MathWorks® Products Solutions Academia Support Community

Examine the Gaussian Mixture Assumption

Documentation Search R2018a Documentation Documentation -**E** CONTENTS Translate This Page

discriminant analysis to be a good classifier. Furthermore, the default linear discriminant analysis assumes that all class covariance matrices are equal. This section shows methods to check these assumptions:

Discriminant analysis assumes that the data comes from a Gaussian mixture model (see Creating Discriminant Analysis Model). If the data appears to come from a Gaussian mixture model, you can expect

Contact Us How to Buy James ▼

R2018a

Trial Software Product Updates **Bartlett Test of Equal Covariance Matrices for Linear Discriminant Analysis**

The Bartlett test (see Box [1]) checks equality of the covariance matrices of the various classes. If the covariance matrices are equal, the test indicates that linear discriminant analysis is appropriate. If not, consider using quadratic discriminant analysis, setting the DiscrimType name-value pair to 'quadratic' in fitcdiscr.

The Bartlett test assumes normal (Gaussian) samples, where neither the means nor covariance matrices are known. To determine whether the covariances are equal, compute the following quantities:

- Sample covariance matrices per class σ_i , $1 \le i \le k$, where k is the number of classes.
- Pooled-in covariance matrix σ . • Test statistic *V*:
- $V = (n k)\log(|\Sigma|) \sum_{i=1}^{k} (n_i 1)\log(|\Sigma_i|)$
- where n is the total number of observations, and n_i is the number of observations in class i, and $|\Sigma|$ means the determinant of the matrix Σ .

• Asymptotically, as the number of observations in each class n_i become large, V is distributed approximately χ^2 with kd(d+1)/2 degrees of freedom, where d is the number of predictors (number of dimensions in

- the data). The Bartlett test is to check whether V exceeds a given percentile of the χ^2 distribution with kd(d+1)/2 degrees of freedom. If it does, then reject the hypothesis that the covariances are equal.
- Check whether the Fisher iris data is well modeled by a single Gaussian covariance, or whether it would be better to model it as a Gaussian mixture.

prednames = {'SepalLength','SepalWidth','PetalLength','PetalWidth'};

Example: Bartlett Test for Equal Covariance Matrices

load fisheriris;

```
L = fitcdiscr(meas, species, 'PredictorNames', prednames);
Q = fitcdiscr(meas, species, 'PredictorNames', prednames, 'DiscrimType', 'quadratic');
D = 4; % Number of dimensions of X
Nclass = [50 50 50];
N = L.NumObservations;
K = numel(L.ClassNames);
SigmaQ = Q.Sigma;
SigmaL = L.Sigma;
logV = (N-K)*log(det(SigmaL));
for k=1:K
    logV = logV - (Nclass(k)-1)*log(det(SigmaQ(:,:,k)));
end
nu = (K-1)*D*(D+1)/2;
pval = 1 - chi2cdf(logV,nu)
```

pval = 0

load fisheriris;

D = 4;

figure;

K = numel(L.ClassNames);

xlabel('Expected quantile');

ylabel('Observed quantile');

load fisheriris;

N = L.NumObservations;

xlabel('Expected quantile');

18

16

ylabel('Observed quantile for QDA');

virginica versicolor

setosa

• *d* is the number of predictors (number of dimensions in the data).

• *n* is the total number of observations.

expKurt = D*(D+2);

pval = 0.7230

Functions

Objects

cvshrink|fitcdiscr

Functions and Other Reference

Release Notes

MathWorks

and scientists.

Discover...

MathWorks is the leading developer of

mathematical computing software for engineers

PDF Documentation

line([0 20],[0 20],'color','k');

gscatter(expQ,mahQ,Q.Y(sorted),'bgr',[],[],'off');

legend('virginica','versicolor','setosa','Location','NW');

line([0 20],[0 20],'color','k');

For linear discriminant analysis, use a single covariance matrix for all classes.

L = fitcdiscr(meas, species, 'PredictorNames', prednames);

expQ = chi2inv(((1:N)-0.5)/N,D); % expected quantiles

gscatter(expQ,mahL,L.Y(sorted),'bgr',[],[],'off');

[mahL, sorted] = sort(mahL); % sorted obbserved quantiles

legend('virginica','versicolor','setosa','Location','NW');

mahL = mahal(L,L.X, 'ClassLabels',L.Y);

prednames = {'SepalLength','SepalWidth','PetalLength','PetalWidth'};

N = L.NumObservations;

```
18
                       virginica
    16
                       versicolor
                       setosa
                       data1
Observed quantile
```

prednames = {'SepalLength','SepalWidth','PetalLength','PetalWidth'}; Q = fitcdiscr(meas, species, 'PredictorNames', prednames, 'DiscrimType', 'quadratic'); Nclass = [50 50 50];

