

Latent Variable Methods Course

Learning from data

Instructor: Kevin Dunn
kevin.dunn@connectmv.com
<http://connectmv.com>

© Kevin Dunn, ConnectMV, Inc. 2011

Revision: 269:35e2 compiled on 15-12-2011

Copyright, sharing, and attribution notice

This work is licensed under the Creative Commons Attribution-ShareAlike 3.0 Unported License. To view a copy of this license, please visit

<http://creativecommons.org/licenses/by-sa/3.0/>



This license allows you:

- ▶ **to share** - to copy, distribute and transmit the work
- ▶ **to adapt** - but you must distribute the new result under the same or similar license to this one
- ▶ **commercialize** - you are allowed to create commercial applications based on this work
- ▶ **attribution** - you must attribute the work as follows:
 - ▶ "Portions of this work are the copyright of ConnectMV", *or*
 - ▶ "This work is the copyright of ConnectMV"

We appreciate:

- ▶ if you let us know about **any errors** in the slides
- ▶ **any suggestions to improve the notes**
- ▶ telling us if you use the slides, especially commercially, so we can inform you of major updates
- ▶ emailing us to ask about different licensing terms

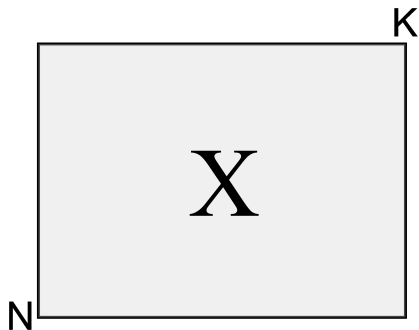
All of the above can be done by writing us at

courses@connectmv.com

If reporting errors/updates, please quote the current revision number: 269:35e2

Data sources

- ▶ PCA considers a single data table (matrix)
- ▶ We will call it \mathbf{X}



- ▶ N observations
- ▶ K variables
- ▶ What goes in the columns of \mathbf{X} ?
- ▶ What goes in the rows?

Visualization

How would you visualize such a data table?

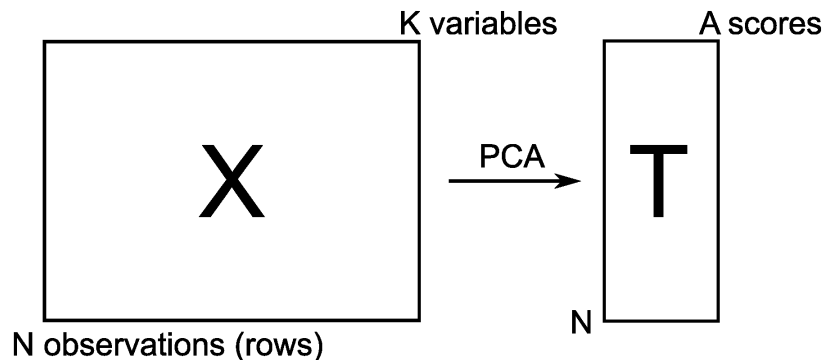
For example: assume $N = 300$ and $K = 50$

- ▶ One column at a time (time-series, histograms, boxplot)
- ▶ One row at a time (e.g. spectral data)
- ▶ Scatterplot matrix, requires $K(K - 1)/2$ pairs of scatterplots

What is PCA (Principal Components Analysis)?

Mathematical objective

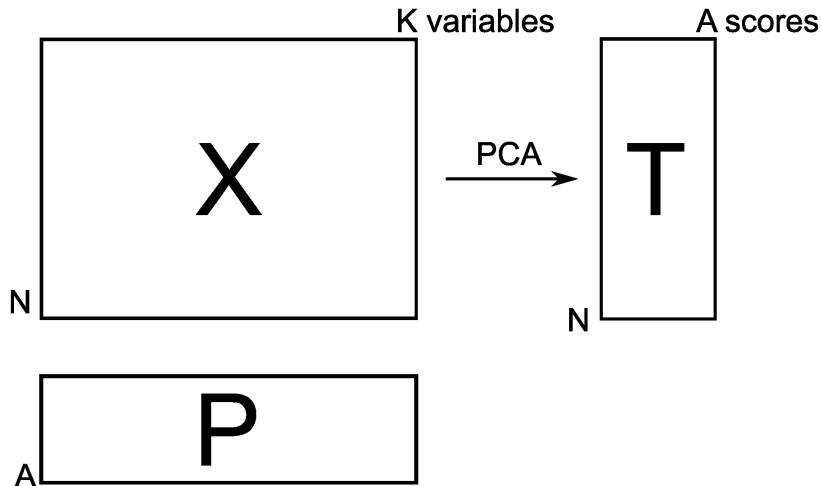
PCA: find me the best summary of my data, \mathbf{X} , with the fewest number of summary variables, called scores, \mathbf{T} .



