# Dimension Reduction

# Principal Components Analysis

# Mechanics

**Input**: $X_1$, $X_2$, ... , $X_p$
**Output**: $PC_1$, $PC_2$, ... , $PC_p$ (Ultimately we'll use a subset)

The i$^{th}$ principal component is a weighted average:

$$PC_i = a_{i1} X_1 + a_{i2} X_2 + ... + a_{ip} X_p$$

Weights chosen such that:
1. PCs are ordered by their variance ($PC_1$ has largest variance)
2. Pairs of PCs have correlation = 0
3. For each PC, sum of squared weights =1

# Example: Business School Programs

| Univ | SAT | Top10 | Accept | SFRatio | Expenses | GradRate |
|---|---|---|---|---|---|---|
| Brown | 1310 | 89 | 22 | 13 | 22,704 | 94 |
| CalTech | 1415 | 100 | 25 | 6 | 63,575 | 81 |
| CMU | 1260 | 62 | 59 | 9 | 25,026 | 72 |
| Columbia | 1310 | 76 | 24 | 12 | 31,510 | 88 |
| Cornell | 1280 | 83 | 33 | 13 | 21,864 | 90 |
| Dartmouth | 1340 | 89 | 23 | 10 | 32,162 | 95 |
| Duke | 1315 | 90 | 30 | 12 | 31,585 | 95 |
| Georgetown | 1255 | 74 | 24 | 12 | 20,126 | 92 |
| Harvard | 1400 | 91 | 14 | 11 | 39,525 | 97 |
| JohnsHopkins | 1305 | 75 | 44 | 7 | 58,691 | 87 |
| MIT | 1380 | 94 | 30 | 10 | 34,870 | 91 |
| Northwestern | 1260 | 85 | 39 | 11 | 28,052 | 89 |
| NotreDame | 1255 | 81 | 42 | 13 | 15,122 | 94 |
| PennState | 1081 | 38 | 54 | 18 | 10,185 | 80 |
| Princeton | 1375 | 91 | 14 | 8 | 30,220 | 95 |
| Purdue | 1005 | 28 | 90 | 19 | 9,066 | 69 |
| Stanford | 1360 | 90 | 20 | 12 | 36,450 | 93 |
| TexasA&M | 1075 | 49 | 67 | 25 | 8,704 | 67 |
| UCBerkeley | 1240 | 95 | 40 | 17 | 15,140 | 78 |
| UChicago | 1290 | 75 | 50 | 13 | 38,380 | 87 |
| UMichigan | 1180 | 65 | 68 | 16 | 15,470 | 85 |
| UPenn | 1285 | 80 | 36 | 11 | 27,553 | 90 |
| UVA | 1225 | 77 | 44 | 14 | 13,349 | 92 |
| UWisconsin | 1085 | 40 | 69 | 15 | 11,857 | 71 |
| Yale | 1375 | 95 | 19 | 11 | 43,514 | 96 |

Use PCA to:

1) Reduce # columns

2) Identify relations between columns

3) Visualize universities in 2D

Source: US News & World Report, Sept 18 1995

# PCA in XLMiner

*Data Reduction & Exploration*

Output specifies whether covariance or correlation matrix used (here – correlation matrix).

**Principal Components**

| Variable | Components 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| SAT | 0.45774868 | 0.03968045 | 0.18703876 | 0.13124055 | 0.02064597 | -0.8580547 |
| Top10 | 0.42714444 | -0.19993152 | 0.49780852 | 0.37489522 | 0.48201644 | 0.39607504 |
| Accept | -0.42430812 | 0.32089293 | -0.15627895 | 0.06128667 | 0.80109364 | -0.21693356 |
| SFRatio | -0.39064837 | -0.43256435 | 0.60608089 | -0.50739086 | 0.07682328 | -0.17204805 |
| Expenses | 0.3625232 | 0.63448638 | 0.20474122 | -0.62340063 | 0.07254726 | 0.17376293 |
| GradRate | 0.37940401 | -0.51555371 | -0.53247261 | -0.43863374 | 0.33810937 | 0.00353743 |

| | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Variance | 4.61208487 | 0.78681612 | 0.28656188 | 0.16378011 | 0.12430621 | 0.02645062 |
| Variance% | 76.86808014 | 13.11360168 | 4.77603149 | 2.72966838 | 2.07177019 | 0.44084364 |
| Cum% | 76.86808014 | 89.98168182 | 94.75771332 | 97.48738098 | 99.5591507 | 100 |
| P-value | 0 | 0.00000004 | 0.00073126 | 0.00263538 | 0.00140999 | 1 |

