

Tratamento nos dados

Redução de Dimensionalidade

PCA

Roseli Ap. F. Romero

Transformação dos dados

- Dados centrados na média e variância 1
 - $\bar{x}_m = 1/n \sum_{i=1}^n x_i$
 - $x_i = (x_i - \bar{x}_m)/\sigma$ (centrados na média)
- Escalamento pela variância (quando temos muito dominante em rel. as demais)
 - $\text{var}_j = 1/(n-1) \sum_{i=1}^n (x_{ij} - \bar{x}_m)^2$
 - $x_{ij}(\text{sv}) = x_{ij} / \text{var}$ desvio = $\sigma = \sqrt{\text{var}}$

PCA

Roseli Ap. F. Romero

Transformação nos dados

- Pre-processamento dos dados
 - log 10 (deve ser usado quando existem grandes disparidades de magnitude nos dados)
 - Normalização
 - $\|x\| = \sqrt{\sum x_{ij}^2}$ para $i, j = 1, \dots$; $x_{ij}(\text{norm}) = x_{ij} / \|x\|$
 - não aconselhável para dados com menos de 10 variáveis
 - Transformação de dados do intervalo [a,b] para [-1,1]:
 $y = (2x - b - a) / (b - a)$

PCA

Roseli Ap. F. Romero

Analise de Componentes Principais - PCA

- Reduz o no. de variáveis
 - $x_1, x_2, \dots, x_p \implies t_1, t_2, \dots, t_m$ ($m < p$)
 - os erros propagados nos valores de w são menores, porque t_1, t_2, \dots, t_m são ORTOGONAIS (sem correlação)

PCA

Roseli Ap. F. Romero

Covariância

- Como é a matriz de covariância?
- Medir c/ var. sobre a média de seus n valores $x_i = 1/n \sum_{i=1}^n x_{ij}$
- $c_{ij} = 1/n \sum_{k=1}^n (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)$
onde c_{ij} é a covariância de x_i e x_j .
- C (p x p)
- A covariância é por vezes chamada de medida de dependência linear entre as duas variáveis aleatórias.

PCA

Roseli Ap. F. Romero

Covariância

$$\begin{pmatrix} 138.32 & -102.12 & -36.20 & -0.94 & -0.14 \\ & 79.74 & 22.38 & 153 & 0.11 \\ & & 13.82 & -0.58 & 0.02 \\ & & & 0.64 & 0.03 \\ & & & & 0.26 \end{pmatrix}$$

PCA

Roseli Ap. F. Romero

PCA

- **Objetivo:** Dadas p variáveis, deseja-se achar combinações lineares dessas, para produzir índices que sejam não correlacionados, de tal forma que:
- Índices Z_i : componentes principais.

PCA

Roseli Ap. F. Romero

PCA

i -ésima componente principal

$$Z_i = a_{i1}X_1 + a_{i2}X_2 + \dots + a_{ip}X_p$$

Com restrição: $a_{i1}^2 + a_{i2}^2 + \dots + a_{ip}^2 = 1$

$Z_1, Z_2, \dots, Z_{i-1}, Z_i$ Não correlacionados

PCA

Roseli Ap. F. Romero

PCA

– PCA: resume-se em encontrar os autovalores e autovetores da matriz C de covariância dos dados

$$- C = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & \vdots & & \vdots \\ c_{p1} & c_{p2} & \dots & c_{pp} \end{bmatrix}$$

PCA

Roseli Ap. F. Romero

PCA

- Supondo que: os autovalores da matriz C estejam ordenados da seguinte forma:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_j \geq \dots \geq \lambda_p$$

- Os auto-vetores associados:

$$a_1, a_2, \dots, a_j, \dots, a_p$$

PCA

Roseli Ap. F. Romero

PCA

- Propriedades

$$\bullet \text{ Para } a_i^T a_j = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

$$Z_i = a_{i1}X_1 + a_{i2}X_2 + \dots + a_{ip}X_p$$
$$a_{i1}, a_{i2}, \dots, a_{ip}$$

- são os elementos do i -ésimo autovetor correspondente

PCA

Roseli Ap. F. Romero

PCA

- a soma dos auto-valores corresponde ao **traço da matriz covariância C**

$$\lambda_1 + \lambda_2 + \dots + \lambda_p = c_{11} + c_{22} + \dots + c_{pp}$$

- $\text{Var}(Z_i) = \lambda_i$

PCA

Roseli Ap. F. Romero

Exemplo:

- Determinar os auto-valores e auto-vetores de:

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$
$$\det(A - \lambda I) = \begin{vmatrix} 2-\lambda & 1 \\ 1 & 2-\lambda \end{vmatrix} = (2-\lambda)^2 - 3 = 0$$

