



A note on Mahalanobis and related distance measures in WinISI and Unscrambler

Journal:	<i>Journal of Near Infrared Spectroscopy</i>
Manuscript ID	Draft
Manuscript Type:	Original Research Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Garrido-Varo, Ana; University of Cordoba, Non-Destructive Spectral Sensors Unit, Faculty of Agriculture and Forestry Engineering Garcia-Olmo, Juan; University of Cordoba, NIR/MIR Spectroscopy Unit, Central Service for Research Support Fearn, Tom; University College London, Statistical Science
Keywords:	Mahalanobis distance, Leverage, Hotelling's T2, Principal components, Near infrared spectroscopy, Outliers
Abstract:	<p>In identifying spectral outliers in near infrared calibration it is common to use a distance measure that is related to Mahalanobis distance. However, different software packages tend to use different variants, which leads to a translation problem if more than one package is used. Here the relationships between squared Mahalanobis distance D_2, the GH distance of WinISI, and the T2 and leverage (L) statistics of Unscrambler are established as $D_2 = T_2 \square L \cdot n \square GH \cdot k$, where n and k are the numbers of samples and variables respectively in the set of spectral data used to establish the distance measure. The implications for setting thresholds for outlier detection are discussed. On the way to this result the principal component scores from WinISI and Unscrambler are compared. Both packages scale the scores for a component to have variances proportional to the contribution of that component to total variance, but the WinISI scores, unlike those from Unscrambler, do not have mean zero.</p>

SCHOLARONE™
Manuscripts

A note on Mahalanobis and related distance measures in WinISI and Unscrambler

A Garrido-Varo¹, J Garcia-Olmo² and T Fearn³

¹Non-Destructive Spectral Sensors Unit, Faculty of Agriculture and Forestry Engineering, University of Cordoba, Spain

²NIR/ MIR Spectroscopy Unit, Central Service for Research Support, University of Cordoba, Spain

³Department of Statistical Science, University College London, UK.

Corresponding author: T Fearn, Department of Statistical Science, UCL, Gower Street, London WC1E 6BT, UK

Email: t.fearn@ucl.ac.uk

Keywords

Mahalanobis distance, Leverage; Hotelling's T^2 , Principal components, Near infrared spectroscopy, Outliers

Abstract

In identifying spectral outliers in near infrared calibration it is common to use a distance measure that is related to Mahalanobis distance. However, different software packages tend to use different variants, which leads to a translation problem if more than one package is used. Here the relationships between squared Mahalanobis distance D^2 , the GH distance of WinISI, and the T^2 and leverage (L) statistics of Unscrambler are established as $D^2 = T^2 \approx L \cdot n \approx GH \cdot k$, where n and k are the numbers of samples and variables respectively in the set of spectral data used to establish the distance measure. The implications for setting thresholds for outlier detection are discussed. On the way to this result the principal component scores from WinISI and Unscrambler are compared. Both packages scale the scores for a component to have variances proportional to the contribution of that component to total variance, but the WinISI scores, unlike those from Unscrambler, do not have mean zero.

Introduction

One of the necessary steps in developing or applying near infrared (NIR) calibrations is to check for spectral outliers. A common way to decide whether a given spectrum is an outlier is to calculate, using a distance measure that takes into account the pattern of spectral variability in the training set, its distance from the mean spectrum of that set. This distance can then be compared with some threshold that is either based on an assumption of some statistical distribution or is simply a rule of thumb based on experience.

1
2
3 So long as the user is faithful to one software package this approach is simple to
4 apply. Problems can arise however if more than one package is used, because
5 different packages tend to use different variants of the same underlying measure,
6 Mahalanobis distance [1, 2]. Then the sort of question that can arise is "A threshold
7 of 3 for GH in WinISI works well for my applications, what is the corresponding
8 threshold for a leverage from Unscrambler?" The investigations reported here were
9 carried out with the aim of answering some of these questions, by comparing the
10 distance measures of WinISI and Unscrambler with Mahalanobis distances
11 calculated from the same data set by Matlab code that implements the textbook
12 formula.
13
14
15

16 An added complication is that the Mahalanobis formula involves the inversion of a
17 variance matrix calculated from the spectra in the training set. This inversion is
18 unstable in high dimensions and so the spectra need to be projected onto a lower
19 dimensional space before the distances can be calculated. The obvious options are
20 to use scores on either principal components (PCs) or partial least squares (PLS)
21 factors. The use of PCs has the advantage that one can use them to screen for
22 outliers before developing calibrations. PCs are also simpler to calculate and much
23 more likely to match between different software packages, and so this is the
24 approach adopted here. Given that PC scores needed to be calculated in each of
25 the three packages, the opportunity was taken to compare the scores also.
26
27
28

