

## Fisher z-transformation:

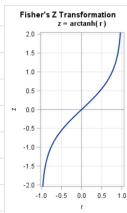
This can be used to inspect and get confidence intervals for the estimated correlation coefficient.

It is based on the variance stabilizing transformation.

$$z = \frac{1}{2} \ln \left( \frac{1+r}{1-r} \right) = \operatorname{arc tanh}(r)$$

This transformation was proposed by Fisher in order to deal with the properties of correlation coefficients which, by definition  $\in [-1, 1]$ .

The issue like that is that when observing and sampling correlation coefficients with a true correlation coefficient lying on the edges you get a very skewed distribution from which it is difficult to make statistical conclusions on the robustness of such coefficients.



Fisher explored the topic further and found that in the case of a bivariate normal distribution of  $(X, Y)$  the above transformation can be shown to be asymptotically normal distributed with

$$\mu = \frac{1}{2} \ln \left( \frac{1+r}{1-r} \right)$$

$$SE = \frac{1}{\sqrt{N-3}}$$

It follows that the above can be used to obtain confidence intervals for the correlation coef. because for large  $N$ :

$$\sqrt{N-3} (z - \mu) \sim N(0, 1)$$

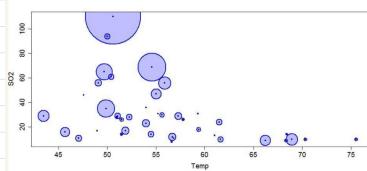
Then once you obtain the confidence interval for  $z$  you can back-transform to  $r$  using the inverse of the  $z$ -transform

$$r = \tanh(z)$$

## Some Multivariate Plotting Techniques:

- Glyph: here you use different plotting attributes (such as size, color, etc.) to represent different variables

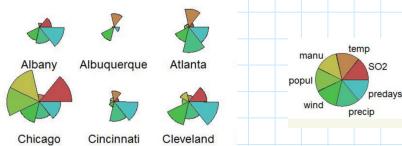
Example: Bubble Plot



- Star plots: Here each star represents one single observation. The spikes then represent the normalized values of each variable.

The normalization that is usually applied is:

$$\frac{x_i - \min(x)}{\max(x) - \min(x)}$$



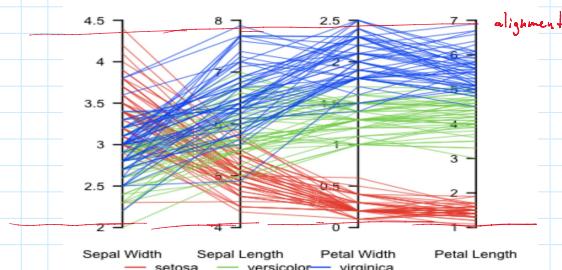
manu temp  
popul SO2  
wind predays  
precip

### - parallel coordinate plots

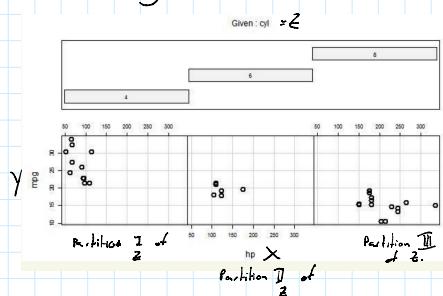
the different variables are expressed on the X-Axis at different positions.

The values for the variables were expressed by the parallel vertical coordinates. This are scaled so that the max and min among all variables are aligned as visible below.

Parallel coordinate plot, Fisher's Iris data



### - Conditioning Plots:

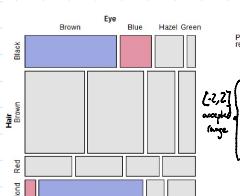


Here the idea is to partition one of the variables (here Z) into equivalent intervals and observe how the scatterplot of other variables behaves for that particular partition.

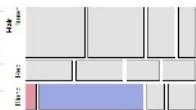
This will help to analyze structural breaks and different patterns in the different partitions.

### - Mosaic Plots:

These are used for visualizing categorical variables where the techniques above fails.



The Pearson residual 3.27 is plotted on the right side of the plot. The text also mentions a 'coupled range' and 'residual'.



frequency with which the categorical combination occurs.

The significance of a given categorical variable being independent with the other. Both the blue and the red colouring above rejects the hypothesis of independence among the variables and reasoning goes as follows:

- Consider the following contingency table.

	$A=1$	$A=2$	Total
$B=1$	$n_{11}$	$n_{12}$	$n_{1\cdot}$
$B=2$	$n_{21}$	$n_{22}$	$n_{2\cdot}$
Total	$n_{\cdot 1}$	$n_{\cdot 2}$	$n$

- Assume that  $A$  and  $B$  are independent. Then it follows:

$$P(A_i \cap B_j) = P(A_i) \cdot P(B_j)$$

$$= \frac{n_{ij}}{n} \cdot \frac{n_{\cdot j}}{n} = \underbrace{\hat{\pi}_{ij}}_{\substack{\text{Expected prob.} \\ \text{of observing } A_i \cap B_j \\ \text{under the Null.}}}$$

It follows now that the expected amount of observations should be  $E_{ij} = n \cdot \hat{\pi}_{ij}$ .

From this you can compute the following statistic:

$$\chi^2 = \sum_{i=1}^2 \sum_{j=1}^2 \frac{(n_{ij} - E_{ij})^2}{E_{ij}}$$

Squared Pearson Residuals.

As under the Null  $\frac{(n_{ij} - E_{ij})^2}{E_{ij}} \sim N(0,1)$

It follows that the above is  $\chi^2_{(2-1)(2-1)}$  degrees of freedom and you can therefore easily get the desired test statistics and significance values.

## Multivariate Outliers.

While it may be easy to recognise outliers in  $\mathbb{R}^2$  and possibly even  $\mathbb{R}^3$  the situation gets tricky in hyperspaces.

A simple solution for finding outliers of multivariate normal distributed data is the one of leveraging the Mahalanobis distance.

Such distance is defined as follows:

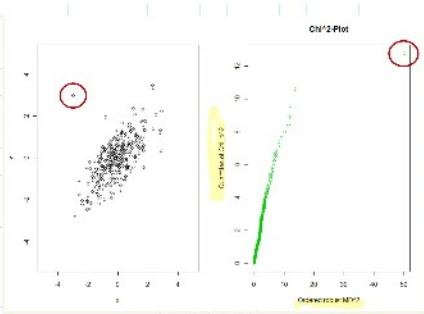
$$\text{1-Dimensional-Case: } \frac{x - \mu}{\sigma} := \text{Z-score}$$

$$\text{N-Dimensional-Case: } \sqrt{(x - \bar{\mu})^T \Sigma^{-1} (x - \bar{\mu})}$$

It is now clear that under the assumption a multivariate normal distribution the squared Mahalanobis distance is  $\chi^2_q$ . This is by the very def of a multivariate normal dist.

It follows now that you can find outliers by checking if the Mahalanobis distance for a particular observation is way in the tail of the corresponding  $\chi^2_q$  distribution.

In order to do so you can even inspect the Q-Q-plot.



Notice that such technique avoids the masking effect where you would not recognize outliers because they appear in clusters.

→ Think in this sense about the dangers of using the Cook's Distance.

Finally, note that albeit in the case of not multivariate distributions you might not be able to make very strong statistical conclusions about outliers high Mahalanobis distances would still point to potential outliers and such points would still be worth of inspection.

## - Principal Component Analysis

The idea of Principal Component Analysis is to reduce the amount of variables used in an application by focussing on a few variables that at best sum up the information content of the original variables.

This solution might be especially useful to deal with the following issues:

1. Visualizing a larger numbers of variables and make qualitative inferences based on the results.
2. To reduce computational bottlenecks.  
Sometimes processing a very large amount of variables might be computationally infeasible or too expensive.  
↳ Reducing the amount of data to be processed while keeping the most information content might be a viable solution in such case.
3. To deal with collinearity that poses a major threat in many statistical techniques.

The question is now the following:

"Is it possible to achieve to reduce the amount of variables used while keeping the most information?"

The solution proposed through PCA is the following:

Create variables  $y_1, \dots, y_m$  being linear combinations of the original variables, such that the variance of  $x_1, \dots, x_m$  is maximal subject to

$$y_1 \perp y_2 \perp \dots \perp y_m$$

Variance of  $\hat{Y}_1, \dots, \hat{Y}_m$  is maximal  
subject to

$$\hat{Y}_1 \perp \hat{Y}_2 \perp \dots \perp \hat{Y}_m$$

Mathematically:

$$Y_1 = a_1 X_1 + a_2 X_2 + \dots + a_q X_q$$

⋮

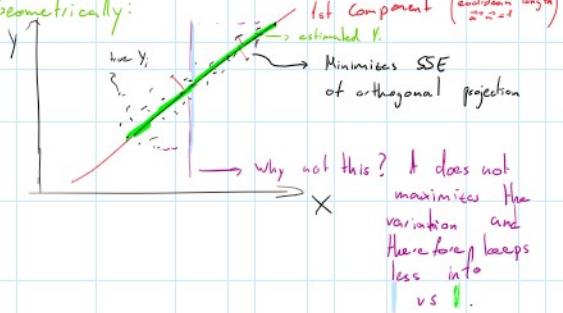
$$Y_m = a_1 X_1 + a_2 X_2 + \dots + a_q X_q$$

or in Matrix notation

$$\vec{Y} = \vec{a}^T \vec{X}, \text{ where}$$

$$\vec{a}^T = \arg \max_{\vec{a}} \vec{a}^T \Sigma \vec{a} \quad \text{s.t. } \vec{a}^T \vec{X} \vec{a} = 0 \quad \forall i \neq j.$$

Geometrically:



It is now clear that the above problem is a convex problem due to the properties of the variance matrix  $\Sigma$ . It is moreover under-specified, as without a further constraint on the  $a_i$ , you might always choose the coefficients to be as high as possible. Therefore we set the further constraint of the euclidean distance of the vector being 1.