# Reducing data dimension

**Principal Components**

| Variable | Components | | | | | |
| --- | ---: | ---: | ---: | ---: | ---: | ---: |
| | 1 | 2 | 3 | 4 | 5 | 6 |
| SAT | 0.45774868 | 0.03968045 | 0.18703876 | 0.13124055 | 0.02064597 | -0.8580547 |
| Top10 | 0.42714444 | -0.19993152 | 0.49780852 | 0.37489522 | 0.48201644 | 0.39607504 |
| Accept | -0.42430812 | 0.32089293 | -0.15627895 | 0.06128667 | 0.80109364 | -0.21693356 |
| SFRatio | -0.39064837 | -0.43256435 | 0.60608089 | -0.50739086 | 0.07682328 | -0.17204805 |
| Expenses | 0.3625232 | 0.63448638 | 0.20474122 | -0.62340063 | 0.07254726 | 0.17376293 |
| GradRate | 0.37940401 | -0.51555371 | -0.53247261 | -0.43863374 | 0.33810937 | 0.00353743 |

| | | | | | | |
| --- | ---: | ---: | ---: | ---: | ---: | ---: |
| Variance | 4.61208487 | 0.78681612 | 0.28656188 | 0.16378011 | 0.12430621 | 0.02645062 |
| Variance% | 76.86808014 | 13.11360168 | 4.77603149 | 2.72966838 | 2.07177019 | 0.44084364 |
| Cum% | 76.86808014 | 89.98168182 | 94.75771332 | 97.48738098 | 99.5591507 | 100 |
| P-value | 0 | 0.00000004 | 0.00073126 | 0.00263538 | 0.00140999 | 1 |

$PC_1$ captures _____ % of the information

The first two PCs capture _____%

# XLMiner: Computing scores

Scores given for each PC

Recall: PC1 & PC2 are uncorrelated ($r = 0$)

| Row Id. | 1 | 2 |
|---|---|---|
| 1 | 0.98947096 | -1.04280615 |
| 2 | 2.76521754 | 2.21340251 |
| 3 | -1.08998942 | 1.5982517 |
| 4 | 0.72675508 | -0.04133511 |
| 5 | 0.30561018 | -0.62240905 |
| 6 | 1.66241097 | -0.33740574 |
| 7 | 1.2216301 | -0.48106378 |
| 8 | 0.33190566 | -0.76930493 |
| 9 | 2.32618284 | -0.37872922 |
| 10 | 1.37492549 | 2.07669187 |
| 11 | 1.69122922 | 0.08645435 |
| 12 | 0.44174835 | -0.01090807 |
| 13 | -0.03942522 | -0.98881435 |
| 14 | -3.168396 | -0.36701241 |
| 15 | 2.19108367 | -0.36428159 |
| 16 | -5.06847715 | 0.76415795 |
| 17 | 1.66530418 | -0.29942313 |
| 18 | -4.48564911 | -0.3405683 |
| 19 | -0.80598319 | -0.68478525 |
| 20 | 0.09578693 | 0.63730472 |
| 21 | -1.92351854 | -0.22022633 |
| 22 | 0.53133249 | -0.07798085 |
| 23 | -0.52146798 | -0.99661624 |
| 24 | -3.47699881 | 0.76273364 |
| 25 | 2.25931191 | -0.11532623 |

# Score plot (score2 vs score1) using Spotfire

# SVD (similar to PCA)

Data Matrix:  **M = U Σ V'**



Using an example from the Wikipedia page:

$$
\begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ \sqrt{0.2} & 0 & 0 & 0 & \sqrt{0.8} \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}
$$

### Books rated by users

|        | User 1 | User 2 | User 3 | User 4 |
|--------|--------|--------|--------|--------|
| Book 1 | 1      | 4      |        | 1      |
| Book 2 | 2      |        |        | 5      |
| Book 3 |        | 4      | 1      | 4      |
| Book 4 | 1      |        |        |        |
| Book 5 | 3      | 5      | 5      | 1      |
| Book 6 |        | 3      | 2      |        |

Original matrix:
$$
\begin{pmatrix} 1 & 4 & 0 & 1 \\ 2 & 0 & 0 & 5 \\ 0 & 4 & 1 & 4 \\ 1 & 0 & 0 & 0 \\ 3 & 5 & 5 & 1 \\ 0 & 3 & 2 & 0 \end{pmatrix}
$$

Rebuilt with $k$=4, the ma[...]