29 A priori the PC scores might be expected to differ between packages, because there
30 is more than one option for scaling them, for example a vector of scores can be
31 scaled to have length 1 or a squared length reflecting the contribution of the PC to
32 total variance, to list just two of the most common options. In addition the sign of the
33 PC loadings, and hence of the scores, is arbitrary, because if v is an eigenvector of
34 the matrix M then so is $-v$. Different packages will often produce scores with
35 different signs, and even using the same package the removal of one spectrum from
36 the calculation can result in the sign of the PC flipping. This is unimportant, but can
37 be disconcerting when a plot appears to change completely after a very small
38 change to the data. None of this should matter so far as the distance calculations
39 are concerned, since Mahalanobis distance is scale invariant, but it is still of interest
40 to compare the scores from the three packages.
41
42
43
44

45 **Materials and Methods**

46 **Data**

47
48 The data set used to compare the results on different software comprised 349
49 spectra of liquid samples of subcutaneous fat of Iberian pigs, measured on a Foss-
50 NIRSystems 6500 monochromator. The wavelength range was 400 to 2498nm in
51 steps of 2nm, and thus the data matrix X was of dimension 349 x 1050. Since the
52 purpose of the current investigation was simply to compare the software, no pre-
53 treatments were applied to the spectra. The data were exported to a .csv file for
54 transfer to Matlab, and to JCAMP-DX format for transfer to Unscrambler.
55
56
57

58 **Software**

The Matlab code in the appendices was run using version R2016b (The MathWorks Inc., Natick, MA, USA). The WinISI software was version 4.8 (FOSS Analytical A/S, Hillerød, Denmark), and the Unscrambler software was Unscrambler X version 10.4.1 (CAMO Software AS, Oslo, Norway). Calculations with much earlier versions of Win ISI and Unscrambler (see the acknowledgement) gave equivalent results.

Mahalanobis distance, Hotelling's T^2 , and leverage

Mahalanobis [3] invented the statistic that bears his name [1] as a way to measure the distance between two groups of observations in k -dimensional space while taking into account the fact that the k -variables may have differing scales and may be intercorrelated. In this case the formula for the squared distance between the groups would be

$$D^2 = (m_1 - m_2)^T S^{-1} (m_1 - m_2)$$

where m_1 and m_2 are the $k \times 1$ vectors of means for the two groups and S is a $k \times k$ within-group variance matrix, all of these quantities being estimated from the data on the two groups. This statistic is closely related to the subsequently developed Hotelling's T^2 , which is a multivariate version of the two-sample t -test. The relationship when there are n_1 observations in group 1 and n_2 in group 2 is

$$T^2 = (n_1 n_2 / (n_1 + n_2)) D^2$$

When the observations come from multivariate normal distributions, the distribution of T^2 , and hence that of D^2 , is known to be a multiple of an F distribution. Full details of all the above can be found in almost any multivariate statistics textbook, for example the one by Krzanowski [2].

In NIR calibration a version of D^2 is commonly used to measure the distance of a single spectrum, or more precisely the scores of this spectrum on a set of PCs or PLS factors, from the centre of the cloud of calibration set spectra. Then the formula becomes

$$D^2 = (x - m)^T V^{-1} (x - m) \tag{1}$$

where x is the $k \times 1$ vector of spectral data for the observation of interest, and the $k \times 1$ mean vector m and $k \times k$ variance matrix $V = X^T X / (n - 1)$ are both calculated from the $n \times k$ matrix X of spectral data for the calibration set.

With m and V both calculated from the full calibration set of n observations, most of the random variability in this version of D^2 comes from x . If we assume x to be randomly sampled from the same multivariate normal distribution as the calibration set, and ignore that fact that m and V are sample estimates of population parameters, then the distribution of D^2 will be approximately chi-squared on k degrees of freedom. This distribution has a mean value of k . This same approximate distribution applies regardless of whether x belongs to the calibration set or is a new observation.