We have therefore the final mathematical problem:

$$a_i = \arg \max_{\vec{a}} \vec{a}^T \Sigma \vec{a};$$

$$\text{s.t. } \vec{a}^T \vec{a} = 1 \\ \vec{a}^T \vec{a}_j = 0 \quad \forall i \neq j$$

This can be solved via Lagrange:

$$\begin{aligned} \frac{\partial}{\partial a_i} \vec{a}^T \Sigma \vec{a} - \lambda_1 (\vec{a}^T \vec{a} - 1) - \lambda_2 \vec{a}^T \vec{a}_j &= 0 \\ 2 \Sigma \vec{a}_i - \cancel{\lambda_1 \vec{a}_i} - \cancel{\lambda_2 \vec{a}_j} &= 0 \quad (\text{because of (iii)}) \\ \dots \quad \lambda_1 \vec{a}_i &= 1 \quad (\text{ii}) \\ \dots \quad \lambda_2 \vec{a}_j &= 0 \quad (\text{iii}) \end{aligned}$$

So that  $\Sigma \vec{a}_i = \lambda_1 \vec{a}_i$ , where it

is easy to see that this equation is a simple eigenvector equation.

The question on the choice of the eigenvector  $\vec{a}_i$  associated with the eigenvalue  $\lambda_1$  is then on which eigenvalue to choose.

The solution is given by transforming the above in our target quantity

$$\vec{a}^T \Sigma \vec{a} = \vec{a}^T \lambda_1 \vec{a}_i = \lambda_1 \underbrace{\vec{a}^T \vec{a}_i}_{=1 \text{ by def}} = \lambda_1$$

so that we choose the eigenvector  $\vec{a}_i$  that corresponds to the highest eigenvalue.

so that we choose the eigenvector  $\vec{x}$ ,  
that corresponds to the highest eigenvalue.

### Computing the PCA loadings.

There are two possible ways to compute the loadings. Both rely on getting the Eigenvectors of the variance matrix  $\Sigma$ .

This can be done by leveraging:

- (i) Eigen decomposition.
- (ii) Singular Value Decomposition.

### Eigen decomposition:

We are looking for a solution to the system:

$$A\vec{x} = \lambda\vec{x}$$

with  $\vec{x} \neq \vec{0}$ .

In order to get a solution we can therefore solve:

$$(A - \lambda I)\vec{x} = \vec{0}$$

The above is a homogeneous system and we therefore look for a solution

$$(A - \lambda I) = \vec{0}$$

Given now the fact that the determinant of a matrix equals 0; its columns are linearly dependent such that a non-trivial solution to a homogeneous system exists.

We therefore look in the above case for a solution

$$\det(A - \lambda I) = 0$$

Characteristic  
Polynomial

Solving the characteristic polynomial above you will get the desired eigenvalues and using the Eigenvector system:

$$A\vec{x} = \lambda\vec{x}$$

you can then insert your eigenvalues and get the desired eigenvectors.

Finally from the above it also holds that for a symmetric matrix  $A$  ( $n \times n$ ) it holds that with

$$\Lambda = \begin{pmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix}$$

$$T = (\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n)$$

$$A \cdot T = T \cdot \Lambda$$

$$A = T \cdot \Lambda \cdot T^{-1}$$

And it furthermore holds for squared matrices A:

$$\det A = \det(T) \cdot \det(T^{-1}) \cdot \det(\Lambda)$$

$$= \det(\Lambda) = \lambda_1 \cdot \lambda_2 \cdots \lambda_n$$

You can use the above to check the correctness of your computation.

Notice that the variance matrix is symmetric and thus the application of the above is straightforward.

As a final note it is possible to show that for the above case

$$\text{trace}(A) = \sum_{j=1}^n \lambda_j$$

In our case it holds

$$\text{trace}(\Sigma) = \sum_{j=1}^N \lambda_j$$

Moreover as  $T^{-1} = T^T$  as the eigenvectors are orthogonal in the case of symmetric matrices,

finally as it was proven before

$$Y_1 = \vec{a}_1 \vec{X}$$

so that As  $\vec{X}$  mean centered in PCA

$$\text{Var}(Y_1) = \vec{a}_1^T \vec{X}^T \vec{X} \vec{a}_1$$

$$= \vec{a}_1^T \sum \vec{a}_1 = \lambda_{\vec{a}_1 \vec{a}_1}$$

And therefore as due to orthogonality of the Principal Components

$$Y_1 \perp Y_2 \perp Y_3 \text{ etc.}$$

we have

$$\underbrace{\text{Var}(Y_1 + Y_2 + \dots + Y_q)}_{\Sigma} = \text{Var}(Y_1) + \text{Var}(Y_2) + \dots + \text{Var}(Y_q) = \lambda_1 + \dots + \lambda_q$$

It is then straight forward to appreciate that:

$$\frac{\sum \lambda_j}{\sum \lambda_j} \quad \left\{ \begin{array}{l} \text{Share of total variance} \\ \text{explained by the first } m \text{ components} \end{array} \right.$$

Finally note that if the observations fit exactly in dimensions it would be indicated by the presence of q-m 0 Eigenvalues.

On the scale of the variables:

An important thing to notice is that PCA is not scale invariant. Moreover, the sign of the Eigenvalues are arbitrary as:

$$a \Sigma \vec{x} = a \lambda \vec{x}$$

It follows that the choice of the scale used has an important effect on the result.

If you have "different" variables, or very different "

the choice of the scale used has an important effect on the result.

If you have different variables on very different scales the one with the higher scale will likely display the highest variance and dominate the first principal component.

↳ This is why it is usually suggested that the principal components should only be extracted from the sample covariance matrix when all of the original variables are roughly on the same scale.

A solution in this sense is to extract the principal components directly from the correlation matrix. This would then be equivalent to extracting the PC after each variable of the covariance matrix has been standardized to unit variance.

Finally note that there is no correspondence between the PC of  $\Sigma$  and the correlation matrix  $R$ . Hence using the latter is indeed an arbitrary decision where you give equal importance to all of the values.

### Rules of thumb for number of PC selection:

When you reduce the dimensionality you keep  $m \leq q$  PC that explain the most of the data.

The question is now how big  $m$  should be. There are three rules of thumb for selecting them:

1.  $\frac{\sum \lambda_j}{\sum \lambda_i} \geq 0.8$ ; i.e. share of variance explained bigger equal than 80%.

2.  $\lambda_j: \lambda_j \geq \frac{1}{q} \sum_j \lambda_j$ ; keep the variable if its variance is higher than the mean variance.

3. Plot  $\lambda_j$  vs  $j$ . Keep  $\lambda_j$  before you experience an elbow in the plot.

### Kernel PCA:

The PCA encountered so far imposed a linear relation between the PC and the origin variables.

Kernel PCA extends PCA by allowing to model a non-linear relation among the variables.

The basic idea is to create a richer set of variables  $\tilde{X}$  by applying a number  $\tilde{q}$  of non-linear transformations to the original variables  $X$ .

$$\tilde{X} = \{\phi_1(X), \phi_2(X), \dots, \phi_{\tilde{q}}(X)\}$$

Where you take  $\tilde{q} > q$ .

Then you continue as seen since here by computing the PC as:

$$Y_i = a_{i1}\phi_1(X) + a_{i2}\phi_2(X) + \dots + a_{i\tilde{q}}\phi_{\tilde{q}}(X)$$

with

$$Y_1 \perp Y_2 \perp Y_3 \perp \dots \perp Y_{\tilde{q}}$$

Albeit it would be theoretically possible to determine the PCA as above by following the schema:

Albeit it would be theoretically possible to determine the PCA as above by following the schema:

1. Transform  $X \in \mathbb{R}^n$  into  $\tilde{X} \in \mathbb{R}^q$  (and scale  $\tilde{X}$ ).
2. Calculate the sample covariance  $\hat{\Sigma}$  of  $\tilde{X}$
3. Calculate the spectral decomposition:

$$\hat{\Sigma} = \tilde{A}^T \Lambda \tilde{A}$$

In practice this is not done but you rather use the kernel trick.

Idea it can be shown that

$$\hat{\Sigma} \vec{x} = \tilde{\lambda} \vec{x}$$

is equivalent to solving

$$\tilde{K} \vec{v} = \tilde{\lambda} \vec{v}$$

with:

$$\begin{aligned}\tilde{K} &= \underline{\phi}(x)^T \underline{\phi}(y) \quad x, y \in \mathbb{R}^q \\ \underline{\phi}(x) &= (\phi_1(x), \phi_2(x), \dots, \phi_q(x))\end{aligned}$$

it is then possible to backtransform

$$\vec{v}, \tilde{\lambda} \rightarrow \tilde{\lambda}, \vec{x}.$$

The benefit of the approach lies in the computational benefit of working with

$$K = n \times q \cdot q \times n = n \times n \quad \text{instead of}$$

$$\Sigma = q \times q$$



and therefore you can work with extremely high dimensional transformations in a very high  $\mathbb{R}^q$ .

## Multidimensional Scaling

The idea of Multidimensional Scaling consists as well on representing points in hyperspaces of high dimension into a lower one.

Here the goal is the one of getting a representative plot in 2-3 dimensions. Then by inspecting the plot you can get meaningful insight in the association among the variables in the hyperspace.

About there might be a link relation with PCA as we will see soon, the approach is different.