Objectives for this class

PCA model will calculate from \mathbf{X} :

- ▶ scores: \mathbf{T}
- ▶ loadings: \mathbf{P}



Objectives for this class

- ▶ Intuitive meaning of the scores, \mathbf{T} and loadings, \mathbf{P} and errors in a PCA model
- ▶ The interpretation of each of these
- ▶ How to start investigating a new data table

Time to break out the math

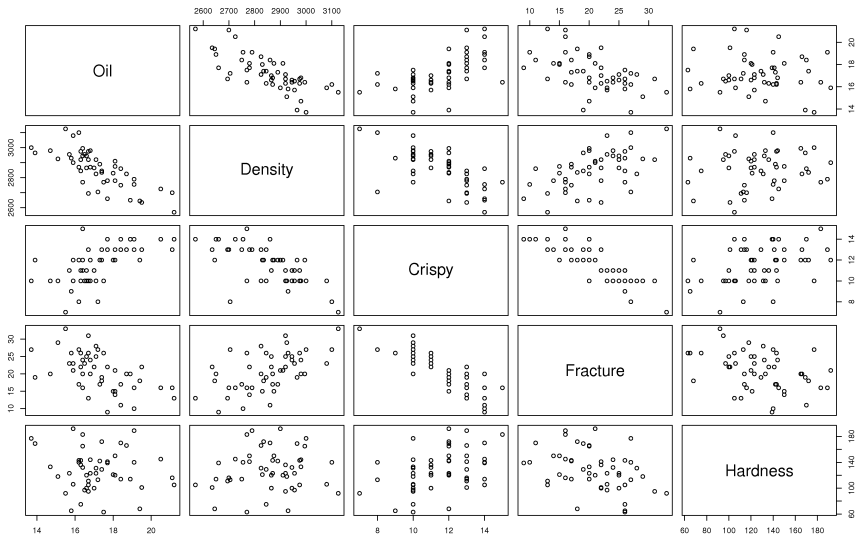
- ▶ Notation for scores: $t_{1,1}, t_{2,1}, \dots, t_{n,1}, \dots, t_{N,1}$
- ▶ Notation for loadings: $p_{1,1}, p_{2,1}, \dots, p_{k,1}, \dots, p_{K,1}$
- ▶ Length of a vector: $\|a\|$

Let's get started

1. Data preprocessing
2. Geometric interpretation (hand waving explanation)
3. Analytical geometry (to understand the hand waving)
4. Algebraic approach (to formalize the notation)
5. Look at applying this all in software

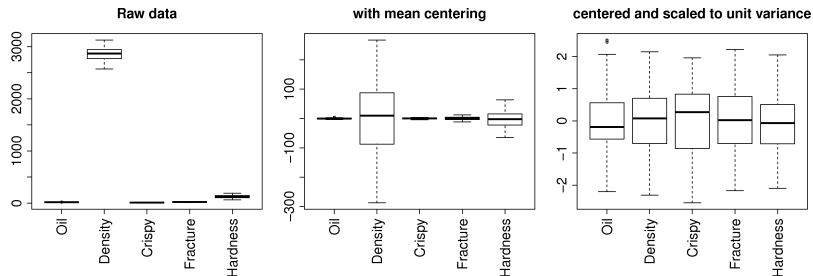
Preprocessing by example

Raw data:



Preprocessing by example

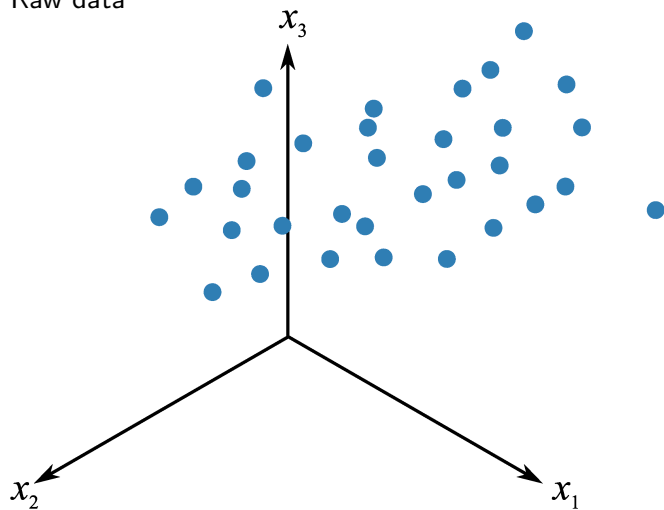
Center and scale the raw data



- ▶ Centering: $\mathbf{x}_{k,\text{center}} = \mathbf{x}_{k,\text{raw}} - \text{mean}(\mathbf{x}_{k,\text{raw}})$
- ▶ Scaling: $\mathbf{x}_k = \frac{\mathbf{x}_{k,\text{center}}}{\text{standard deviation}(\mathbf{x}_{k,\text{center}})}$
- ▶ Does *not change* relationships between variables