Leverage L is a statistic developed for identifying influential observations in multiple linear regression [4]. It is also used to identify outliers in NIR calibration [5, 6]. Starting from the standard statistical definition its relationship with D^2 should be

$$L = 1/n + D^2/(n-1). \quad (2)$$

The factor of $(n-1)$ arises because leverage uses $(X^T X)^{-1}$ in place of V^{-1} in a formula analogous to Equation 1, and the $1/n$ represents the influence of x on the mean m . Because of the context in which it was designed to be used, there is an assumption here that x is one of the n rows of X and so has contributed to the estimation of m . In cases where it is not, for example in comparing the spectra of prediction samples with those of a calibration set, it would make sense to omit the $1/n$, though if this makes any practical difference the calibration set is too small.

PCA calculations

PCA scores were calculated in Matlab using the code in Appendix 1 with the scaling parameter 'scal' set to 1. This uses the Matlab SVD function to decompose X and scales each vector t of scores so that its squared length $t^T t$ is equal to the corresponding eigenvalue of $X^T X$. In other words, the variance of the scores for a component is proportional to the contribution of that component to the total variance in X . The columns of X are centered but not rescaled in the computation. This would correspond to the 'covariance' option in most standard statistical packages. In WinISI there are no options; in Unscrambler mean centering and the SVD algorithm were selected.

Distance calculations

Squared Mahalanobis distances D^2 were calculated in Matlab using the code in Appendix 2. This implements the textbook formula in Equation 1. In WinISI and Unscrambler the desired statistics were selected from the appropriate menus, choosing to base the calculations on 10 PCs in each case.

Results and discussion

The aim of this investigation was to relate the formulas used by the packages compared, not to establish the accuracies of the computations. Thus statements like 'the scores matched' should be interpreted as meaning only that the correspondence was good enough to establish equivalence beyond reasonable doubt, not as a claim of identity to the level of machine precision. In any case the rounding errors involved in the transfers of data probably dominate the errors in any internal computations.

Comparison of PC scores

As expected, there were differences in signs between some of the scores returned by the three programs. With these differences resolved the scores produced by Unscrambler matched those from Matlab. The scores from WinISI had the same scaling as those from the other two programs but, unlike the other two sets, were not centred on zero. Figure 1 shows the Matlab and WinISI scores on the first two PCs after the directions of the Matlab scores have been reversed so that the plots match.

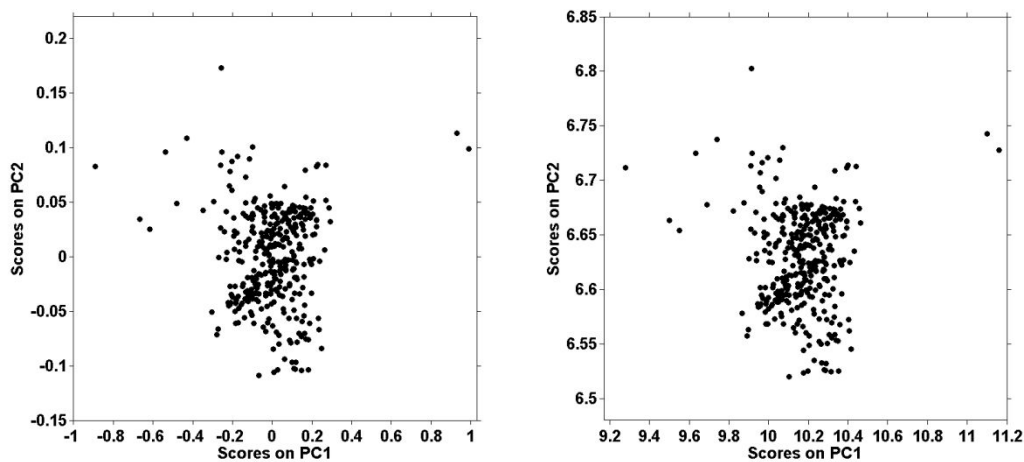


Figure 1. Scatter plots of first 2 PC scores from Matlab (left) and WinISI (right).

Further comparisons revealed that although WinISI centres the columns of X before carrying out the PCA, the scores it produces correspond to applying the loadings to an uncentred X . Running the Matlab code in Appendix 1 and then calculating $X \cdot L$ reproduces the WinISI scores. This is presumably done because it simplifies the calculation of scores for future samples by eliminating the need to subtract the mean spectrum of the set used to carry out the PCA.

Comparison of Mahalanobis and other distances

The T^2 results from Unscrambler correspond to the squared Mahalanobis distances D^2 from the Matlab program. The leverages L correspond to $D^2/(n-1) + 1/n$ as in Equation 2. The Unscrambler reference manual [7], which is generally quite precise, clearly defines leverage as the standard statistic, but is uncharacteristically vague about T^2 .