While PCA acts directly on the multivariate matrix  $X$ , the idea of multidimensional scaling is the one of working on the information contained in the distances among the  $n$  data points in the hyperspace. Such information is summed up in proximity matrices.

The idea is then the one to find the best set of coordinates in the lower dimensional space  $\mathbb{R}^m$  such that the distances among the  $n$  points respects at best the original hyperspace distances.

Interesting is moreover the case of Non-Metric Multidimensional Scaling. This allows, in contrast to PCA, to deal with a spatial low-dimensional representation of ordinal data.

interesting is moreover the case of Non-Metric Multidimensional Scaling. This allows, in contrast to PCA, to deal with a spatial low-dimensional representation of ordinal data as we will shortly see.

### Classical Multidimensional Scaling - The Mathematics

The objective of MDS is to determine both the dimensionality  $m$  of the lower dimensional space  $\mathbb{R}^m$  as well as the  $m$ -dimensional coordinates  $x_1, \dots, x_n$  so that the model gives a good fit for the observed proximities.

In classical multidimensional scaling the distance used to register in the proximity matrix is the euclidean distance.

The goal here is then to create a representation

$$X \in \mathbb{R}^m \quad \text{for the } n \text{-data points from the proximity matrix} \quad D \in \mathbb{R}^{n \times n} \quad \begin{matrix} \text{dim orig.} \\ \text{hyperspace} \end{matrix}$$

This is possible as we will shortly prove that it is possible to express:

$$B := X^T X \quad \text{in terms of } D.$$

Given that it is possible to synthesize  $B$  from  $D$  and given  $B$  use spectral decomposition to get

$$B = A^T \Lambda A \quad \text{so that}$$

$$X = A \Lambda^{1/2}.$$

Intuition for expressing  $B$  in terms of  $D$ :

This will be shown in the case of  $\mathbb{R}^2$ . The generalization of it is straight.

Assume a matrix

$$X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix} \quad \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix}$$

and that

$$\begin{aligned} B &= X^T X \\ &= \begin{pmatrix} x_{11}^2 + x_{21}^2 & x_{11}x_{21} + x_{12}x_{22} \\ x_{21}x_{11} + x_{22}x_{12} & x_{21}^2 + x_{22}^2 \end{pmatrix} \end{aligned}$$

Moreover for the distance among  $X$  row vectors, i.e. data points, we have

$$\begin{aligned} d_{12}^2 &= (x_{11} - x_{21})^2 + (x_{12} - x_{22})^2 \\ &= x_{11}^2 + x_{21}^2 - 2x_{11}x_{21} + x_{12}^2 + x_{22}^2 \\ &\quad - 2x_{12}x_{22} \\ &= b_{11} + b_{22} - 2b_{12} \end{aligned}$$

where

$b_{ij}$  comes from the  $i,j$  entry in the matrix  $B$  above so that

$$b_{ij} = x_{ii}x_{jj} + x_{ij}x_{ji} \quad \forall i, j \in \{1, \dots, n\}$$

$$b_{ij} = x_{ii}x_{jk} + x_{ij}x_{jj} \quad \forall i, j \in \{1, \dots, n\}$$

$$b_{ij} = \sum_{k=1}^n x_{ik}x_{jk}$$

Now as the location of your plot does not change the distance among the data we can set the location restriction:

$$\bar{x}_{ij} = 0 \quad \text{so that}$$

$$\sum_{k=1}^n x_{ik} = 0$$

And consequently:

$$\sum_{i=1}^n b_{ij} = \sum_{k=1}^n \sum_{i=1}^n x_{ik}x_{jk} = 0 \quad \forall i$$

$$\sum_{i=1}^n b_{ij} = 0$$

$$\text{So that } \sum_{i=1}^n b_{ij} = \sum_{i=1}^n b_{ii} + b_{jj} - 2b_{ij}$$

$$= \sum_{i=1}^n b_{ii} + q b_{jj} \quad \stackrel{i \neq j}{\Rightarrow} b_{ij} = \frac{1}{n} \sum_{j=1}^n b_{ij} - \frac{\text{Trace}(B)}{n}$$

$$\Rightarrow \text{Trace}(B) + q b_{jj}$$

$$\text{Similarly } \sum_{j=1}^n b_{ij} = \text{Trace}(B) + q b_{jj} \Leftrightarrow b_{ij} = \underbrace{\frac{1}{n} \sum_{j=1}^n b_{ij}}_{:= d_{ij}} - \frac{\text{Trace}(B)}{n}$$

$$\text{And } \sum_{i=1}^n \sum_{j=1}^n b_{ij} = \sum_{i=1}^n (\text{Trace}(B) + q b_{ii})$$

$$= 2n \text{Trace}(B) \quad \Leftrightarrow 2 \frac{\text{Trace}(B)}{n} = \underbrace{\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n b_{ij}}_{:= d..}$$

so that putting everything together we can rearrange

$$d_{ij}^2 = b_{ii} + b_{jj} - 2b_{ij}$$

$$\begin{aligned} b_{ij} &= \frac{1}{2} (d_{ij}^2 - b_{ii} - b_{jj}) \\ &= \frac{1}{2} (d_{ij}^2 - (d_{ij} - \frac{\text{Trace}(B)}{n}) - (d_{ii} - \frac{\text{Trace}(B)}{n})) \\ &= \frac{1}{2} (d_{ij}^2 - d_{ij} - d_{ii} + d..) \end{aligned}$$

Given this derivation we can obtain multidimensional scaling by:

Step 1: Derive B from the proximity matrix D

Step 2: Obtain the spectral decomposition

$$B = A^T \Lambda A$$

Step 3:

Set

$$X = A^T \Lambda^{1/2}$$

Notice that when D  $n \times q$  is of full rank then  $\Lambda$  will have  $n-q$  zero eigenvalues so that you can represent B by

$$B = A_1^T \Lambda_1^{-1} A_1$$

where

$A_1$  = eigenvectors associated to  $\lambda_1, \dots, \lambda_q$

$\Lambda_1$  = diagonal matrix with non-zero eigenvalues

non-zero eigenvalues

You can see how easily from the previous section that you can obtain a low-dimensional representation of your data points by using

$$X = A_n \Lambda_n^{1/2}$$

using the highest  $m$  Eigenvectors. Using all of the  $q$  Eigenvectors you will be able to recover the complete proximity Matrix  $D$ .

So far we have assumed the proximity matrix to represent euclidean distances. When this is the case  $D$  is positive definite and the result of Multidimensional Scaling corresponds to the one obtained by PCA.

However, when the proximity matrix is obtained using different distance metrics  $D$  might not be positive definite and you might get negative Eigenvalues and complex coordinates. If number of negative Eigenvalues is small a meaningful representation of the proximity matrix may still be possible using the Eigenvectors associated with the highest Eigenvalues.

Two rules of thumb in this sense are:

**Trace criterion:** Choose the number of coordinates so that the sum of the positive eigenvalues is approximately equal to the sum of *all* the eigenvalues.

**Magnitude criterion:** Accept as genuinely positive only those eigenvalues whose magnitude substantially exceeds that of the largest negative eigenvalue.

If the number of negative Eigenvalues is high some other methods than classical Multidimensional Scaling, such as non-metric Multidimensional Scaling might be advisable.

### Non-metric Scaling:

When you have ordinal data that express preference the unit in which such preferences are expressed does not matter and a proximity matrix that focuses on some absolute distance metric does not make sense.

In non-metric Multidimensional Scaling we focus therefore on some representation that preserves the ordinal order:

$$\delta_{ijk} < \delta_{iik} < \dots < \delta_{\frac{n(n-1)}{2}}$$

The idea is then that the lower dimensional coordinates of the data should not change when monotonically transforming the original data

$$f(\delta_{ij}), \text{ with } \frac{df(x)}{dx} > 0 \quad \forall x$$

This is achieved by finding a geometrical representation of your data points such that you minimize the following stress criterion:

$$S^2 = \frac{\sum_{i,j} (f(\delta_{ij}) - \tilde{\delta}_{ij})^2}{\sum_{i,j} \tilde{\delta}_{ij}^2} \quad \begin{cases} \text{controls for} \\ \text{the scale} \\ \text{of the distances.} \end{cases}$$

Where

$\delta_{ij}$  = represent the original distances among the ordered data

$f(\delta_{ij})$  = is a monotone transformation.

$\tilde{\delta}_{ij}$  = are the fitted distances that result from the multidimensional scaling.

Non-Metric MDS works then by finding

- (i) the monotonic transformation  $f(\cdot)$
- (ii) that minimizes the stress criterion above.

In order to do that an iterative procedure is used where:

- (i) it searches for  $f(\cdot)$  using monotonic regression
- (ii) it minimizes the stress criterion and goes back to (i).

One of the major flaws with the above is that you might get stuck in local minima and do not find global solutions.

As a rule of thumb you can use the following:

- $S = 20\%$ : poor
- $S = 10\%$ : fair
- $S = 5\%$ : good
- $S = 0$ : perfect

## Non-linear Techniques

### ISOMAPS

So far we have only addressed linear scaling techniques. In fact, both PCA and Multidimensional scaling aims at representing the data points by applying rotations - linear operators - so that the new geometrical representation respects some desired properties - for instance having the axes in the direction of maximum variance.

It is thus clear that PCA and classical MDS just apply linear transformations and plot the data into linear subspaces.

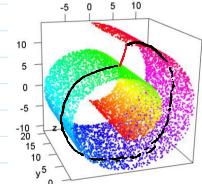
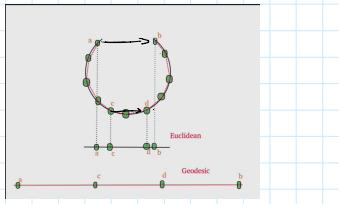
The question is now on how to accomodate some different structure. Our data might well not lie into a linear space. Some different structure among the data might hold and metrics computed on the notion of a linear space might be flawed.

