

MLA

8/05/15. PRINCIPLE COMPONENT ANALYSIS

- For data with many variables/dimensions it is often difficult to comprehend or visualise inherent associations
- PCA thought of as a method for re-expressing the data so as to reveal its internal structure and explain its variation through the use of a few linear combinations of the original variables
- Used as either a dimension reduction technique or as a method for identifying associations among variables
- Aim of PCA is to describe the variation in a set of correlated variables X_1, \dots, X_m in terms of a new set of uncorrelated variables Y_1, \dots, Y_p hopefully $p \ll m$ and where each Y_i is a linear combination of the X_1, \dots, X_m .
- New variables 'Principal components', derived in decreasing order of importance so that first PC Y_1 accounts for more variation in the original data than any other possible linear combination of X_1, \dots, X_m .
- Second PC is chosen to account for as much of the remaining variation as possible subject to constraint that it is uncorrelated with Y_1 and so on...
- Hope is that the first few PCs will account for a substantial amount of variation in the original data, and as such can be used as a convenient lower dimensional summary of it.

Constrained Optimisation:

Find a to maximise $a^T Q a$ subject to $a^T a = 1$

Lagrange multipliers

- Given a function $f(x)$, gradient of f , ∇f (the collection of partial derivatives) indicates the direction of the steepest slope. Level curves (lines where $f(x)$ is constant) run perpendicular to the gradient.

- Here we want location where $\nabla f = \lambda \nabla g$.
- In other words, find x so that $\nabla f - \lambda \nabla g = 0$.
- Gives us $m+1$ eqns down with $m+1$ unknowns.

- Let $p = a^T Q a - \lambda (a^T a - 1)$
- Solution found by computing partial derivatives dp/da_i $i=1, \dots, m$ and $dp/d\lambda$
- Seek to solve $m+1$ eqns when they are set to 0.
- Note: $dp/d\lambda = -(a^T a - 1)$, so long as constraint is satisfied

See sheet IN NOTES!

- First PC of data is linear combination of the variables that have greatest variance
- Corresponds to taking a linear combination of the variables where the weights are given by the eigenvector λ with largest e -value. This e -value represents the variance of the linear combination.
- All PCs are orthogonal to each other.

- Requires an eigenvalue decomposition of the covariance matrix in order to find the linear combinations of the data variables with greatest variance.
- Each e -value of cov matrix can be interpreted as the variance of a linear combination of the variables with weights given by e -vector.
- Value of a particular e -value divided by the total sum of e -values is the proportion of the variance explained by the associated PC.

- As PCA seeks to maximize variance it can be sensitive to scale differences across variables.
- Standardizing ensures that the data be expressed in comparable units.
- Divide by sample standard deviation \Rightarrow variance equal to one \rightarrow correlation matrix.
- Replace e -value by $-e$.

- Keep adding until a fixed proportion of variance is included.
- Find a kink in the scree plot (variance explained against PC number).
Means that marginal additional variance explained is reducing as a function of PC.
i.e. benefit of including an additional PC may no longer be worth the extra cost of model complexity.

09/05/15 MLA
EXAM PAPER 2014 Q1

DAVID WEBERRECHT

I A. Usefulness of PCA

- Dimension reduction - included dimensions are orthogonal and ordered to account for as much variation as possible.
- Visualisation
- Lower Dimensional Summary.
- Helps identifying relationships in data
- Potentially assists other clustering or classification algorithms (performance)

B. Lagrange Multiplier

$$P = \mathbf{a}^T \mathbf{z} - \lambda (\mathbf{a}^T \mathbf{a} - 1)$$

Solve $dP/d\mathbf{a}_k = 0$ for $k=1 \dots m$

$$dP/d\mathbf{a} = 0 \Rightarrow \mathbf{a}^T \mathbf{a} = 1 \Rightarrow \text{constraint satisfied}$$

$$\Rightarrow \sum_{i=1}^m \mathbf{a}_i \mathbf{b}_i^T - \lambda \mathbf{a} \mathbf{a}^T = 0$$

$$\Rightarrow \mathbf{a}^T \mathbf{a} = 1 \Rightarrow \mathbf{a} \text{ is an eigenvector of } \mathbf{S}$$

C. Standardisation of data

- If not standardised, one component will account for most variation with figure
- more prominently in the solution, swamping out all other values
- PCA will focus on that variable

ii. Number of PCs

- Could select 2 or 3 3-account for more variation (96.9%) but 2 dimension is easier to graph/visualise and interpret.

iii. PC1: high score = high birth rate, high death rate, high ID and small LEM and LEF.

