,31
Typicality	+0,38	+0,09	+0,04
Maturation bouquet	+0,10	-0,30	-0,44
Purity	+0,35	+0,23	+0,13
Acidity	+0,06	+0,29	-0,60
Astringency	-0,05	+0,43	-0,46
Bitterness	+0,16	-0,03	+0,26
Fullness	+0,34	-0,15	-0,15
Flavour	+0,43	0,00	+0,04
Harmony	+0,44	0,00	-0,06

Chenin blanc, 1978 vintage: Three factors, explaining 69,3 per cent of the variation were extracted from the data. Factor loadings for the different sensory parameters for this wine are given in Table 5.

Summary of results for the different judges over all four reference wines: When all penalising points are accumulated, the total scores are as set out in Table 6. To retain only the most consistent judges, an arbitrary cut-off point of 60 points was decided upon, leaving 7 panellists for the calculation of the final scores and the evaluation of the Nietvoorbij score card. Before continuing the analysis, therefore, the sensory evaluations of judges 3, 9, 10, 12, 13, 17, 18, 20, 21, 24 and 25 were deleted from the data file.

Analysis of the 1979 sensory evaluation exercise using the selected panel: *Comparison of factor loadings for three reference wines after PCA:* Scores of individual wine characteristics for the selected panel were again subjected to PCA as set out before. To examine individual wine descriptors and their correlation to overall quality, the latter property was not included as an independent variable. In order to simplify explanation of the factor structure, a Varimax rotation (program KAVARI) was applied (McBoyle, 1971). This has the effect of maximizing the variance of the loadings of the squared elements (Preston-Whyte, 1974). The factor loadings of the rotated solutions for the four wines, as well as for three cultivars combined, are given in Table 7.

TABLE 5
Factor loadings for the first three principal components:
1978 Chenin blanc

Factor	I	II	III
Variation explained	38,8 per cent	15,5 per cent	15,0 per cent
Overall quality	-0,40	+0,02	-0,23
Clarity	-0,26	+0,28	+0,20
Colour	-0,31	+0,21	+0,11
Typicality	-0,31	+0,08	+0,40
Maturation bouquet	-0,15	-0,07	+0,54
Purity	-0,19	+0,24	+0,14
Acidity	+0,20	+0,33	+0,39
Astringency	+0,05	+0,60	-0,27
Bitterness	-0,14	-0,57	+0,20
Fullness	-0,14	-0,57	-0,20
Flavour	-0,43	-0,03	-0,07
Harmony	-0,43	+0,02	-0,08

TABLE 6
Accumulated inconsistency scores for the 19 judges over all four wines

Judge number	Score
1	15
2	40
3	96
7	33
9	144
10	81
11	58
12	68
13	174
16	40
17	73
18	207
20	164
21	96
23	48
24	159
25	93
26	48

TABLE 7
Matrix of rotated factor loadings

Wine	Chenin blanc 1978			Colombard 1979			Cabernet 1979			Combination of three wines		
	I	II	III	I	II	III	I	II	III	I	II	III
Cumulative variance per cent	19,8	39,0	56,7	28,0	41,5	54,6	19,0	37,9	56,3	25,6	43,4	56,2
Variable												
Clarity	0,00	0,09	-0,02	0,02	-0,14	-0,01	-0,01	0,10	-0,03	-0,03	0,04	-0,08
Colour	0,11	0,66	-0,02	0,15	-0,50	-0,41	-0,06	-0,01	0,02	-0,02	0,35	-0,62
Typicality	0,06	0,16	0,35	-0,49	-0,13	-0,14	-0,46	0,45	0,13	-0,50	0,25	-0,12
Maturation bouquet	0,22	0,15	0,05	0,07	-0,07	-0,16	0,01	0,03	-0,07	-0,02	0,07	-0,06
Purity	-0,08	-0,02	0,08	-0,53	0,07	-0,15	-0,69	0,05	-0,01	-0,58	-0,19	0,10
Acidity	-0,05	-0,09	-0,64	-0,05	-0,05	0,01	0,01	0,04	-0,02	0,02	-0,09	0,21
Astringency	-0,68	-0,07	-0,06	0,08	0,04	-0,11	0,36	-0,07	-0,41	0,04	0,01	-0,07
Bitterness	0,68	0,07	0,06	-0,03	0,78	-0,12	0,09	0,11	0,68	-0,06	0,15	0,75
Fullness	0,08	0,04	0,13	-0,16	0,08	-0,75	0,01	0,66	0,10	-0,13	0,66	-0,02
Flavour	0,07	0,27	0,63	-0,49	0,15	-0,22	-0,29	0,48	0,29	-0,48	0,36	0,05
Harmony	0,05	0,64	0,17	-0,41	0,25	-0,36	-0,28	0,29	0,49	-0,40	0,44	0,13