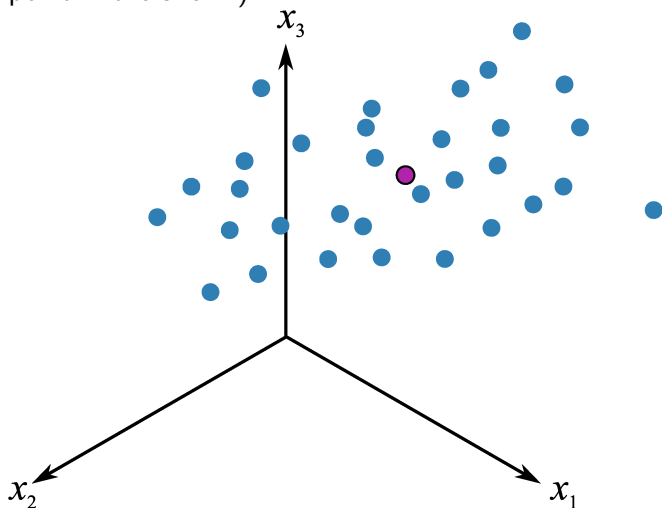
Geometric explanation of preprocessing

Raw data



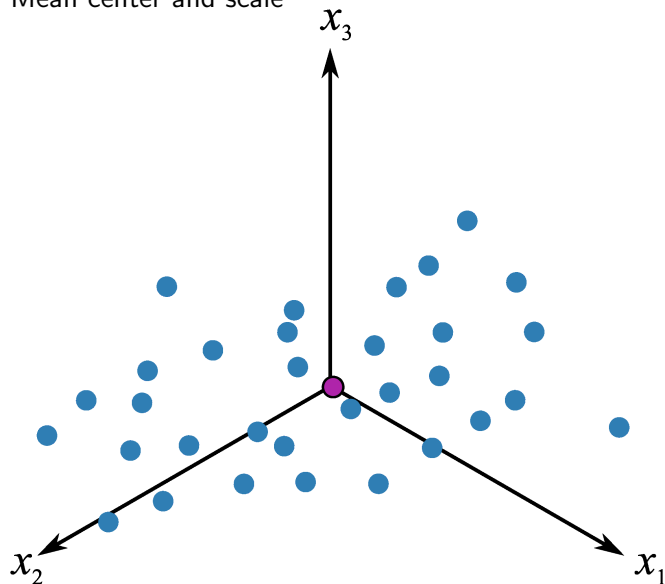
Geometric explanation of preprocessing

Calculate the mean of each variable (creates a “new” reference point in the swarm)