Low score = opposite

Interpreted as health care - independent of GNP

PC2: High score = High loading on GNP \Rightarrow wealth

High score for poor countries, low score for rich countries

C. iv. $(-0.3, -1.1, -0.5, 0.8, 0.9, -0.6)$

$$0.43(-0.3) + 0.37(-1.1) + 0.46(-0.5) - 0.47(0.8) - 0.48(0.9) + 0.08(-0.6) = -1.622$$

$$-0.03(-0.3) + 0.14(-1.1) + 0.06(-0.5) - 0.03(0.8) - 0.04(0.9) - 0.99(-0.6) = 0.282$$

$$(-1.622, 0.282)$$

15

Mr A.

Which linear combination of $\vec{a}^T X$ of the variables in the data has maximum variance, subject to the constraint $\vec{a}^T \vec{a} = 1$?

Lagrange Multiplier

- Given a function $f(x)$, gradient of f , ∇f (the collection of partial derivatives) indicates the direction of steepest slope.
- Level curves (line where $f(x)$ is constant) run perpendicular to the gradient.
- A constraint $g(x) = c$ represents a plane that cuts through the variable space.
- When direction of gradient of constraint equals direction of gradient of f function, then that location constitutes a local maxima.

- Here we work location where $\nabla f = \lambda \nabla g$

\Rightarrow Find x so $\nabla f - \lambda \nabla g = 0$

- Gives a m+1 equation and m+1 unknowns

- Use Lagrange Multiplier to find it. $\vec{a}^T \vec{Q} \vec{a}$ $\vec{a}^T \vec{a} = 1$

- Let $p = \vec{a}^T \vec{Q} \vec{a} - \lambda (\vec{a}^T \vec{a} - 1)$

- Compute partial derivative $\frac{dp}{da_k}$ for $k=1, 2, \dots, m$ and $\frac{dp}{d\lambda}$

- Solve m+1 equations when setting them to 0

- Notice $dp/d\lambda = -(\vec{a}^T \vec{a} - 1)$, guaranteed that constraint is satisfied

- Write $p = \sum_{i=1}^m a_i^2 q_{ii} + \sum_{i=1}^m \sum_{j \neq i} a_i a_j q_{ij} - \lambda (\sum_{i=1}^m a_i^2 - 1)$

$$\frac{dp}{da_k} = 2 a_k q_{kk} + \sum_{j \neq k} a_j q_{kj} + \sum_{i \neq k} a_i q_{ik} - 2 \lambda a_k$$

$$= 2 a_k q_{kk} + 2 \sum_{i \neq k} a_i q_{ik} - 2 \lambda a_k$$

$$\text{So } \frac{dp}{da_k} = 0 \quad \sum_{i=1}^m a_i q_{ik} - \lambda a_k = 0$$

For $k=1, \dots, m$ $\frac{d}{dt} \leq 0$ implies

$$(q_1, q_2, \dots, q_m) \begin{pmatrix} q_{1k} \\ \vdots \\ q_{mk} \end{pmatrix} = \lambda_k q_k$$

we have $Q^T Q = \Lambda Q^T$

Taking transpose of both sides $Q^T Q = Q Q^T = \Lambda Q$

- Thus Q is an eigenvector of Q . (constraint tells it is a unit vector)

- We know the maximum part of the problem is to find the largest eigenvalue

$$Var(a^T x) = a^T Q a = a^T \lambda a = \lambda a^T a = \lambda$$

MLA

15.