The relation between WinISI's GH and D^2 was found to be

$$GH = (n/(n-1)) \cdot D^2/k$$

where k is the number of PCs used for the calculation of D^2 and GH. The factor $n/(n-1)$ is presumably due to the use by WinISI of a divisor of n rather than the more usual $n-1$ in the calculation of the variance matrix in the Mahalanobis formula. The division by k scales GH to have typical values of around 1 whatever the value of k .

Ignoring subtleties like factors of $n/(n-1)$, which is 1.003 for the data set used here for example, the relationships may be summarised as

$$D^2 = T^2 \approx L \cdot n \approx GH \cdot k. \quad (3)$$

Implications for thresholds

Using the approximate relationships in Equation 3, a threshold of T_M for squared Mahalanobis distance corresponds to thresholds of T_M/k for GH, T_M for Unscrambler's T^2 statistic and T_M/n for Unscrambler's leverage statistic. So, for example, the GH rule of thumb of 3 would convert to $3k$ for squared Mahalanobis distance or for T^2 , and to $3k/n$ for leverage. In fact $3k/n$, along with $2k/n$, is a commonly suggested rule of thumb for leverage [6].

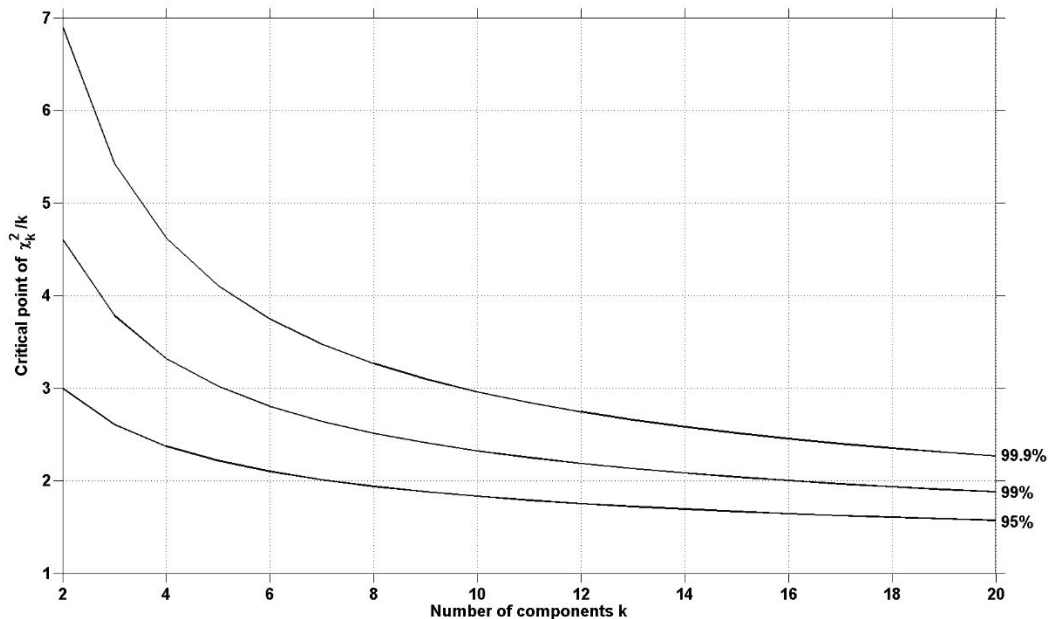


Figure 2. Thresholds for GH at three probability levels based on a chi-squared distribution with k degrees of freedom for D^2

The alternative to a rule of thumb is to base a threshold on a probability distribution. If we assume a multivariate normal distribution for the PC scores, the approximate distribution for D^2 is chi-squared on k degrees of freedom. Figure 2 shows thresholds at probability levels of 95, 99 and 99.9% for GH based on this distribution. While these probabilities should not be taken too seriously, since the normality assumption is unlikely ever to be correct, the figure does suggest that 3 is a sensible choice for a fixed threshold for GH, at least for modest k .

Acknowledgements

The authors would like to dedicate this article to the memory of Prof. Dr. Tomas Isaksson, who passed away on 12 July 2012. This work began as a collaboration between the first two authors and Tomas, paused upon his death, and has only now been completed.

This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

Declaration of conflicting interests

The authors declare that there is no conflict of interest.