This section addresses such a case. It considers the option of using **Manifolds** to represent different local structures among the data so that it will be possible then to apply different transformation based on the derived manifold structure with a different focus than the linear transformations so far encountered.

For a quick and dirty intro to manifolds check at: <https://bilqis.github.io/posts/manifolds/>

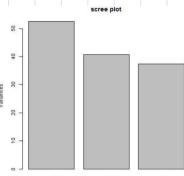
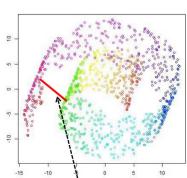
! Key good guy.

To further understand Isomap and to see the issues of relying on linear mappings when data lies on non-linear manifold consider the following examples:



- Geodesic Distance
- Euclidean distance.

In both of the examples linear distances are highly misleading. Moreover when considering the right hand plot it is easy to see that a linear method such as PCA will hardly achieve an effective variance decomposition of the data.



Do drop in elbow plot. Two dimensions do not suffice to explain the higher dimensional variability in the data.

This is where ISOMAP kicks in. This is a technique that leverages the fact that a manifold locally resembles an Euclidean Space such that "traditional" distance metrics can be computed.

ISOMAP works then by computing distance metrics between two points by approximating the global geodesic distance in the following way.

Step 0: Find neighbors of each point.  
This can be done:

- either by taking the K-nearest-neighbor
- or by taking all the points within radius  $\epsilon$ .

Step 1: Form a weighted graph, where

- each point is connected with its neighbours.
- The weights are given by the euclidean distances among the points.

Step 2: Compute the graph distance between two points as the shortest path (in terms of edges weights) between the two points.

In such a way you effectively approximated the geodesic distance of your -possibly unknown- manifold.

You can then compute above to obtain a full distance matrix and finally apply classical MDS on it.

In such a way ISOMAP is able to represent both the local and the global structure of your data and will therefore be able to retain the information of the two in lower dimensional representations.

both the local and the global structure of your data and will therefore be able to retain the information of the two in lower dimensional representations.

Important is moreover to stress the following when working with Isomaps:

- You must work on a convex space for extracting the local distances among neighbours. Otherwise you might jump over holes and short-circuit the true distance.



- The choice of the K-neighbours or the radius  $\epsilon$  is of first order importance.
  - ↳ Too big  $K/\epsilon$  will result in short circuits
  - ↳ Too small  $K/\epsilon$  will lead to zig-zag roads considering the discrete amount of data
- Short-circuits might occur even if the manifold is too curvy or if the collected data-points are noisy.

## T-SNE

This is a recent model for visualizing high-dimensional data that explicitly sets a high importance in preserving the local structure of the data. It ultimately allows for lower-dimensional plots preserving both the local structure of the data as well as the global one.

The method achieves that by the following scheme:

- Compute the probability function for a data point  $x_i$  that states the probability of observing another one  $x_j$  assuming a gaussian kernel around  $x_i$ :

$$f(x_j) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_j - x_i)^2}{2\sigma^2}}$$

- Compute the conditional probability:

$$p_{ji} = \frac{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_j - x_i)^2}{2\sigma^2}}}{\sum_k \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_k - x_i)^2}{2\sigma^2}}}$$

Stating the probability of observing  $x_j$  as a neighbour of  $x_i$  among all of the  $x_k$  given the gaussian distribution of  $x_k$  centered around  $x_i$ .

- Compute the similarity metric:

$$p_{ij} := \frac{(p_{ji}) + (p_{ij})}{2N} \quad \text{Not symmetric}$$

This will then be compared to the lower dimensional similarity metric:

$$q_{ij} := \frac{f(\|y_b - y_i\|)}{\sum_{k \neq i} f(\|y_b - y_k\|)} \quad \text{with } \underbrace{f(z) = 1 - \dim + \text{dist.}}_{\text{this for mitigating.}}$$

$$q_{ij} = \frac{1 \|y_i - y_j\|}{\sum_{k \neq i} f(\|y_k - y_i\|)}$$

with  $f(z) = 1 - \text{dim } + \text{dist.}$

this for mitigating  
the crowding problem  
in lower dimensions where  
you would get all of the  
points to lie close together  
without the occurrence of  
clusters.

The comparison is done via the Kullback-Leibler information criteria  
that would ultimately account for the preservation of  
both the local and the global structure of the data.

$$KL(P||Q) = \sum_{i,j} p_{ij} \log \left( \frac{p_{ij}}{q_{ij}} \right)$$

↑  
when the neighbors  
are equally likely  
selected in the higher and  
lower dimension  $KL=0$ .

When different, You scale  
according to  $p_{ij}$  such that  
when

$$\|x_j - x_i\|^2 \gg 0 \Rightarrow p_{ij} \text{ is small}$$

(think of the gaussian  
kernel).

and if

$$\|x_j - x_i\|^2 \sim 0 \rightarrow p_{ij} \text{ is big so}$$

that you give  
higher weight to  
the  $\log \left( \frac{p_{ij}}{q_{ij}} \right)$  if the  
local structure is not  
preserved.

T-SNE works then by finding the lower dimensional  
representation that minimizes such Kullback-Leibler information  
criteria.

The question on how to set the  $\sigma^2$  of the Gaussian  
kernel remains open. Of course setting a very  
high  $\sigma^2$  will set a bigger weight on preserving  
the **global structure** while a smaller  $\sigma^2$  will set higher  
weight on maintaining the **local structure**.

In order to set such important parameter the authors of  
the method suggest to choose it indirectly by the choice of  
the perplexity measure.

The perplexity measure is an information measure that  
can be interpreted roughly as:

"This thing is right about as often  
as an  $x$ -sided dice, where  
 $x$  is the perplexity measure, and  
where right as a 6-sided dice means  
right  $1/6$  of the time."

Mathematically it can be written as:

$$H(p_{xx}) = - \sum p_{xx} \log p_{xx}$$

i.e. a geometric average of the inverses of the  
probabilities.

Once you choose the perplexity parameter you choose  $\sigma$   
so that this is reached. The authors of the papers that  
presented the method advise on  $H(p_{xx}) \in [5, 50]$ .

# Clustering

The goal of clustering is the one of creating groups of observations such that

- Elements within the cluster are similar
- Elements among different clusters are different.

This notes will treat:

- ① Agglomerative Clustering: Distance-Based Aggregation
- ② Divisive Clustering (Partitioning): K-Nearest Neighbours
- ③ Model Based: Finite Mixture Models: Gaussian Mixture Models
- ④ Density Based Clustering

Before going into the nitty gritty, a brief comparison of the methods:

## ① Agglomerative Clustering:

- Idea: Start with  $N = k$  clusters. Iteratively merge the two nearest clusters based on some distance metric.

## ② Divisive Clustering:

- Idea: Given an a priori defined number of clusters, find the best topological position for the cluster centers such that some small cluster distance metric will be minimized in relation to such centers.

## ③ Model Based Clusters

- Idea: Data are generated from different probability distributions. The number of probability distributions as well as their means/parameters are defined minimizing some information criteria measure such as the Bayesian Information Criteria

## ④ Density Based Clustering:

here the assumption is that clusters are characterized by having a high amount of data density in the cluster region. And no or little data within clusters regions. You find then the clusters by specifying two hyperparameters that will be discussed later by a heuristic argument that is as intuitive as elegant.

### ①.1 Agglomerative Clustering:

This algorithm here is rather simple and the details are omitted.

Important is however to understand that the choice of the distance metric by which the closest clusters are merged.

Some of the possible choices are

#### ①.1.1 Single linkage $\Leftrightarrow$

minimal distance among two datapoints pairs in different clusters.

$$d_{B1} = \min_{i \in [3]} d_{ij}$$

$\downarrow$  two clusters

#### ①.1.2 Complete linkage:

minimal distance among two datapoints pairs coming from the two different clusters

$$d_{AB} = \max_{i \in A, j \in B} d_{ij}$$

(3) Average-linkage  $\Leftrightarrow$  Average among all of the datapair's distances among the two datasets



$$d_{AB} = \frac{1}{N+M} \sum_{i,j}^{N+M} d_{ij}$$

(4) Ward's Method:  $\Leftrightarrow$  works by leveraging the notion of within group sum of squares defined as follows:

The Ward's Method defines the distance between two clusters as follows

$$d_{AB} = \frac{WGSS_{AB}}{(WGSS_A + WGSS_B)}$$

$$WGSS = \sum_{i \in C} \|x_i - \bar{x}_C\|^2$$

$\underbrace{\text{simply the sum of the } L^2 \text{ distance for the points}}$   
within a cluster.

It is clear that different methods are more or less robust against outliers and have in general bigger merits/ flaws. Notice, that the distance/dissimilarity metric can be chosen as one of the many and has not to be confined to the usual euclidean one.

## ② Divisive Clustering: K-MV.

This method consists in finding the cluster centers that minimizes the WGSS given the fixed amount of chosen clusters.

To find the clusters centroids is computationally infeasible. Hence, an approximation algorithm is used: the Lloyd's algorithm.

This works as follows:

- ① Start with  $k$  random centroids for the  $k$  clusters
- ② Assign each datapoint to the center, that is closest given some distance metric
- ③ Recompute a new center for the cluster as the average of the observations assigned to the cluster in ②.
- ④ Stop if the cluster mean did not change much.  
Go back to ② otherwise.

As there is no guarantee that the Lloyd's algorithm will reach the global optimum it is usual to compute the KNN multiple times and observe how the results change.