INTRODUCTION AND PRINCIPLE COMPONENT ANALYSIS

- Multivariate data arise when two or more attributes are recorded for each set of objects (e.g. height/weight/age). Variables can be continuous, binary, categorical etc.

$$X = \begin{pmatrix} x_{11} & \dots & x_{1m} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nm} \end{pmatrix} \quad \begin{array}{l} n \text{ rows by } m \text{ columns} \\ n - \text{number of observational units} \\ m - \text{number of variables} \end{array}$$

- $\mu = E[X] = \sum x P(X=x)$ if X is discrete random variable with probability mass function $P(X=x)$

- let x_1, \dots, x_m be random variables.

$$\text{Var}[X_i] = E[X_i^2] - [E[X_i]]^2$$

$$\text{Cov}[X_i, X_j] = E[(X_i - E[X_i])(X_j - E[X_j])]$$

$$\text{Corr}[X_i, X_j] = \frac{\text{Cov}[X_i, X_j]}{\sqrt{\text{Var}[X_i] \text{Var}[X_j]}} \quad \text{correlation is "normalized" covariance}$$

$$S = \begin{pmatrix} s_{11} & \dots & s_{1m} \\ \vdots & & \vdots \\ s_{m1} & \dots & s_{mm} \end{pmatrix} \quad \begin{array}{l} \text{covariance matrix with } s_{ij} = \text{Cov}(X_i, X_j) \\ s_{ii} = \text{Var}[X_i] \end{array}$$

- Two random variables X_1, X_2 are independent iff

$$P(X_1=x_1, X_2=x_2) = P(X_1=x_1) P(X_2=x_2) \text{ and so } E[X_1 X_2] = E[X_1] E[X_2]$$

Linear Combination

$$E[aX+b] = aE[X] + b = a\mu + b$$

$$\text{Var}[aX+b] = a^2 \text{Var}[X] = a^2 \sigma^2$$

- If X_1, X_2 independent then

$$E[a_1 X_1 + a_2 X_2] = a_1 E[X_1] + a_2 E[X_2] = a_1 \mu_1 + a_2 \mu_2$$

$$\text{Var}[a_1 X_1 + a_2 X_2] = a_1^2 \text{Var}[X_1] + a_2^2 \text{Var}[X_2] = a_1^2 \sigma_1^2 + a_2^2 \sigma_2^2$$

- If X_1, X_2 not independent then

$$\text{Var}[a_1 X_1 + a_2 X_2] = a_1^2 \sigma_1^2 + a_2^2 \sigma_2^2 + 2a_1 a_2 \text{Cov}[X_1, X_2]$$

$$-a_1x_1 + \dots + a_mx_m = (a_1 \dots a_m) \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix} = a^T X$$

$$E[a^T X] = a^T \mu \quad \text{with } \mu = (\mu_1 \dots \mu_m)$$

$$- \text{Similarly } \text{Var}[a_1x_1 + \dots + a_mx_m] = \text{Var}[a^T X] = a^T \Sigma a$$

- Let A be a $m \times m$ matrix. Then λ is an eigenvalue of A if $\exists v \neq 0$: $Av = \lambda v$. The vector v is said to be an eigenvector of A corresponding eigenvalue λ

- Can find eigenvalues by solving $\det(A - \lambda I) = 0$

- If v is an eigen vector corresponding to eigenvalue λ , then $u = \alpha v$ will also be an eigenvector corresponding to λ .

- v is a unit eigenvector if $\sum_{i=1}^m v_i^2 = v^T v = 1$. We can turn any eigenvector v into a unit eigenvector by multiplying it by $\alpha = 1/\sqrt{v^T v}$

- 2 vectors u, v are orthogonal if $u^T v = v^T u = 0$

- 2 vectors are orthonormal if they are orthogonal and $u^T u = v^T v = 1$

- If λ is an eigenvalue of covariance matrix Σ then $\Sigma v = \lambda v$ where v is the eigenvector corresponding to λ . Hence

$$v^T \Sigma v = v^T \lambda v = \lambda v^T v \quad \text{so}$$

$$\lambda = \frac{v^T \Sigma v}{v^T v}$$

- An $m \times m$ covar matrix Σ has m orthonormal eigenvectors

4/15.