In the case of Chenin blanc it is interesting to note the difference in relative factor loadings after removal of overall quality as a variable, and the elimination of inconsistent panelists. In Table 5 the main emphasis falls on overall quality, flavour and harmony for the most important component, whereas astringency and bitterness seem to come forward in Table 7. However, it must be borne in mind that the more even distribution of percentage explained variation after the application of the Varimax routine, might have affected the results in a quantitative way, but the principle remains the same – a shift from flavour and harmony towards bitterness and astringency – when the selected panelists were used. This could indicate that these properties play an important part in the variation of Chenin blanc wine scores, and may warrant further investigation. Harmony and colour, together with flavour and acidity, are prominent in the second and third factors respectively (Table 7).

The Colombard wine was evaluated in a more conventional way by the selected panel, with a notable absence of astringency as an important factor, especially when compared with the above results for Chenin blanc. Panelists were depending heavily on purity when evaluating the 1979 Cabernet Sauvignon wine, with typicality also figuring prominently in the first principal factor.

In combining the three data sets, component one is heavily weighted by typicality, purity and flavour, with harmony less important. Factor two is dominated by fullness, which came forward especially in the Cabernet Sauvignon evaluation, with further emphasis on harmony. Colour is important in factor three, but bitterness is the most heavily loaded.

Variables which were never loaded to a significant extent in Table 7 were clarity and maturation bouquet, with astringency only important in the Chenin blanc wine.

Correlation of sensory parameters to overall quality: Least square multi-linear regression analysis (Harper *et al.*, 1977) was used in the sense that scaled data for the 11 individual sensory variables were fitted to overall quality data for the three reference wines as a combined data set. The weights judges placed on the individual parameters are given in Table 8. Extreme discrepancies as reported by Kwan & Kowalski (1980 a) do not exist between real and theoretical score card weights, but it must be borne in mind that only three reference wines are evaluated in this case, and not reference plus other wines (albeit all of the same cultivar) as reported by them. However, it does appear that a parameter like fullness is not used or understood properly by the panelists, although some correlation with overall quality does exist. As pointed out in discussing Table 7, astringency seems to be important regarding the Chenin blanc evaluations, but in this data set this parameter appears relatively unimportant. The same applies to acidity and maturation bouquet. The remainder of the parameters have weights more or less in line with those on the score card, but this data set is too small to comment within reason upon the possible modification of the score card. What does become clear, is that an investigation should be made into the use of the abovementioned parameters. Decisions must be made

TABLE 8
Least square fitting of 106 wine evaluations to overall quality

Feature	Weight		Correlation to quality
	Weight in equation	Weight according to score card	
Clarity	-3	5	-0,10
Colour	2	5	0,26
Typicality	5	5	0,76
Maturation bouquet	-0,9	5	0,07
Purity	5	3	0,66
Acidity	-0,3	5	-0,13
Astringency	0,1	5	-0,13
Bitterness	2	3	0,29
Fullness	0,2	5	0,44
Flavour	3	5	0,75
Harmony	5	5	0,79
Fit correlation			R = 0,92

about further training of panelists in the evaluation of wines when using these descriptors, or the possible elimination of some of the descriptors from the score card. The importance of purity in the assigned weights deserves consideration, and could indicate that the judges attach too much importance to this parameter, its weight being more than in the actual score card. Conversely, it could be reasoned that flavour does not have enough weight, and that panelists should reconsider their evaluation of this important parameter.

Fitting of sensory evaluation data to quality scores for 485 experimental wines: Previous studies (Kwan & Kowalski 1980 a) have used only one variety to evaluate or compare different score cards. When the weights allocated to different wine descriptors for the three different cultivars are compared, it becomes clear that weights could differ in accordance with this factor. Furthermore, depending upon the relative abundance in a data set of a certain wine variety, a certain subset of descriptors could attain larger relative weights. To evaluate a score card, it must be tested using a large variety of wine cultivars, but it must always be borne in mind that the outcome depends heavily on the abovementioned factors. The application of LEAST to the data for 485 experimental wines evaluated with the Nietvoorbij score card in 1979, resulted in the output given in Table 9. It can be seen that all the descriptors have much bigger weights than would be expected from the actual score card weights, except for bitterness and astringency. Large deviations are found regarding harmony, flavour, purity, maturation bouquet, typicality and colour. It is clear that if the wine descriptors in this case are intended to give a total quality picture of the wine and should, therefore, correlate with overall quality, the score card weights should be changed. This change should be combined with specific training in the use of descriptors like bitterness and astringency. The wide range of characteristics found in experimental wines could be the main cause for the deviation from expected weights. For a descriptor like purity to attain such high relative importance could only be due to related problems with many of the wines. The high correlation of this descriptor to overall quality supports this view to a certain extent. Other high correlations to

overall quality for typicality, fullness, flavour and harmony indicate, especially for the latter two, an inclination in panelists to relate them subconsciously to the overall quality of the wine. The same results, in terms of high relative weights, were also found by Kwan & Kowalski (1980 a).