The question about the cluster centroids remain open. How should such center be determined? There are essentially two options:

- ① Means: Here the center of the cluster does not have to be an observed datapoint. It is simply the average of all points.
- ② Medoid: Here the center is defined as the observation for which the distance/dissimilarity to all other datapoints is minimized. This is more robust toward outliers. However, it is computationally more demanding.

### (3) Model Based Clustering: Finite Mixture Models

So far clustering was based on heuristics based on some distance metrics.

An other option consists on specifying a full statistical distribution for our data.

The idea as mentioned above is to specify different PDFs and assign the different datapoints to the distribution they belong to with the highest posterior probability.

The idea works as follows using different gaussian distributions with different parameterization to express the different clusters and to express the probability of observing a datapoint as the mixture of different clusters.

$$f(x) = \sum_{j=1}^K p_j g(x | \theta_j)$$

↑ gaussian finite mixture with  $\sum p_j = 1$

Number of classes/clusters  
 gaussian probability of observing  $x$  given cluster  $j$  parameterization  $\theta_j$

Given the above setting it is clear that once the number of clusters / mixture distributions  $K$  and their parameterization has been defined, you can classify observation  $x_i$  by means of bayesian arguments as:

$$\{x_i \in \theta_j : P(x \in \theta_j | x_i) = \frac{p_j \cdot g(x_i | \theta_j)}{\sum_j p_j \cdot g(x_i | \theta_j)} > P(x \in \theta_l | x_i)\}$$

posterior of observing  $x$  given  $x_i$

The central question now is on how to define the number of clusters  $K$  and the parameterization of the finite mixture distributions:

$$[p_j \quad \wedge \quad \theta_j]$$

This is achieved by means of the Expectation-Maximization algorithm. To understand that watch the two videos following

(ML 16.3) Expectation-Maximization (EM) algorithm

```

Given:  $x = (x_1, \dots, x_n)$  sample sap for  $X = (X_1, \dots, X_n)$ 
Model  $(X, \theta) \sim f_\theta$  for some (unknown)  $\theta \in \Theta$ .
Goal:  $\theta^* \in \arg \max_\theta f_\theta(\theta)$ 
Info:  $f_\theta(x) = \prod_i f_\theta(x_i)$  > want to maximize.
Alg: Initialize  $\theta_0 \in \Theta$ .
For  $t = 0, 1, 2, \dots$ :
  + E-step:  $Q(\theta, \theta_t) = \dots$ 

```



There you can get as well the intuition as well as the details for the algorithm.

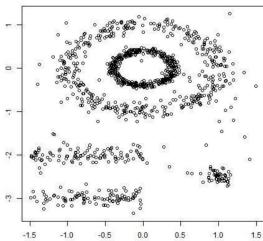
Moreover to make the optimization problem more tractable it is common practice to impose restrictions on the covariance structure of the different gaussian distributions. Finally to bound the solution the parameters are chosen minimizing some information criteria as a pure maximization of

$$f(x) = \sum_{j=1}^K p_j g(x | \theta_j)$$

will always increase in the number of selected clusters  $K$ .

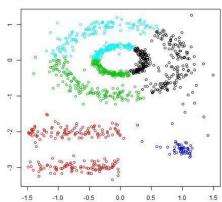
### (4) Density Based Clustering

This was introduced in the slides of the lecture in a very interesting way.  
Consider the following non-convex clusters scenes.

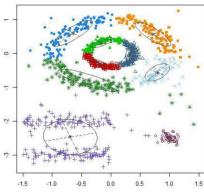


Here the few clusters can be easily recognized by humans through eyeballing. However, the mathematical constructs so far seen fail brutally when working on non-convex structures:

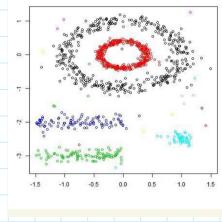
K-means



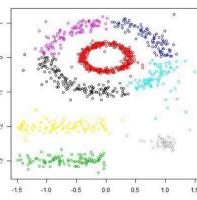
Gaussian Mixture



Agglomerative (Single Linkage)



Agglomerative (Average Linkage)



DBSCAN - density based spatial clustering of applications with noise - works now as defining clusters as gaps of high-density points.

↳ Different regions with a high density of points will then be classified as different clusters.

The idea works as follows:

- Define a radius  $\epsilon$

- Define  $M \in \mathbb{N} \Leftrightarrow$  minimum number of points necessary for the definition of a core point.

Given the two parameters you can then label three types of different points:

1. Core Points
2. Border Points
3. Noise Points

These three categories are defined as follows: for each point check at the # of points in their neighborhood.

Then:

- If the number of points in the neighborhood is  $\geq M$  then classify it as a core point.

Given then your core points:

is  $\geq M$  then classify it as a core point.

Given then your core points:

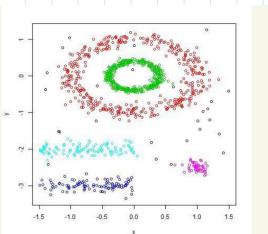
- Classify points in the neighborhood of core points that are themselves not classified as core points as border points.

Classify points that are neither core nor border points as noise-points. These might be then further analyzed by looking whether they are outliers or not.

Finally you would classify a cluster as follows:

- Connect all the core points that lie within distance  $\epsilon$  with an edge to each other. The connected core points will form a cluster.
- Assign border points to one of the clusters they belong to.

The result for DBSCAN would look as follows in the example before:



The open question that stays is on how to deal with the hyperparameter selection.

If:

- $\epsilon$  too big and  $M$  too small then all of the points are labeled as core-points.
- if  $\epsilon$  too small and  $M$  too big all noise points

Some way to select data points:

Possible strategy for finding tuning parameters: look at the  $k$ -nearest neighbor distance ( $k$ -NN distance), the distance from a point to its  $k$ -nearest neighbor.

1. Pick an initial  $k$ . Rule of thumb  $k \geq d + 1$ , where  $d$  is the dimension of the data. For  $d = 2$ ,  $k = 4$  often works well.
2. Compute the  $k$ -NN distance for each point.
3. Sort them in increasing order and plot them.
4. Look for sharp increase, and choose this distance as  $\epsilon$  and the  $k$  as  $MinPts$ . If there is no "elbow", change  $k$ .

There are then a few ways to interpret the clusters and one quantitative metric to analyze its results:

To interpret:

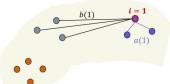
1. Analyze the centroids or representative center of the clusters (this especially for method (2)). What are the characteristics of the data?
2. Plot in lower dimension via PCA. Look at the linear combination of the axis. What do they represent? Where do the clusters lie w.r.t. these representative axis?

cluster lie 1 w.r.t. these representative axis!

To quantify the goodness of the clustering:

Consider the silhouette  $S(i) \in [-1; 1]$ .  
This is defined as:

$$S(i) = \frac{b(i) - a(i)}{\max(b(i), a(i))}$$



where:

$b(i)$  = minimal average distance between point  $i$  and all of the points belonging to another cluster.

$a(i)$  = average distance of point  $i$  and all of the other points belonging to the cluster.

It is now clear that for a high silhouette the clustering performed well for the point. For silhouette  $\approx 0$  the clustering was ok. For silhouette  $\approx -1$  the point was likely badly classified.

## Explanatory Factor Analysis

This technique deals with the situation where we observe some manifest variables and based on that we try to make inferences on an unknown latent variable.

The situation can be formulated mathematically as:

$$\begin{aligned} X_1 &= \lambda_{11}f_1 + \dots + \lambda_{1k}f_k + u_1 \\ X_2 &= \lambda_{21}f_1 + \dots + \lambda_{2k}f_k + u_2 \\ \dots \\ X_q &= \lambda_{q1}f_1 + \dots + \lambda_{qk}f_k + u_q \end{aligned}$$

or in Matrix notation:

$$\vec{X} = \Lambda \vec{f} + \vec{u}$$

where the major difference lies on the fact that  $\vec{f}$  represent the latent unobserved factors so that it is not possible to estimate the factor loadings  $= \Lambda$  by means of LS.

It is now clear that without posing some structure on the relation above it will be hardly possible to make inference on the factors and their loadings from the manifest variables.

Such structure can be posed in the following way:

- $E(u) = 0$   $\text{Cov}(u) = \Psi$ , with  $\Psi$  being a diagonal matrix, meaning that the error terms are uncorrelated.
- $\text{Cov}(\vec{u}, \vec{f}) = 0_{q,k}$  being the zero matrix such that the factors and the error terms are uncorrelated.
- $E(\vec{X}) = \vec{0}$ , mean centered manifest variables
- $E(\vec{f}) = 0$ ,  $\text{Cov}(\vec{f}) = I$ , so that factors are standardized variables.

It follows now that given such structures there is a 1:1 relation among the manifest variables covariance structure and the factor loading. as:

$$\text{Cov}(X_1, \dots, X_q) = \Sigma$$

$$\begin{aligned}\Sigma &= \text{Cov}(\Lambda \tilde{F} + \psi) \\ &= \text{Cov}(\Lambda \tilde{F}) + \text{Cov}(\psi) \quad \begin{array}{l} \text{(as } \tilde{F} \text{ and } \psi \text{ uncorr.)} \\ \downarrow \end{array} \\ &= \Lambda^T \text{Cov}(\tilde{F}) \Lambda + \psi \\ &= \Lambda^T \Lambda + \psi\end{aligned}$$

So that

$$\sigma_j^2 = \sum_{i=1}^k \lambda_{ji}^2 + \psi_j$$

and it is possible to see that the variance is decomposed into

an a

- common term  $h_i^2 = \sum_{j=1}^k \lambda_{ji}^2$

and a

- specific term  $\psi_j$

The question consists now on obtaining some estimates for the factor loadings above by leveraging the above decomposition, where in order to do so we have to fix first:

- the number of latent factors  $k$ .