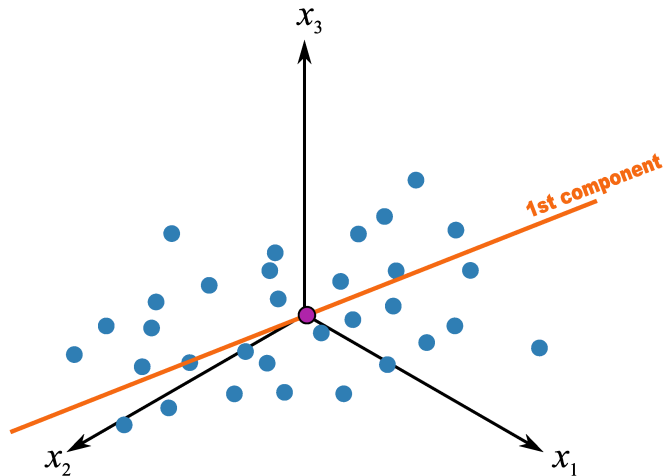


Geometric explanation of preprocessing

Mean center and scale

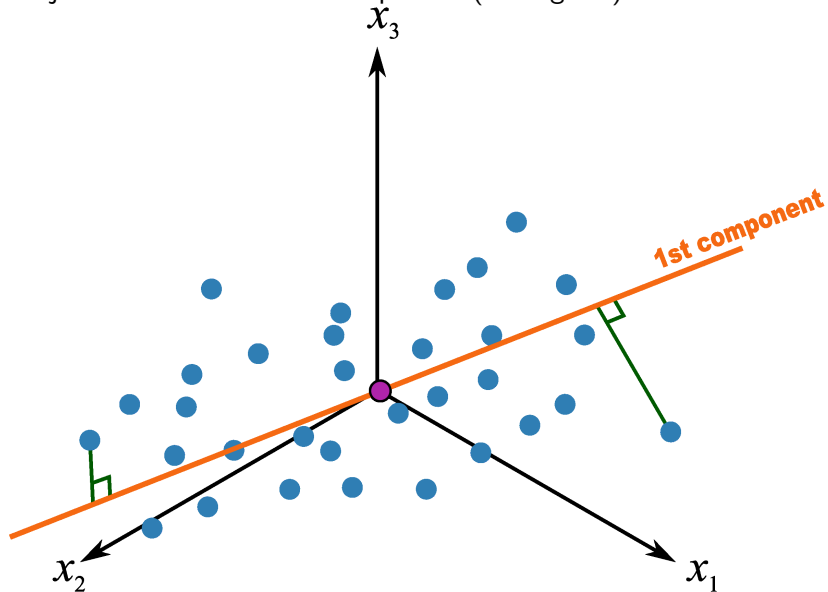


Geometric explanation of PCA



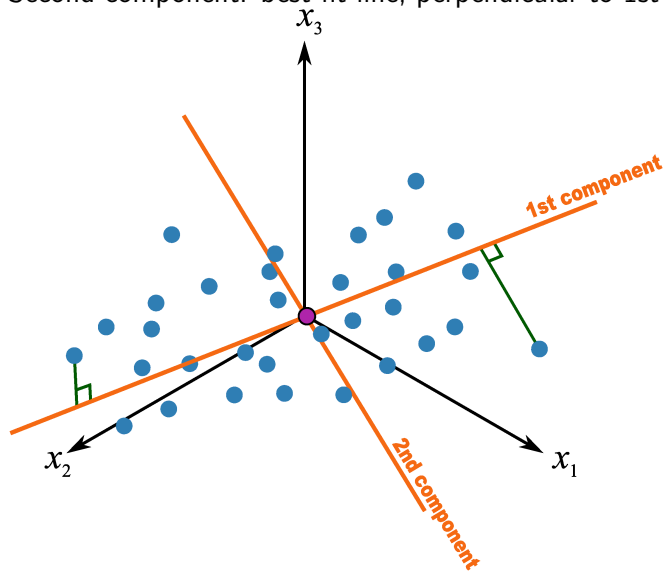
Geometric explanation of PCA

Project observations onto component (90 degrees)



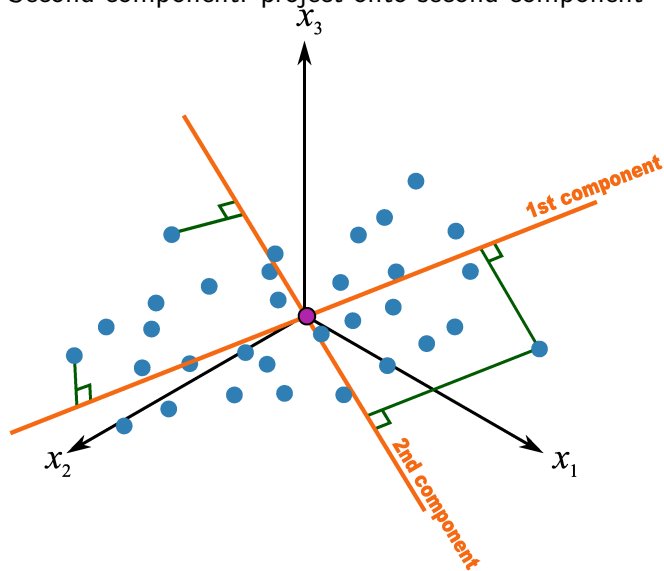
Geometric explanation of PCA

Second component: best-fit line; perpendicular to 1st component



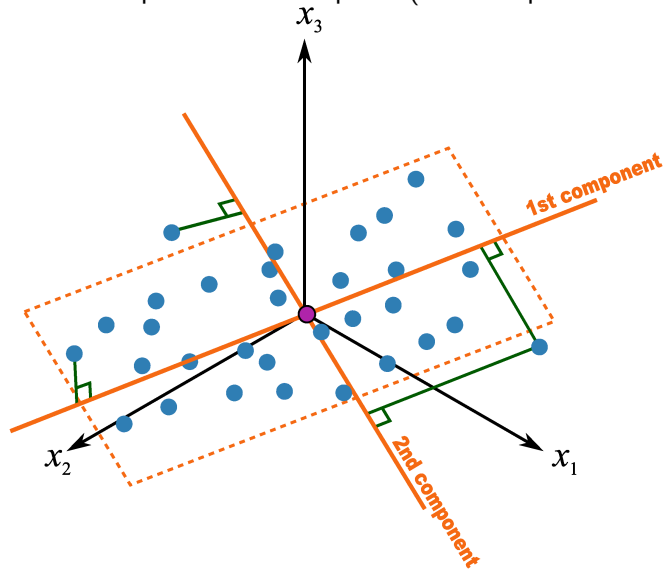
Geometric explanation of PCA

Second component: project onto second component



Geometric explanation of PCA

The 2 components form a plane (2-D subspace inside a 3D space)



Geometric explanation of PCA

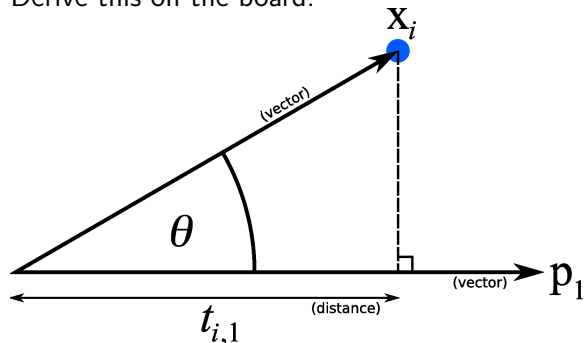
What have we done here?