Approximation in rounde[...]
$$
\begin{pmatrix} 0 & 4 & 1 & 1 \\ 2 & 0 & 0 & 5 \\ 0 & 4 & 1 & 4 \\ 0 & 0 & 0 & 0 \\ 3 & 5 & 5 & 1 \\ 0 & 3 & 2 & 0 \end{pmatrix}
$$

Need only
**k columns of U
k values of Σ
k rows of V'**

http://journal.batard.info/post/2009/04/08/svd-fun-profit

# This week's online discussion

**Discussion: Reducing Dimension of Mobile Survey Data**

Prof. Galit Shmueli

Mar 21 at 3:29pm

16 | 20

Consider a sample from a dataset on mobile us... mobile users in India. Take a look at the data sa...

**Choose one of the points below and post a ... thread with the adequate point.**

1. What approaches would you take to reduce ... which method you would apply to which col...

2. Suppose the goal is *describing* the relation... potential factors (customer demographics, t... apply PCA? Which columns would you app... describe the relationship?

3. For *predicting* service switching (churn), if w... space, what information would we need to p... PCA reduce the number of questions that w...

| Column Name | Description |
| --- | --- |
| serialnum | ID of respondent |
| StartDate | Survey start date/time |
| EndDate | Survey end date/time |
| SurveyDuration (Hrs.Min) | Survey Duration |
| Completed | Whether the survey was completed (1=yes) |
| Num Mobiles | Q1 Do you currently own one or more Mobiles?<br>☐ No. I use landline only. (1)<br>☐ One - with single SIM (2)<br>☐ One handset with two SIMs (3)<br>☐ Two handsets (4)<br>☐ More than two handsets or more than two SIMs (5) |
| Mobile Type (Primary) | What type of mobile phone handset do you own? If you own more than one, ple...<br>☐ Basic phone without internet capability (1)<br>☐ Smartphone (non-touchscreen) (2)<br>☐ Smartphone with touch-screen (3)<br>☐ Tablet with phone features (4) |
| Service Provider (Primary) | Who is your current service provider for your primary mobile phone?<br>☐ Airtel (1)<br>☐ Reliance (2)<br>☐ Idea (3)<br>☐ Vodafone (4)<br>☐ Tata DOCOMO (5)<br>☐ Tata Indicom (6)<br>☐ Aircel (7)<br>☐ BSNL / MTNL (8)<br>☐ Uninor (9)<br>☐ Virgin Mobile (10)<br>☐ Other (11) |
| Network duration | Q41 How long have you been on this network?<br>☐ Less than 6 months (1)<br>☐ 6 months to 1 year (2)<br>☐ 1 to 2 years (3)<br>☐ More than 2 years (4) |
| Mobile Service | Q42 Which mobile service type do you use for your primary mobile phone?<br>☐ GSM (1)<br>☐ CDMA (2)<br>☐ Not sure/Don't know. (3) |
| Provider-Network Coverage | Rate your service provider on this area (1=poor, 2=below average, 3=average, ... |
| Provider-Call quality | Rate your service provider on this area (1=poor, 2=below average, 3=average, ... |
| Provider-Call charges | Rate your service provider on this area (1=poor, 2=below average, 3=average, ... |
| Provider-Roaming charges | Rate your service provider on this area (1=poor, 2=below average, 3=average, ... |
| Provider-Customer support | Rate your service provider on this area (1=poor, 2=below average, 3=average, ... |
| Provider-Offers and promotions | Rate your service provider on this area (1=poor, 2=below average, 3=average, ... |
| Provider-Easy bill payment, varied recharge options, etc. | Rate your service provider on this area (1=poor, 2=below average, 3=average, ... |

# Use compressed data in modeling

For predicting?

For explaining?

Data Mining Contest (Crowdanalytics)

Each restaurant, each year

Lots of variables!

MODELING - Olive Garden Restaurant Comparison Analysis

OPEN

📅 15 days left
👥 116 Solvers
🏆 USD $5000
🔓 Public

Perform store comparison analysis to uncover drivers that explain why some Olive Garden ...

Visualization | Statistical Modeling | Consumer Insight
Store-Level | Restaurant | Food and Beverages
Structured | Moderate

Go to Competition

Uncover drivers that impacted performance of Olive Garden restaurants. Drivers identified should explain why some restaurants in the chain perform worse when compared to other restaurants in the chain.

https://crowdanalytix.com/contests/modeling---olive-garden-restaurant-comparison-analysis