In order to do that we will explore two different methods:

1. Principal Factor Analysis.

2. Maximum Likelihood Estimation.

### 1. Principal Factor Analysis:

This is based on an iterative algorithm to determine the factor loadings and the variable specific variance  $\psi$

The issue is namely the following: given the determined  $K$  based on the observation  $\Sigma$  you will need to be able to discern

$$\Lambda^T \Lambda \text{ and } \psi$$

Now given one of  $\Lambda$  and  $\psi$  you could determine the other. However in the specific situation you do not know neither of the two.

The idea of Principal Factor Analysis is then to mediate the issue by:

1. making an educated guess of  $\hat{\psi}$
2. determine  $\Lambda$  via

$$\Lambda^T \Lambda = \Sigma - \hat{\psi}$$

i.e. by finding the spectral decomposition

$$\underbrace{\Sigma - \hat{\psi}}_{\text{symmetric}} = A^T \Phi A$$

And set  $\Lambda = A^T \phi^{1/2}$

3. Compute  $\psi$
4. Iterate (2)-(3.) until convergence.

The educated guess on  $\phi$  is then based on the following reasoning:

- First the Factor Analysis is scale invariant as:

$Y = CX$  where  $C$  is a diagonal matrix representing the loadings and  $X$  is the specific variance. It is then true that:

$$\begin{aligned} \text{Var}(Y) &= \text{Var}(CX) \\ &= C^T \text{Var}(X) C \\ &= C^T (\Lambda^T \Lambda + \psi) C \\ &= \Lambda^T \Lambda C + \psi C \\ &= \tilde{\Lambda}^T \tilde{\Lambda} + \tilde{\psi} C \end{aligned}$$

with:  $\tilde{\Lambda} = \Lambda C$   
 $\tilde{\psi} = \psi C$

So that applying the factor model to  $X$  and then scaling the loadings and the specific variance is the same as scaling  $X$  first and then applying the factor model to  $Y$ .

It is then possible to see that instead of the covariance matrix  $\Sigma$  of the manifest data you can use the correlation matrix  $R$  thereby having

$$R = \tilde{\Lambda}^T \tilde{\Lambda} + \tilde{\psi}$$

$$\tilde{\psi} = 1 - \sum_{j=1}^K \tilde{\lambda}_{jj}^2$$

so that you can obtain the estimates of  $\tilde{\psi}$ , compute  $\tilde{\Lambda}$  and finally scale the two back; where the initial estimates  $\tilde{\psi}$  are

$$\tilde{\psi} = 1 - \sum_{j=1}^K \tilde{\lambda}_{jj}^2, \text{ where}$$

$$\sum_{j=1}^K \tilde{\lambda}_{jj}^2 = \left\{ \begin{array}{l} \text{squared of correlation coefficient of} \\ \text{manifest variable } j \text{th with all the other manifest variables} \\ \text{or} \\ \text{largest correlation coeff of the manifest variable} \\ \text{and the other variables} \end{array} \right.$$

## 2. Maximum Likelihood Estimator

Another possibility when estimating the factor loadings is the one of relying on MLE and its asymptotic properties.

The idea works as follows:

- Assume a multivariate normal distribution of the manifest data

Then under this assumption and the factor analysis model assumptions you can derive the likelihood:

of the manifest data

Then under this assumption and the factor analysis model assumptions you can derive the likelihood:

$$L(\boldsymbol{\psi}, \Lambda) = -\frac{N}{2} (\log(2\pi\Sigma) + \text{trace}(\Sigma^{-1}\Lambda^T\Sigma))$$

It is then possible to find the arguments  $\boldsymbol{\psi}$ ,  $\Lambda$  that maximize the above.

Notice that a sensible advantage of using MLE for estimating the factor loadings is that you can test for the hypothesis of the goodness of fit of  $k$  factors that result by comparing the restricted and unrestricted likelihoods, which follow a chi-squared distribution. When you persistently fail to reject some restricted version with  $k$  being increased at each run, then there is evidence that the model fails to represent the data or that your sample is too large.

Notice finally that once the factor loadings are obtained it is common practice to apply factor rotation to them to better interpretability of the factors.

Rotation does not in fact alter the structure of factor analysis, which is in fact invariant to such transformations.

To see that consider the rotation  $M$ , so that  $M^T M = I$ . It follows that

$$f = M^T f \quad \tilde{\Lambda} = \Lambda M$$

so that:

$$X = \Lambda M M^T f + \tilde{\nu} = \tilde{\Lambda} f + \tilde{\nu}$$

and

$$\Sigma = \tilde{\Lambda} M M^T \tilde{\Lambda}^T + \Psi = \tilde{\Lambda} \tilde{\Lambda}^T + \Psi$$

So practically this means that you have to set some restriction at first to get a solution; you can read about in the book - the first is however to set the restriction

$$G = \tilde{\Lambda} \tilde{\Lambda}^T - \tilde{\Lambda}$$

to be diagonal with its element sorted in decreasing order, leading to a factor structure where the first factor explains the most of the common variation of  $X$ , the second the most being uncorrelated to the first effect. It is a similar reasoning to the PCA.

In order to facilitate the interpretation of the factors obtained with the above restriction, two types of rotation are possible

- orthogonal rotation: which preserves the uncorrelated structure of the factors as above
- oblique rotation: frees up the uncorrelated factors restriction.

Given the two it is clear that to facilitate the interpretation we need to find rotations such that ideally:

- Each variable is loaded on at most one factor
- All factor loadings are either large and positive or close to zero.

Two examples of reaching that consists in:

- Varimax rotation (a type of orthogonal rotation; few large and many small loadings)
- Promax rotation (oblique rotation that results in a low correlation among factors).

Notice finally how in the case of orthogonal rotations, the factor loadings can be interpreted as the correlation coefficients in the case of factor analysis with scaled manifest variables.

Proof:

$$\begin{aligned}\text{Corr}(X, f) &= \text{Corr}(\tilde{\Lambda}f + u, f) \\ &= \text{Corr}(\tilde{\Lambda}f, f) + \text{Corr}(u, f) \\ &= \tilde{\Lambda} \cdot \text{Corr}(f, f)\end{aligned}$$

## Classification

The general problem of this chapter is of the following form; given a set of features  $X$  you have to classify a variable  $y$ .

There are two possibilities in such a case:

(I) Directly model  $\Pr(Y=j | X=x)$

(II) Use a soft bayesian approach modeling first  $\Pr(X=x | Y=y)$  and then using Bayes theorem to obtain  $\Pr(Y=j | X=x)$ . The soft adjective used before refers to the fact that you will not use a rigorous approach when defining a prior; in the sense that you will not use a non-informative prior.

### On the second approach: LDA and QDA

The general approach of the two techniques belonging to the method is then, essentially the following:

1. Assign a prior  $p_j$  that each variable belongs to a class
2. Assume that the feature variables are multinomial distributed conditional on observing the categories  $y$ ; i.e.  
 $X|Y_j \sim N(\mu_j, \Sigma_j) = g(X, \vec{\theta}_j)$  for all classes  $j$ .

Such that for the model above, the unconditional PDF of  $X$  is given by:

$$\Pr(X) = \sum_{j=1}^K p_j \cdot g(X, \vec{\theta}_j)$$

Unmarginalizing,

By applying Bayes Theorem it is now possible to see that

$$\begin{aligned}\Pr(Y=j | X=x) &= \frac{\Pr(X=x | Y=j) \cdot \Pr(Y=j)}{\Pr(X=x)} \\ &\propto \frac{g(x, \vec{\theta}_j) \cdot \Pr(Y=j)}{\sum_j p_j g(x, \vec{\theta}_j)}\end{aligned}$$

The classification is then chosen as

*Not dependent  
on  $j$ .*

$$\operatorname{argmax} \Pr(Y=j | X=x) = \operatorname{argmax} \log(g(x, \vec{\theta}_j)) + \log(p_j) - \log\left(\sum_j p_j g(x, \vec{\theta}_j)\right)$$

Given then the multivariate normal assumption it is clear that:

$$g(X, \beta_j) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_j|}} \exp(-\frac{1}{2} (X - \mu_j)^T \Sigma_j^{-1} (X - \mu_j))$$

so that it follows:

$$\underset{j}{\operatorname{argmax}} P_r(Y=j | X=x) = -\log(\varepsilon_j) - \frac{1}{2} ((x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j)) + \log(P_r(y=j))$$

This is in fact called quadratic discriminant analysis due to the quadratic term  $(x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j)$ .

A final question is on how to estimate the necessary parameters of the model:  $\varepsilon_j; P_r(y=j), \mu_j$ :

These are estimated as:

- $P_r(Y=j)$  = fractions of observations belonging to  $j$   
(recall that we are here in supervised learning domain).
- $\mu_j$  = sample mean  $\bar{x}_j$  of all observations belonging to  $j$
- $\Sigma_j$  = sample covariance  $S_j$

A problem that might arise with the Quadratic Discriminant Analysis is that if the feature space is particularly large the number of parameters to be estimated might be massive. Think for instance at the sample covariance  $\Sigma_j$  and how its size increase in the features dimension. In such case if the available dataset is not too large you might quickly run into issues.

To obviate such problem a more restrictive model is proposed - the **Linear Discriminant Analysis**.

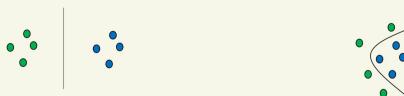
This consists of the same approach presented before with the only difference of restricting the covariance matrix to be equal across the classes  $j$ . So that  $\text{Var}(X) = \Sigma \forall j$ .