Broken **X** down into 2 parts:

- ▶ projected points "*on the plane*"
- ▶ residual distance "*off the plane*"

Mathematical derivation for PCA

Derive this on the board:



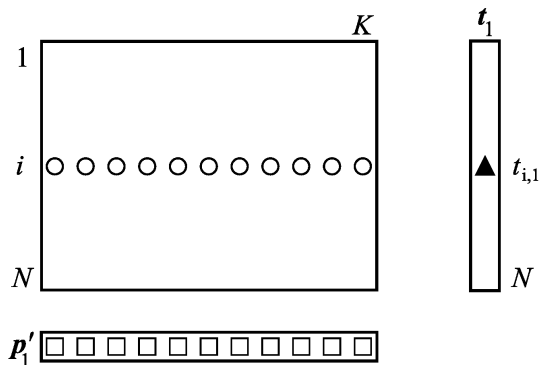
$$\cos \theta = \frac{\text{adjacent length}}{\text{hypotenuse}} = \frac{t_{i,1}}{\|\mathbf{x}_i\|}$$

$$\text{and also } \cos \theta = \frac{\mathbf{x}_i^T \mathbf{p}_1}{\|\mathbf{x}_i\| \|\mathbf{p}_1\|}$$

$$\frac{t_{i,1}}{\|\mathbf{x}_i\|} = \frac{\mathbf{x}_i^T \mathbf{p}_1}{\|\mathbf{x}_i\| \|\mathbf{p}_1\|}$$

$$\begin{aligned} t_{i,1} &= \mathbf{x}_i^T \mathbf{p}_1 \\ (1 \times 1) &= (1 \times K)(K \times 1) \end{aligned}$$

Mathematical derivation for PCA



$$\begin{aligned} t_{i,1} &= \mathbf{x}_i^T \mathbf{p}_1 \\ &= x_{i,1}p_{1,1} + x_{i,2}p_{2,1} + \dots + x_{i,k}p_{k,1} + \dots + x_{i,K}p_{K,1} \end{aligned}$$

- ▶ K individual terms add up (i.e. linear combination) to give t_1
- ▶ Stack $t_{i,1}$ values from N rows: $\mathbf{T} = \mathbf{XP}$

Interpreting $t_i = \mathbf{x}_i^T \mathbf{p}_1$

Given that:

- ▶ values in \mathbf{x}_i^T are centred and scaled, and
- ▶ entries in \mathbf{p} are between -1 and $+1$

using

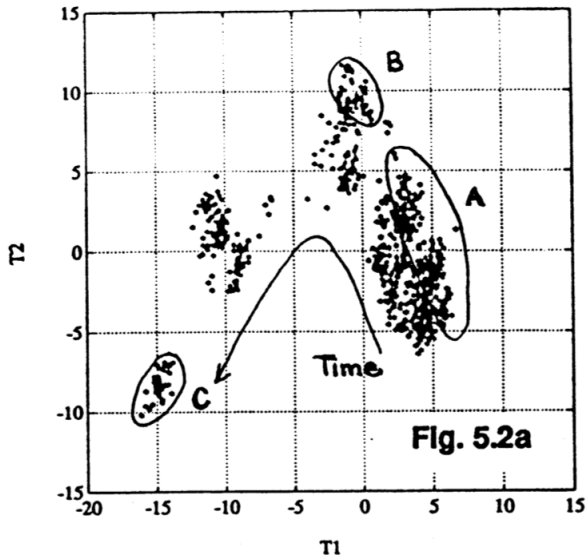
$$= x_{i,1}p_{1,1} + x_{i,2}p_{2,1} + \dots + x_{i,k}p_{k,1} + \dots + x_{i,K}p_{K,1}$$

how would you

- ▶ get a large positive value of $t_{i,1}$?
 - ▶ get a large negative value of $t_{i,1}$?
 - ▶ get a value of $t_{i,1} \approx 0$?
-
- ▶ What can you say about observation (row) 13 and 22 if $t_{13,1} \approx t_{22,1}$?