References

[1] Mahalanobis PC. On tests and measures of group divergence. *Journal of the Asiatic Society of Bengal* 1930; 26; 541-588.

[2] Krzanowski WJ. *Principles of Multivariate Analysis: A User's Perspective*. 2nd ed. New York: OUP, 2000.

[3] Ghosh JK and Majumdar PP. Mahalanobis, Prasantra Chandra. In: Armitage P and Colton T (eds) *Encyclopaedia of Biostatistics*. New York: Wiley, 1998, pp. 2372-2375.

[4] Weisberg S. *Applied Linear Regression*. 3rd ed. New York: Wiley, 2005.

[5] Martens H and Næs T. *Multivariate Calibration*. Chichester UK: Wiley, 2001.

[6] Næs T, Isaksson T, Fearn T and Davies T. *Multivariate Calibration and Classification*. Chichester UK: NIR Publications, 2002.

[7]

<https://www.camo.com/downloads/U9.6%20pdf%20manual/The%20Unscrambler%20Method%20References.pdf>

Appendix 1. Matlab code for PCA

This function uses Matlab's singular value decomposition on the $n \times p$ matrix X , thus avoiding the need to compute the matrix product $X^T X$. The 'econ' option stops the decomposition when the number of eigenvalues extracted corresponds to the smaller of n and p , the maximum number of nonzero eigenvalues for a matrix of this size. The PCA scores and loadings are easily computed from this decomposition.

```
function [S,L,v,m] = pcomp(X,k,scal)
% PCA
% Usage [S,L,v,m] = pcomp(X,k,scal)
% Inputs
%   X ..... n x p matrix of spectra (in rows)
%   k ..... number of components to return
%   scal .. scalar, options for scaling the scores
%           0 - orthonormal scores
%           1 - scores scaled to have squared length equal to
%               the corresponding eigenvalue of X'X
% Outputs
%   S .... n x k matrix of scores
%   L .... p x k matrix of loadings
%   v .... k x 1 vector of eigenvalues of X'X
%   m .... 1 x p vector, column means of X
%
% Note: to calculate scores for a new data matrix M use
%       Scor = (M-ones(size(M,1),1)*m)*L
```



```

1
2
3
4 % Tom Fearn, February 2019
5
6 m = mean(X,1); % column means of X
7 Xc = X - ones(size(X,1),1)*m; % centre columns of X
8 [U,E,V] = svd(Xc,'econ'); % decompose Xc as U*E*V'
9 e = diag(E); % all the eigenvalues of X
10 e = e(1:k); % first k eigenvalues of X
11 v = e.^2; % first k eigenvalues of X'X
12
13 if scal==0
14     S = U(:,1:k); % k scores, orthonormal
15     L = V(:,1:k)*diag(1./e); % k loadings, scaled to give
16 % orthonormal scores
17
18 else
19     S = U(:,1:k)*diag(e); % k scores, scaled by eigenvalues
20     L = V(:,1:k); % k loadings, orthonormal
21
22 end
23
24

```

Appendix 2. Matlab code for Mahalanobis distance

Two separate functions were used. The first, MVinv, calculates the mean and inverse variance matrix from a set of data, the second takes these as inputs and calculates Mahalanobis distances. To use the PCA function above and the two functions below to calculate Mahalanobis distances using k PCs from a data matrix X, the code would be

```

32 [S,L,v,mx] = pcomp(X,k,1);
33 [ms,Vi] = MVinv(S);
34 D = MD2(S,ms,Vi);

```

Function MVinv

```

39 function [m,Vi] = MVinv(X)
40 % Calculates mean vector and inverse variance matrix of data set X
41 % Usage[m,Vi] = MVinv(X)
42 % Input
43 % X ... n x p matrix of data, cases in rows
44 %
45 % Outputs
46 % m ... 1 x p vector, column means of X
47 % Vi .. p x p symmetric matrix, inverse of covariance matrix of X
48 %
49 % Notes
50 % 1. If p is very large, and in particular if p>n-1, computing the
51 % inverse of the covariance matrix will give unstable results
52 % 2. The divisor in the computation of the covariance matrix is n-1
53 %
54 %
55 % Tom Fearn, February 2019
56 %
57 m = mean(X);
58 V = cov(X);
59 Vi = inv(V);
60