Inserting this in the expression above

$$\underset{j}{\operatorname{argmax}} P_r(Y=j | X=x) = \log(P_r(y=j)) + x^T \Sigma^{-1} \mu_j - \frac{1}{2} \mu_j^T \Sigma^{-1} \mu_j$$

Summing up:

LDA	QDA
+ Only few parameters to estimate; accurate estimates	- Many parameters to estimate; potentially less accurate estimates
- Less flexible (linear decision boundary)	+ More flexible (quadratic decision boundary)



On Naive Bayes: Note that this method is widely used in Machine Learning is nothing more than a typical very restrictive form of QDA, where you have:

$$\Sigma_j = \text{diag}(\sigma_1^2, \dots, \sigma_k^2), \text{ where } k = \# \text{ features}$$

so that Naive Bayes assumes uncorrelation among the different features.

On Method 1: Directly modeling  $P_r(Y=j)$  through logistic regression

As discussed the aim of the logistic regression is to model directly the quantity of interest

$$\Pr(Y=1|X=x)$$

This is done through the inverse-logit function. Recall the logic of it with two-categories case. The logit is a GLM, i.e. it tries to model:

$$g(E(y|x)) = \underbrace{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k}_{\text{link function}} \quad \underbrace{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k}_{\text{linear model}}$$

In the case of a two-classes case:

$$E(y|x) = \Pr(Y=1|x) \cdot 1 + 0 \cdot \Pr(Y=0|x)$$

Notice how that to model

$$\Pr(Y=1|x) = \beta_0 + \beta_1 x_1 + \dots \quad \text{does not make sense as you might well}$$

end up with  $\Pr < 0$  and  $\Pr > 1$ .

The solution is to use the link function:  $g(x) = \log\left(\frac{x}{1-x}\right)$

this allows to model the log-odds ratio in a linear way through the linear model

so that by rearranging:

$$\log\left(\frac{\Pr(Y=1|x)}{1 - \Pr(Y=1|x)}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

$\Leftrightarrow$

$$\begin{aligned} \frac{\Pr(Y=1|x)}{1 - \Pr(Y=1|x)} &= e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k} \\ \Pr(Y=1|x) &= \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k}} \quad \left. \begin{array}{l} \text{cannot be negative;} \\ \text{bounded to 1.} \end{array} \right\} \end{aligned}$$

It is now clear that you can assign the class according to the maximum  $\Pr(Y=1|x)$  so far obtained.

An interesting note moreover is the one of looking at differences between LDA and the logistic model. In fact it is straightforward to see that the underlying model of the two is essentially the same as:

LDA:

$$\log\left(\frac{\Pr(Y=1|x)}{\Pr(Y=0|x)}\right) = \underbrace{\log\left(\frac{p_1}{p_0}\right)}_{\beta_0} - \frac{1}{2} (\mu_0 + \mu_1)^T \Sigma^{-1} (\mu_1 - \mu_0) + \underbrace{x^T \Sigma^{-1} (\mu_1 - \mu_0)}_{\beta^T}$$

Logit:

$$\log\left(\frac{\Pr(Y=1|x)}{\Pr(Y=0|x)}\right) = \beta_0 + \beta^T X$$

However, despite the similar underlying model the LDA model estimates its parameter in a much less efficient way as it maximizes:

$$\prod_i p(x_i, y_i) = \prod_i \underbrace{p(x_i|y_i)}_{\text{gaussian}} \underbrace{p(y_i)}_{\text{bernoulli}}$$

while in the logit model

$$\prod_i p(x_i, y_i) = \underbrace{\prod_i p(y_i|x_i)}_{\text{Bernoulli-logit}} \cdot \underbrace{p(x)}_{\text{ignored}}$$

the conditional probability is maximized directly and involves therefore less assumptions staying more flexible.

## Fisher Discriminant Analysis

In contrast to the previously seen methods which rely on statistical methods, Fisher discriminant analysis is another important classification technique that relies on heuristic methods.

The idea is the one of creating a linear combination separating the classes to

Discriminant analysis is another important classification technique that relies on heuristic methods.

The idea is the one of creating a linear combination separating the classes to such a way that the Rayleigh quotient - i.e. the ratio among the between-class variance and the within class variance is maximized.

Such vector  $\vec{a}$  can be calculated:

$$\arg \max_{\vec{a}} \frac{BCV}{WCV} = \frac{\vec{a}^T B \vec{a}}{\vec{a}^T W \vec{a}},$$

where  $B$  = between-group sample variance

$W$  = within-group sample variance

It is then possible to show that such maximum is given by the eigenvector  $W^{-1}B$  corresponding to the largest eigenvalue.

## Regression Tree and Random Forrest

This chapter explores tree-based classification and regression models. The idea of the model is to partition the space according to the features characteristics in a data driven way, and according to the partition space fit a local regression or make a classification decision.

A major point drawing the line between plain vanilla regression trees and random forests is the way to which the two models deal with overfitting. As we will soon see, the plain vanilla regression/classification trees will use cross-validation for pruning the tree according to a cost-complexity criterion similar in spirit to the information criteria; in contrast random forests leverage the idea of MSE reduction by aggregating - ideally uncorrelated - trees with a reduction of the variance and an increase of the bias - i.e. a general reduction of overfitting.

### Regression And Classification Trees:

Both rely on finding the partitioned feature space that leads to the best classification and regression fit!

In order to define such best possible fit some quantitative metrics have to be defined. In the case of trees:

- Regression Tree:

→ you fit according to the local average of your predictor variable,  $y_i$ :

$$\hat{c}_m = \text{ave}(y_i | x_i \in R_m) = \frac{\sum_{i:x_i \in R_m} y_i}{N_m}$$

•  $N_m = \#\{x_i \in R_m\}$ : number of observations in  $R_m$

→ you estimate the goodness of the fit according to the well known residual sum of squares.

- Classification Trees

→ You classify according to the probability of each class in the end nodes.

Such probability is given by the share of each class in each node. The classification decision is then based on the most probable class.

- Proportion of class  $k$  observations in node  $m$   

$$\hat{p}_{mk} = \frac{1}{N_m} \sum_{i: x_i \in R_m} 1_{(y_i=k)}$$
- Majority vote: predict class  $k(m)$  for which  $\hat{p}_{mk}$  is largest, i.e.,  

$$k(m) = \operatorname{argmax}_k \hat{p}_{mk}.$$

~ You then classify the goodness of the fit according to a classification error measure. Some well known measures are:

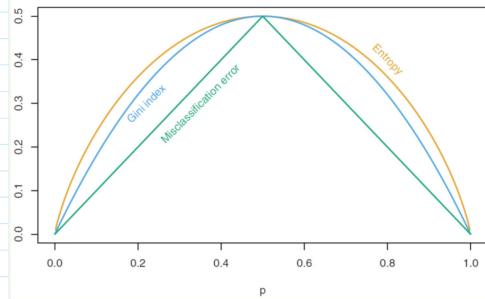
- Misclassification Error:

$$\text{share of misclassified} \left\{ \frac{1}{N_m} \sum_{i: x_i \in R_m} 1_{(y_i \neq k(m))} = 1 - \hat{p}_{mk}(m) \right.$$

- Gini Index:  $\overset{\text{GCE}}{\sim}$   

$$\sum_{k \neq k'} p_{mk} p_{mk'} = \sum_k p_{mk} (1 - p_{mk})$$
- Cross-entropy or deviance:  

$$-\sum_{k=1}^K \hat{p}_{mk} \log(p_{mk})$$



Given such characterization of the goodness of fit it is possible to understand the data driven partitioning of the feature space maximizing the fit:

1. Start from the trivial partition  $R = \{R\}$
2. For a variable / feature  $j$  create two partitions  
 $R_1(j,s) = \{x | x_j \leq s\}$  and  $R_2(j,s) = \{x | x_j > s\}$   
 Choose the splitting point such that the goodness of fit defined above is maximized
3. Iterate (2.) until a stopping criterion is reached - for instance minimal number of obs. in one node.

Important Note: Such algorithms are greedy in the sense that once a partition has been selected there is no modification of it deeper in the tree. This would be computationally infeasible.

Obviously an algorithm as above will lead to a strong overfit. In order to control for it the usual applied technique is the one of the following cost-complexity criterion

$$C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|$$

where:  $\alpha$  = complexity penalization coefficient  
 This is chosen through cross-validation.

$N_m$  = # observations in partition  $m$ .

$|T|$  = length of the tree (number of leaves).

$Q_m$  = goodness of fit for partition  $m$ .

Summing it up:

1. Use greedy algorithm to grow a large tree.
2. Apply cost complexity pruning to the large tree in order to obtain the best subtree for each  $\alpha$ .<sup>1</sup>
3. Use cross-validation to choose  $\alpha$ . i.e., divide the data into training and test data sets. Then, for each division:
  - (a) Repeat steps 1 and 2 on the training data.
  - (b) Evaluate the mean squared prediction error on the test data for each  $\alpha$ .
  - (c) Pick the  $\alpha$  which minimizes the average out-of-sample error on all test sets.
4. Return the subtree from step 2 that corresponds to the value of  $\alpha$  chosen in step 3.

<sup>1</sup>Only a finite number of  $\alpha$ 's needs to be considered since there are only finitely many subtrees. Moreover, one can show that the subtrees obtained in step 2 are nested.