Score plots: interpretation

Clustering



Score plots: interpretation

Also look for:

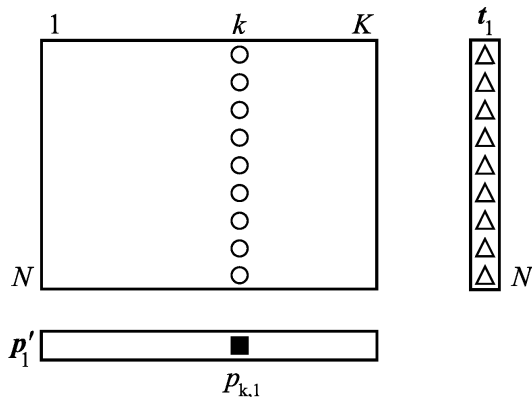
- ▶ outliers
- ▶ patterns in the sequence order (if time-based row order)
- ▶ colour-code score plots by another variable (good/bad)

We'll see more of these tips as we work with the software.

Key point: anything you would normally have done to visualize a column can be done with a score.

Mathematical derivation for PCA

We'll look at this in next class:



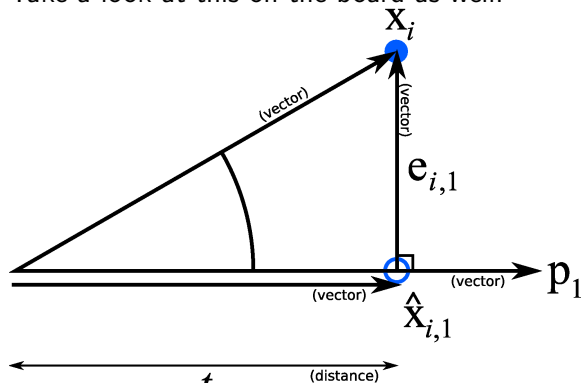
For now though:

- ▶ Columns that are related have similar loadings
- ▶ “Direction vectors” = “Loadings”
- ▶ Link between the real-world (K) and latent-variable world (A)

$$\begin{aligned}\mathbf{T} &= \mathbf{X}\mathbf{P} \\ (N \times A) &= (N \times K)(K \times A)\end{aligned}$$

Predicted values for each observation

Take a look at this on the board as well:



$$\begin{aligned}\hat{\mathbf{x}}_{i,1}^T &= t_{i,1} \mathbf{p}_1^T \\ (1 \times K) &= (1 \times 1)(1 \times K) = \text{best prediction of } \mathbf{x}_{i,1} \text{ with 1 component}\end{aligned}$$

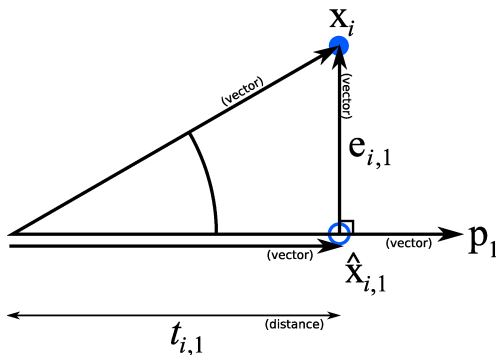
In this case: “best” = “smallest error”

The residuals

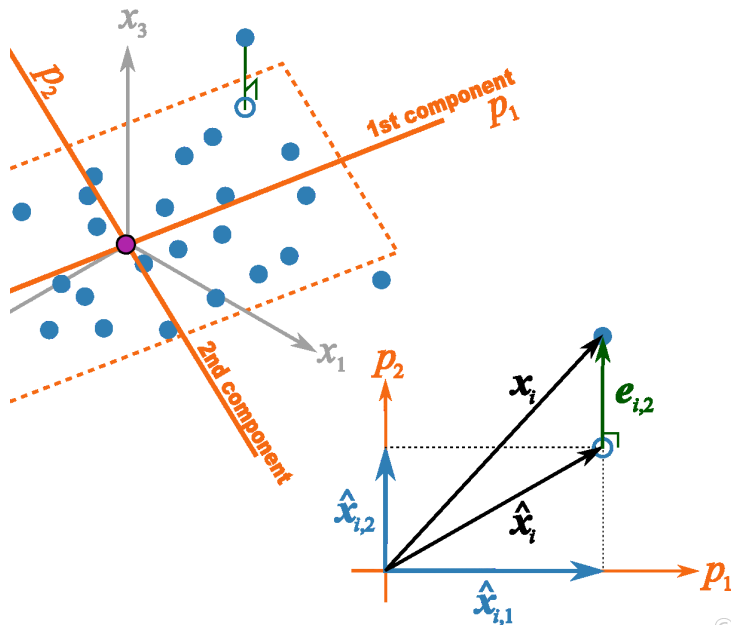
Residuals for row i after extracting one component = $\mathbf{e}_{i,1}$

$$\mathbf{e}_{i,1}^T = \mathbf{x}_i^T - \hat{\mathbf{x}}_{i,1}^T \quad (\text{each is a } 1 \times K \text{ vector})$$