```

1
2
3 end
4

5 Function MD2

```

6
7 function D = MD2(X,x0,Vi);
8 % Calculates the squared Mahalanobis distance of each row of X from
9 % each of the rows in x0 using the inverse covariance matrix Vi
10 % Usage D = MD2(X,x0,Vi)
11 % Inputs
12 % X ... n x p, data matrix, observations in rows
13 % x0 .. k x p, centres for calc of MD, in rows
14 % Vi .. p x p, inverse variance matrix for calc of MD
15 %
16 % Output
17 % D ... n x k matrix, squared Mahalanobis distances between each
18 % row in X and each row in x0
19 %
20 % Notes
21 % 1. To get distances from a calibration set mean set x0=m where
22 % m is the 1 x p mean vector of the cal set as given by MVinv
23 % 2. To get distances from several individual observations
24 % set these observations as the rows of x0
25 % 3. This will obviously crash if the dimensions p of the inputs
26 % do not match!
27 %
28 %
29 % Tom Fearn, February 2019
30 %
31 % make sure that if x0 is a vector it is a row vector
32 if size(x0,2)==1; x0=x0'; end;
33 %
34 % set up storage for results
35 n = size(X,1); k = size(x0,1);
36 D = zeros(n,k);
37 %
38 % loop over rows of x0
39 for i = 1:k
40     m = x0(i,:); % set i'th row as as the centre
41     Xc = X - ones(n,1)*m; % center X
42     D(:,i) = sum((Xc*Vi).*Xc,2); % calculate squared MDs
43     % Note: The more obvious code would be diag(Xc*Vi*Xc') but
44     % this would be less efficient
45 end
46 end
47
48
49
50
51
52
53
54
55
56
57
58
59
60
```

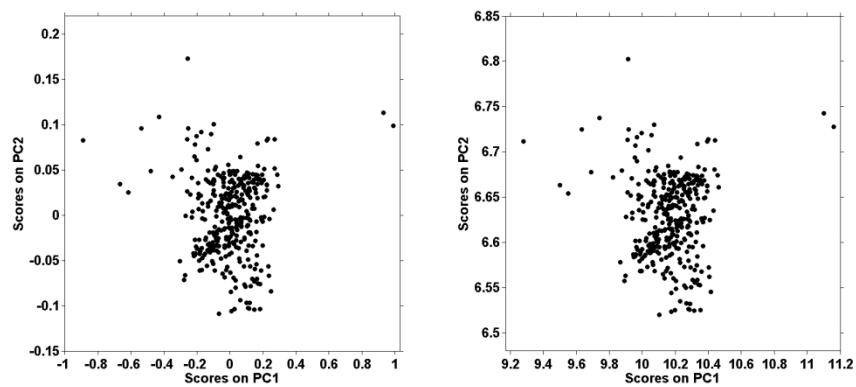


Figure 1. Scatter plots of first 2 PC scores from Matlab (left) and WinISI (right).

508x287mm (96 x 96 DPI)

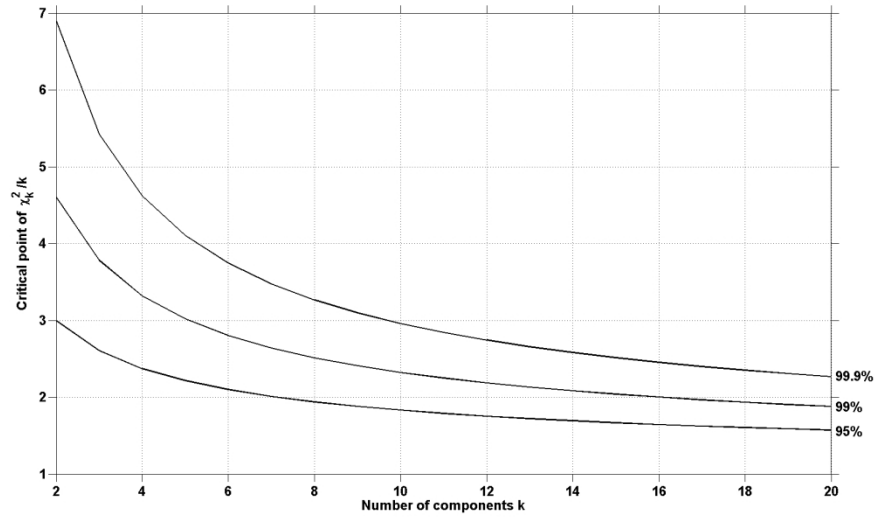


Figure 2. Thresholds for GH at three probability levels based on a chi-squared distribution with k degrees of freedom for D2

508x287mm (96 x 96 DPI)