Fabio Sigrist

Applied Multivariate Statistics

ie if you have different trees  
even if you have different trees  
the subtrees will share the same tree structure

## Random Trees:

This is essentially the same procedure of standard regression trees with the difference that  $N$  different trees are grown according to a later bootstrapping procedure. The prediction is then based on an ensemble - aggregation - of the individual trees:

1. For  $b = 1$  to  $B$ :

- (a) Draw a bootstrap sample  $Z^*$  of size  $N$  from the training data.
- (b) Grow a random-forest tree  $T_b$  to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size  $n_{min}$  is reached.
  - i. Select  $m$  variables at random from the  $p$  variables.
  - ii. Pick the best variable/split-point among the  $m$ .
  - iii. Split the node into two daughter nodes.
2. Output the ensemble of trees  $\{T_b\}_1^B$ .

To make a prediction at a new point  $x$ :

$$\text{Regression: } \hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x).$$

Classification: Let  $\hat{C}_b(x)$  be the class prediction of the  $b$ th random-forest tree. Then  $\hat{C}_{\text{rf}}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$ .

Important is that with this second technique we do not prune.

The idea is in fact the one of reducing the MSE by reducing the variance as previously mentioned. It is in fact possible to see that through the aggregation the variance of the predictor is:

$$\begin{aligned} \text{Var}\left(\frac{1}{B} \sum_{b=1}^B T_b(x)\right) &= \frac{1}{B^2} \left( \sum_{b=1}^B \text{Var}(T_b(x)) + \sum_{i=1}^B \sum_{j \neq i}^B \text{Cor}(T_i(x), T_j(x)) \right) \\ &= \frac{1}{B^2} \sum_{b=1}^B \sigma^2 + \sum_{i=1}^B \sum_{j \neq i}^B p \cdot \sigma^2 \xrightarrow{\text{corr.}} \\ &= \frac{1}{B^2} (B \sigma^2 + (B-1) B p \sigma^2) \\ &= \frac{1}{B} (\sigma^2 + B p \sigma^2 - p \sigma^2) \\ &= p \sigma^2 + \frac{\sigma^2(1-p)}{B} \end{aligned}$$

So that it is clear that the lower the correlation  $p$  among the trees the lower the variance. This is in fact the reason why for each tree you select randomly  $m$  features/variables.

Estimating generalized error for the Random Forest can be done by averaging over the out-of-bag error for each tree.

The idea is the one of splitting the dataset into training and testing sets. Then you would obtain the dataset for growing the tree based on the bootstrapped training dataset.

R. I. U. I. . . . . 1. I. U. I. L. the L. I. L. . . . . 1. I. L. at L. I. . . .

Then you would obtain the dataset for growing the tree based on the bootstrapped training dataset.

Based on that you would then obtain the **out-of-bag** error by looking at the share of testing set observations that were correctly classified.

An important field of research in the area of random forests is moreover the one of determining the relevance of a particular variable/feature for the classification case. Two major approaches were proposed:

① Mess up with a particular variable for which you want to get the importance.

↳ Understand how big the impact of messing up with this variable was.

That is more concretely, permute variable  $i$  across all of the trees. Then get the new out-of-bag errors after the permutation.

Look at the average out-of-bag error increase after the permutation.

The bigger the error increase the more important the variable.

To compare the effect among different variables you can then look at the average out-of-bag error rate corrected by the standard deviation.

② While growing the tree keep track of the decrease in impurity measure (be it RSS, Gini, MCE etc.) after each variable split.

↳ This gives an indication on the variable importance.

## Random and Fixed Effects Models

The idea of random effects models is the one of information pooling. The idea here is that when dealing with panel data, i.e. data collected in a cross-sectional way over a span of time.

Then one approach for dealing with the dataset is to estimate a fixed effect model of the form:

$$\vec{y}_t = \vec{\beta}_0 + \vec{\beta}^T \vec{X}_t + \vec{\epsilon}$$

where you could set:

$$\vec{\beta}_0 = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_0 \end{pmatrix}, \text{ i.e. a } \underline{\text{universal}} \text{ intercept for all the cross-sectional cases } i.$$

or in the case you might want to specify a different intercept for all of the states, this might even be necessary as if the underlying model requires one and you do not specify it you get correlated models and therefore you would get a biased model.

$$\vec{\beta}_0 = \begin{pmatrix} \beta_{0,1} \\ \vdots \\ \beta_{0,n} \end{pmatrix} \quad i = 1, \dots, n \text{ states.}$$

Important to see is that in such a model you do not put any constraints on the intercepts  $\beta_{0,i}$  and each one might differ sensibly from the other.

↳ Such an approach is therefore highly flexible but might suffer sensibly in the case of many different states  $i$  and just a few observations across time for them.

↳ In such a case the estimation of the fixed effects coefficients is highly variable and can be just poorly estimated.

A different approach in such a case is the one of using random effect to model the states coefficient. The model would then look as follows:

$$\vec{Y}_t = (\vec{\beta}_0 + \vec{v}) + \vec{\beta}^T \vec{X}_t + \vec{\epsilon}_t$$

where,

$\vec{v}$  and  $\vec{\epsilon}_t$  are independent

$$v_i \sim N(0, \sigma^2) \quad \text{if } v_i \in \vec{v}, \quad \epsilon_{it} \sim N(0, \sigma^2) \quad \text{if } \epsilon_{it} \in \vec{\epsilon}_t$$

Such model is a mixed effect model in that it has a fixed effect component  $\vec{\beta}^T$  and a random one:

$$\vec{\beta}_0 + \vec{v} = \begin{pmatrix} \beta_0 + v_1 \\ \vdots \\ \beta_0 + v_n \end{pmatrix}$$

The concept of the model is the one of information pooling:

- you estimate the  $\sigma^2, \beta_0, \vec{\beta}, \sigma^2$ .
- you obtain an estimate of  $v_i$  given the  $\sigma^2$ , for each state.
- The model is parsimonious and it uses the concept that you allow a specific state effect while keeping a pooling to the overall mean intercept as the mean and it is  $0$ .

Finally, notice that with such a model you account for the correlation of the observations across time

$$\text{Proof: } \text{Var}(y_t) = \text{Cov}(\beta_0 + v_i + \vec{\beta}^T \vec{X} + \epsilon_{it}, \beta_0 + v_i + \vec{\beta}^T \vec{X} + \epsilon_{it}) = \text{Cov}(v_i + \epsilon_{it}, v_i + \epsilon_{it}) \stackrel{\sigma^2 \text{ due to independence}}{=} \text{Var}(v_i) + \text{Var}(\epsilon_{it}) + \text{Cov}(v_i, \epsilon_{it})$$

$$\text{Cov}(y_t, y_{t'}) = \text{Var}(v_i) \quad \text{as correlation in time} \neq 0.$$

$$\text{Correlation in time} = \rho(y_t, y_{t'}) = \frac{\text{Var}(v_i)}{\text{Var}(v_i) + \text{Var}(\epsilon_{it})}$$

Notice moreover the interesting property that correlation is the same for units close and distant in time.

The model above can be further extended to include a random term in the independent variables for the different states  $i$ .

Then in its simplest form (i.e. in the case of a single independent variable coefficient), the model will look as follows:

$$y_{it} = (\beta_0 + v_{it}) + (\beta_1 + v_{i2}) X_{it} + \epsilon_{it}$$

with:

$$\vec{v}_i = \begin{pmatrix} v_{i1} \\ v_{i2} \end{pmatrix} \sim N(\vec{0}, \Sigma) \quad \Rightarrow \text{i.e. normal multivariate dist.}$$

$$\epsilon_{it} \sim N(0, \sigma^2)$$

so that now the parameters to be estimated within the model are:  $\beta_0, \beta_1, \Sigma, \sigma^2$

And the covariance structure would look as follows:

$$\begin{aligned} \text{Var}(y_{it}) &= \text{Cov}(\beta_0 + v_{it} + (\beta_1 + v_{i2}) X_{it} + \epsilon_{it}, \beta_0 + v_{it} + (\beta_1 + v_{i2}) X_{it} + \epsilon_{it}) \\ &= \text{Cov}(v_{it}, v_{it}) + 2\text{Cov}(v_{it}, \epsilon_{it}) + 2\text{Cov}(v_{it}, v_{i2}) X_{it} + 2\text{Cov}(v_{i2}, \epsilon_{it}) X_{it} + \text{Var}(\epsilon_{it}) \\ &\quad + \text{Cov}(v_{i2}, v_{i2}) X_{it}^2 \end{aligned}$$

$$\begin{aligned}
&= \text{Cov}(v_{it}, v_{it}) + 2\text{Cov}(v_{it}, \epsilon_{it}) + 2\text{Cov}(v_{it}, v_{it})X_{it} + 2\text{Cov}(v_{it}, \epsilon_{it})X_{it} + \text{Var}(\epsilon_{it}) \\
&\quad + \text{Cov}(v_{it}, v_{it}) X_{it}^2 \\
&= \text{Var}(v_{it}) + \text{Var}(v_{it}) X_{it}^2 + \text{Var}(\epsilon_{it}) + 2\text{Cov}(v_{it}, v_{it})X_{it} \\
\text{Cov}(y_{it}, y_{it'}) &= \text{Cov}(v_{it}, v_{it}) + \text{Cov}(v_{it}, v_{it'})X_{it} + \text{Cov}(v_{it}, v_{it'})X_{it'} + \\
&\quad \text{Cov}(v_{it}, v_{it'}) X_{it} \cdot X_{it'}
\end{aligned}$$

here the  
 correlation  
 time also depends on  
 time itself.

As a final note, notice that the above models can be estimated through the standard ways:

- ① Maximum likelihood
- ② Restricted Maximum Likelihood (the iterative approach).

- Maximum likelihood (ML):
  - Variance estimates are **biased**
  - + Tests between two models with differing fixed and random effects are possible
- Restricted maximum likelihood (REML):
  - + Variance estimates are **unbiased**
  - Can only test between two models that have the **same fixed effects**

Recommended  
(default in R)