Another way of stating this: $\mathbf{x}_i^T = \hat{\mathbf{x}}_{i,1}^T + \mathbf{e}_{i,1}^T$
 $\mathbf{x}_i^T = t_{i,1} \mathbf{p}_1^T + \mathbf{e}_{i,1}^T$

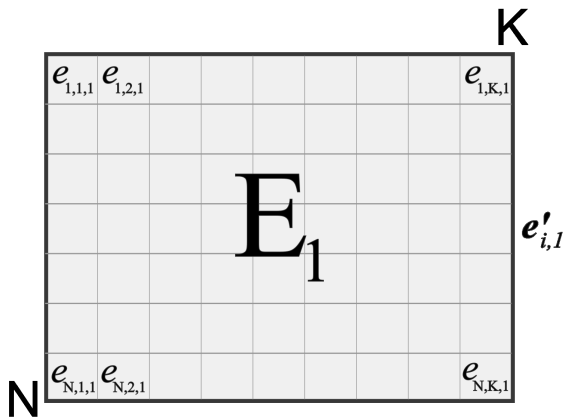


Predictions, residuals, vectors: explained



The residuals

Assemble the residuals for every row in a matrix, \mathbf{E}_1



The residuals

The next few slides discuss the residuals

- ▶ important part of fitting a model
- ▶ ideally, contains no information (just noise)

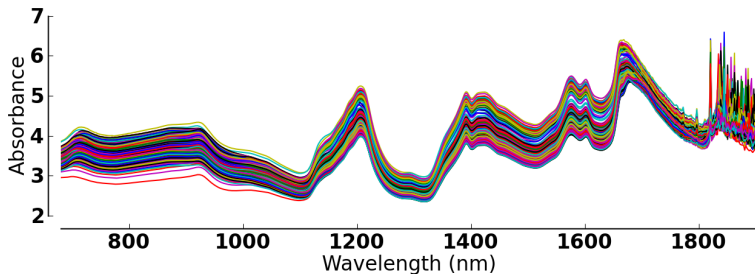
We will consider

- ▶ whole matrix residuals
- ▶ column residuals (per variable)
- ▶ row residuals (per observation)

Main way of quantifying residuals:

- ▶ calculate their sum of squares (ssq)
- ▶ in this case the ssq = variance
- ▶ and $R^2 = \frac{\text{variance explained by model}}{\text{initial variance}}$

Residuals: spectral example



- ▶ Data on course website
- ▶ Try it yourself: <http://datasets.connectmv.com/info/tablet-spectra>
- ▶ $N = 460$
- ▶ $K = 650$

Whole matrix residuals

- ▶ $\mathbf{X} = \mathbf{TP}' + \mathbf{E} = \hat{\mathbf{X}} + \mathbf{E}$

- ▶ Quantify how well the model (\mathbf{TP}') fits the data

- ▶ $R_a^{2(\text{overall})} = 1 - \frac{\text{Var}(\mathbf{X} - \hat{\mathbf{X}}_a)}{\text{Var}(\mathbf{X})} = 1 - \frac{\text{Var}(\mathbf{E}_a)}{\text{Var}(\mathbf{X})}$

- ▶ $R_{a=0}^2 = 0.0$ (no components, means no variance explained)

- ▶ R^2 increases with every component added

- ▶ $R_{a=1}^{2(\text{overall})} < R_{a=2}^{2(\text{overall})} < \dots < R_{a=A}^{2(\text{overall})} = 1.0$

Matrix residuals: spectral example

$$\begin{matrix} & K \\ \begin{matrix} N \\ \boxed{X} \end{matrix} & = & \begin{matrix} K \\ \boxed{\hat{X}} \end{matrix} & + & \begin{matrix} K \\ \boxed{E} \end{matrix} \end{matrix}$$

Variance = 100%	a=1: $R^2 = 73.7\%$ [73.7%]	a=1: $1 - R^2 = 26.3\%$
	a=2: $R^2 = 18.5\%$ [92.2%]	a=2: $1 - R^2 = 7.8\%$
	a=3: $R^2 = 1.99\%$ [94.2%]	a=3: $1 - R^2 = 5.8\%$

- ▶ $R_{a=1}^2 = 73.7\%$
- ▶ $R_{a=2}^2 = 92.2\%$ (an additional 18.5%)
- ▶ $R_{a=3}^2 = 94.2\%$ (an additional 2.00%)