PRINCIPLE COMPONENT ANALYSIS

- For data with many variables/dimensions it is often difficult to comprehend/visualise internal associations - two or more variables could be highly correlated.
- PCA can be thought of as a method for re-expressing the data so as to reveal its internal structure and explain its variation through the use of a few linear combinations of the original values.
- The aim of PCA is to describe the variation of a set of correlated variables x_1, \dots, x_m in terms of a new set of uncorrelated variables y_1, \dots, y_p , hopefully with $p \ll m$ and each y_i a linear combination of x_1, \dots, x_m .
- The new variables or principle components are derived in decreasing order of importance, so that the 1st PC accounts for more variation in the original data than any other possible linear combination of x_1, \dots, x_m .
- The second PC y_2 is chosen to account for a much of the remaining variation as possible subject to the constraint that it be uncorrelated with y_1 and so on.
- Hope is that first few PC's will account for a substantial amount of the variation in the original data and as such can be used as a convenient lower dimension summary of it.
- If we have a dataset of n observations, each consisting of m measurements then the sample mean of each variable is $\bar{x}_i = \frac{1}{n} \sum_{i=1}^n x_{i1}$ and the sample covariance matrix is a matrix whose terms are defined to be
$$q_{ij} = \frac{1}{n-1} \sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)$$
 for $i=1 \dots n$ $j=1 \dots m$
- Then calculate the eigenvalues and corresponding eigenvectors of the sample covariance matrix.

- We use Lagrange multipliers to find the maximum value of $a^T Q a$ subject to $a^T a = 1$. We eventually get $a^T Q = \lambda a^T$, which taking the transpose gives $Qa = \lambda a \Rightarrow a$ is an eigenvector of Q and by the constraint is a unit eigenvector. Also due to maximization part of problem we know it is eigenvector with largest eigenvalue

$$\text{Var}[a^T X] = a^T Q a = a^T \lambda a = \lambda a^T a = \lambda^T = \lambda$$

- The 1st PC of the dataset is the linear combination of the variables that has the greatest variance. This corresponds to taking a linear combination of the variables, where the weights are given by the eigenvector of Q with largest eigenvalue. This eigenvalue also represents the variance of the linear combination.

- The 2nd PC is the linear combination of the variables where the weights are given by the eigenvector of Q corresponding to the 2nd largest eigenvalue and this eigenvalue represents the variance of this linear combination and so on.

- PCA requires an eigenvalue analysis of the covariance matrix in order to find the linear combination of the data variables with greatest variance.

- These linear combinations are called PC's - constructed so they are uncorrelated to each other. PC's have decreasing variance.

- The proportion of the variation in the data variables that is explained by a PC is equal to that component's associated eigenvalue divided by the sum of all eigenvalues. We say that the value of a particular eigenvalue divided by the total sum of eigenvalues is the proportion of the variance explained by the associated PC.

How Do We Interpret PCA Output?

- The 1st row in the table shows the standard deviation of each PC (the standard deviation is square root of e-value obtained with component).

- 2nd row shows proportion of the variance in the data explained by each PC: $SV^2 / \sum SV^2$

PRINCIPLE COMPONENT ANALYSIS (PCA)

How to interpret PCA output out.

- 3rd row show the cumulative proportion of the variance accounted for each PC
- An indicator of the number of PC's required to adequately summarise the data can be inferred by examining the proportion of the variance explained by the PC's
- Consider the new variables (PC) which capture variation of data. PC (column 1, 11) elements are the coefficients/loadings of each original variable on the PC. It matters if the loading has opposite sign but not which is positive and negative. Magnitude of loading important.
- Results differ if scale differs (mm vs. cm). Standardising ensures that the data are expressed in comparable units.
- One way to do this is to make each variable have variance equal to 1. To do this, divide value of each variable by its sample standard deviation. We could then do PCA on covariance matrix of the transformed data. The covariance matrix of a set of variables with variance equal to 1 is a correlation matrix.

How to choose appropriate number of PC's?

- No correct answer \rightarrow rule of thumb
 - Keep adding until a fixed proportion of variance is included
 - Find a kink in scree plot. (variance v. PC numbr!)
- This means that the marginal additional variance explained is reduced as a function of PC. i.e. that the added benefit of including an additional PC may no longer be worth the extra cost of model complexity (remember we are seeking dimension reduction)