TABLE 9
Least square fitting of 485 experimental wine evaluations to overall quality

Feature	Weight in equation	Weight according to score card	Correlation to quality
Clarity	- 7,8	5	-0,03
Colour	10,9	5	0,47
Typicality	16,1	5	0,60
Maturation bouquet	-17,1	5	0,09
Purity	23,9	3	0,60
Acidity	- 5,3	5	-0,21
Astringency	- 2,2	5	-0,01
Bitterness	- 0,4	3	0,15
Fullness	7,8	5	0,60
Flavour	21,0	5	0,74
Harmony	32,2	5	0,76
Fit correlation			R = 0,88

SUMMARY AND CONCLUSIONS

The data processing methods used in this study show great promise for the evaluation of sensory evaluation scores as a means of selecting consistent panelists. Other related methods should be explored, but sufficient proof of usefulness of PCA for this purpose has been demonstrated, and it should, therefore, be employed on a regular basis. It became clear in the study that wine cultivar can affect panelists differently, and such tests should preferably be made for both red and white wines to ensure an unbiased selection of panelists for the main wine types to be evaluated eventually. Specific shortcomings of individual panelists can be indicated by this method, making it useful for training as well as selection.

The application of the least-squares multi-linear regression method to relate wine descriptors to overall quality scores, produced interesting results when applied to three reference wines as a combined data set. Weights attributed to the different wine descriptors make possible an evaluation of the relative importance of the parameters as used by a panel selected by PCA

for consistency. It became clear that certain parameters are not used properly, whereas others are weighted out of proportion to actual score card weights. Evidence was found that properties like harmony, flavour, typicality and to a lesser extent purity, correlate highly with overall quality, indicating a possible subconscious equation of especially the first three parameters with overall wine quality.

The above inferences apply to, and have great usefulness in evaluating scoring systems when the wines are all in a healthy condition with no large deviations in quality and purity. Applying the same principle to a large and heterogeneous data set containing nearly 500 wines, produced results which were not consistent with the scoring system, especially with regard to the weights judges placed on individual parameters. Purity, almost certainly because many experimental wines are deficient in this respect, gained an excessive relative weight in the equation, whereas bitterness was hardly used on a large scale. It is difficult, therefore, to propose specific score card modifications at this stage. A point to debate seems to be whether one should try to relate wine descriptors to overall quality in the case of widely diverging experimental wines. Further investigations along similar lines, using alternative score cards, should provide a better understanding of this specific aspect of wine evaluation, and should be allocated a research priority alongside oenological and viticultural investigations using these very data to evaluate their own results.

LITERATURE CITED

- HARPER, A. M., DUEWER, D. K., KOWALSKI, B. R. & FASCHING, J. L., 1977. ARTHUR and experimental data analysis: The heuristic use of a polyalgorithm. "ACS Symposium No. 52" American Chemical Society, Washington D C.
- KWAN, W. O. & KOWALSKI, B. R., 1980 a. Data analysis of sensory scores. Evaluation of panelists and wine score cards. *J. Food Sci.* **45**, 213-216.
- KWAN, W. O. & KOWALSKI, B. R., 1980 b. Correlation of objective chemical measurements and subjective sensory evaluations. Wines of *Vitis vinifera* variety Pinot noir from France and the United States. *Anal. Chim. Acta* **122**, 215-222.
- McBOYLE, G. R., 1971. Climatic classification of Australia by computer. *Austr. Geogr. Studies* **9**, 1-14.
- PRESTON-WHYTE, R. A., 1974. Climatic classification of South Africa: A multivariate approach. *South Afr. Geogr. J.* **56**, 79-86.
- TROMP, A. & CONRADIE, W. J., 1979. An effective system for sensory evaluation of experimental wines. *Am. J. Enol. Vitic.* **30**, 278-283.
- WU, L. S., BARGMANN, R. E. & POWERS, J. J., 1977. Factor analysis applied to wine descriptors. *J. Food Sci.* **45**, 213-216.