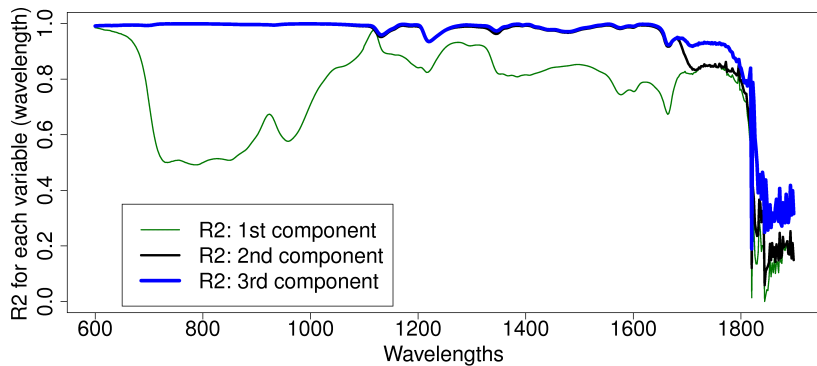
Column residuals

- ▶ R^2 can be calculated for each column

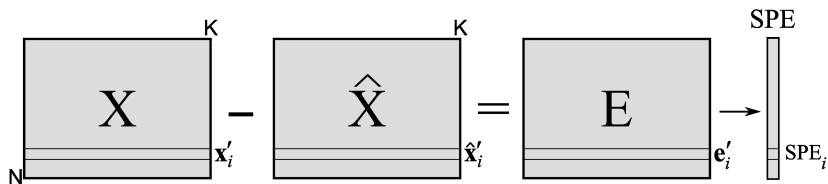
$$\begin{matrix} & & K \\ \begin{matrix} N \\ \text{---} \end{matrix} & \boxed{X} & - & \boxed{\hat{X}} & = & \boxed{E} \\ & \text{---} & & \text{---} & & \text{---} \\ & \mathbf{x}_k & & \hat{\mathbf{x}}_k & & \mathbf{e}_k \end{matrix} \longrightarrow R_k^2$$

- ▶ $R_k^2 = 1 - \frac{\text{Var}(\mathbf{x}_k - \hat{\mathbf{x}}_k)}{\text{Var}(\mathbf{x}_k)}$
- ▶ indicates how well each column is explained by the model
- ▶ is 0.0 when there are no components
- ▶ increases for every every component added

Column residuals: spectral example



Row residuals



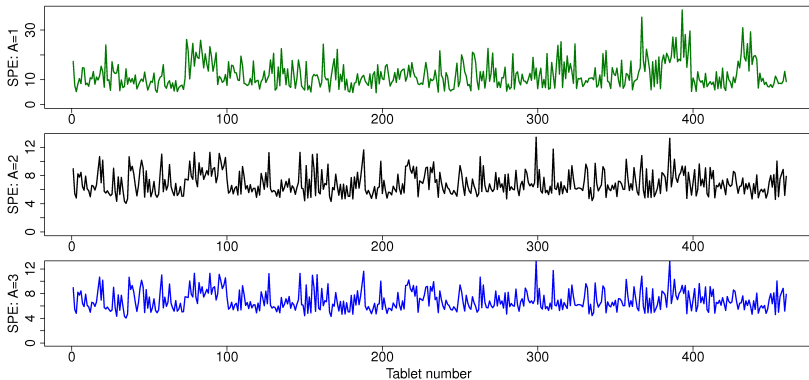
- ▶ $\mathbf{e}'_i = \mathbf{x}'_i - \hat{\mathbf{x}}'_i$
 $\mathbf{e}'_i = [(x_{i,1} - \hat{x}_{i,1}) \quad (x_{i,2} - \hat{x}_{i,2}) \quad \dots \quad (x_{i,k} - \hat{x}_{i,k}) \quad \dots \quad (x_{i,K} - \hat{x}_{i,K})]$
- ▶ Variance of residuals in a row = $e_{i,1}^2 + e_{i,2}^2 + \dots + e_{i,K}^2$
- ▶ Call this SPE = squared prediction error = $\mathbf{e}_i^T \mathbf{e}_i$
- ▶ Square root of SPE = “distance to model’s X-space”
- ▶ “DModX” (used in some software) is related to $\sqrt{\text{SPE}_i}$

Square prediction error

Distance from each observation to the model's plane:

- ▶ Smallest SPE ?
- ▶ Distribution of SPE values
- ▶ If $\text{SPE} > 95\%$ limit:
 - ▶ poorly explained by the model
 - ▶ something new in this observation
 - ▶ new phenomenon?

Row residuals: spectral example



Data sets to look at in class

- ▶ **Website:** <http://datasets.connectmv.com>
- ▶ Click on the “Peas” link and download CSV file
- ▶ Click on the “Food consumption” link and download CSV file
- ▶ Click on the “Food texture” link and download CSV file

For next class

1. Assignment: instructions will be posted on course website
2. Read the paper by Wold (item 12)
 - ▶ <http://literature.connectmv.com/item/12>
 - ▶ This will help you understand the material in next class
 - ▶ No Q&A: but strongly recommend you are familiar with the concepts
3. Next class will cover
 - ▶ *how* we calculate the components
 - ▶ *how many* components should be calculated
 - ▶ *using* the PCA model on new data