Overall, the agreement between the expected and observed quantiles is good. Look at the right half of the plot. The deviation of the plot from the 45° line upward indicates that the data has tails heavier than a

K = numel(L.ClassNames); mahQ = mahal(Q,Q.X,'ClassLabels',Q.Y);

As shown in Bartlett Test of Equal Covariance Matrices for Linear Discriminant Analysis, the data does not match a single covariance matrix. Redo the calculations for a quadratic discriminant.

```
data1
        Observed quantile for QDA
            12
            2
                          2
                                            6
                                                    8
                                                             10
                                                                      12
                                                                               14
                                                                                        16
                                           Expected quantile
Mardia Kurtosis Test of Multivariate Normality
The Mardia kurtosis test (see Mardia [2]) is an alternative to examining a Q-Q plot. It gives a numeric approach to deciding if data matches a Gaussian mixture model.
In the Mardia kurtosis test you compute M, the mean of the fourth power of the Mahalanobis distance of the data from the class means. If the data is normally distributed with constant covariance matrix (and is thus
```

The Mardia test is two sided: check whether M is close enough to d(d + 2) with respect to a normal distribution of variance 8d(d + 2)/n. **Example: Mardia Kurtosis Test for Linear and Quadratic Discriminants**

load fisheriris; prednames = {'SepalLength','SepalWidth','PetalLength','PetalWidth'};

N = L.NumObservations; obsKurt = mean(mahL.^2);

mahQ = mahal(Q,Q.X,'ClassLabels',Q.Y); obsKurt = mean(mahQ.^2); [~,pval] = ztest(obsKurt,expKurt,sqrt(varKurt))

Because pval is high, you conclude the data are consistent with the multivariate normal distribution.

See Also

ClassificationDiscriminant | gmdistribution

Related Topics

Statistics and Machine Learning Toolbox Support **Documentation** MATLAB Answers Examples

Explore Products Try or Buy Accelerating the pace of engineering and science MATLAB Downloads Trial Software Simulink

Student Software

Hardware Support

File Exchange

Contact Sales

Pricing and Licensing

Installation Help

Get Support Learn to Use Documentation Tutorials Answers Examples Videos and Webinars Training

About MathWorks Installation Help Careers Newsroom **Social Mission** About MathWorks

Was this topic helpful?

Yes

Join the conversation

United States Patents | Trademarks | Privacy Policy | Preventing Piracy © 1994-2018 The MathWorks, Inc.

discriminant analysis, as opposed to linear discriminant analysis. **Q-Q Plot** A Q-Q plot graphically shows whether an empirical distribution is close to a theoretical distribution. If the two are equal, the Q-Q plot lies on a 45° line. If not, the Q-Q plot strays from the 45° line. **Check Q-Q Plots for Linear and Quadratic Discriminants**

The Bartlett test emphatically rejects the hypothesis of equal covariance matrices. If pval had been greater than 0.05, the test would not have rejected the hypothesis. The result indicates to use quadratic

8 10 12 14 16 Expected quantile

normal distribution. There are three possible outliers on the right: two observations from class 'setosa' and one observation from class 'virginica'.

expQ = chi2inv(((1:N)-0.5)/N,D);[mahQ, sorted] = sort(mahQ); figure;

The Q-Q plot shows a better agreement between the observed and expected quantiles. There is only one outlier candidate, from class 'setosa'.

L = fitcdiscr(meas, species, 'PredictorNames', prednames); mahL = mahal(L,L.X, 'ClassLabels',L.Y); D = 4;

suitable for linear discriminant analysis), M is asymptotically distributed as normal with mean d(d + 2) and variance 8d(d + 2)/n, where

varKurt = 8*D*(D+2)/N; [~,pval] = ztest(obsKurt,expKurt,sqrt(varKurt)) pval = 0.0208The Mardia test indicates to reject the hypothesis that the data is normally distributed. Continuing the example with quadratic discriminant analysis: Q = fitcdiscr(meas, species, 'PredictorNames', prednames, 'DiscrimType', 'quadratic');

References [1] Box, G. E. P. A General Distribution Theory for a Class of Likelihood Criteria. Biometrika 36(3), pp. 317–346, 1949. [2] Mardia, K. V. Measures of multivariate skewness and kurtosis with applications. Biometrika 57 (3), pp. 519–530, 1970.

Discriminant Analysis Classification Bartlett Test of Equal Covariance Matrices for Linear Discriminant Analysis Check Q-Q Plots for Linear and Quadratic Discriminants

> **Bug Reports Product Requirements Software Downloads** Free eBook: Machine Learning with **MATLAB** Download now

> > Consulting **Application Status** License Center