- Sol: $\lambda_1=3$ e $\lambda_2=1$

PCA

Roseli Ap. F. Romero

Exemplo

- Determinar o auto-vetor assoc. a $\lambda_1=3$

$$-v_1 + v_2 = 0$$

$$v_1 - v_2 = 0$$

$$\text{Sol: } v_1 = v_2$$

$$\text{Portanto } v = [1 \ 1]^T$$

$$\text{Analogamente, } v = [1 \ -1]^T \text{ para } \lambda_2=1$$

PCA

Roseli Ap. F. Romero

Operações em PCA

- Det. dos auto-valores e auto-vetores: calculo do determinante e sol. de um sistema de eq. lineares.

PCA

Roseli Ap. F. Romero

Exercicio

- Implementar a técnica PCA (use o Matlab ou algum pacote estatístico, ou 1 ling. de progr.) e verifique as 2 comp. Principais obtidas para o conj. iris

PCA

Roseli Ap. F. Romero

Reconstrução dos dados originais

$$Z = [Z_1, Z_2, \dots, Z_p]^T$$

$$= [X^T a_1, X^T a_2, \dots, X^T a_{p-1}]^T$$

$$= A^T X$$

$$X = A Z = \sum_{i=1}^p Z_i a_i \quad A^T = A^{-1}$$

PCA

Roseli Ap. F. Romero

Redução da Dimensionalidade

- Sejam $\lambda_1, \lambda_2, \dots, \lambda_m$
- os m auto-valores da matriz C
- Então $X' \sim X$ onde:

$$X' = \sum_{i=1}^m Z_i a_i \quad m < p$$

- o erro: $e = X - X'$ onde:

$$e = \sum_{i=m+1}^p Z_i a_i$$

PCA

Roseli Ap. F. Romero

- O vetor de erro e é ortogonal ao vetor X' , que aproxima X .
Verifique isto!!!
A equação: $e^T X' = 0$: princípio da ortogonalidade
- Existe uma rede neural que implementa PCA, proposta por RUBN-89 (vamos ver!!!)

PCA

Roseli Ap. F. Romero

• Algoritmo da PCA

1. Padronize os dados correspondentes às variáveis para que estes tenham média igual a 0 e variância igual a 1.
2. Calcule a matriz de covariância C .
3. Encontre os autovalores da matriz C e seus correspondentes auto-vetores. Os coeficientes da i -ésima componente principal são dados pelo auto-vetor associado ao i -ésimo auto-valor

PCA

Roseli Ap. F. Romero

4. Descarte as componentes que acumulem uma pequena proporção da variação dos dados. Por exemplo, se os dados originais tiverem 20 variáveis e as três primeiras componentes principais forem responsáveis por 90% do total da variância, as outras 17 componentes principais podem ser ignoradas.

PCA

Roseli Ap. F. Romero

Algoritmo PCA - Matlab

```

%%%
% Componentes Principais
% PCA - Roseli Romero
%%%
% carregamento do arquivo iris.dat
load iris.txt
% verificando a dimensão do conjunto
[n,p] = size(iris);
% tirando a primeira coluna da matriz de dados
iris(:,1) = [];
p = p-1;
X = iris; % matriz X contém os dados
% centrando os dados na média 0
S = std(X); % S armazena os desvios-padrão de cada coluna de X
M = mean(X); % M armazena as médias de cada coluna de X
X = X - ones(n,1) * M;
% transformando os dados para ter variância 1
X = X ./ (ones(n,1) * S);
% calculando a matriz de covariância C
C = (X'*X)/n;
% auto-valores (A) e auto-vetores (V) da matriz de covariância
[V,A] = eig(C);
A = diag(A);
PCA

```

PCA

Roseli Ap. F. Romero

```

% ordenando auto-valores e auto-vetores por ordem crescente de auto-valores
V = V'; % coloca auto-vetores nas linhas de V
A = [A,V]; % concatena A e V (cada linha de A contém um auto-valor e auto-vetor)
A = sortrows(A,1); % ordena os auto-valores em ordem crescente
A = A(:,1) % separa auto-valores em A e imprime na tela
V % imprime V na tela
% calculando as componentes principais
Z = [];
for i=1:p
    Vi = V(p+1-i,:); % pega os auto-vetores em ordem decrescente dos auto-valores
    Zi(:,i) = X * Vi; % obtém-se a i-ésima componente principal
end
% salvando as componentes principais em um arquivo
save('componentes_iris.txt')
Z1=Z(:, 1:2);
L1=Z1(1:50,:);
L2=Z1(51:100,:);
L3=Z1(101:150,:);
% fazendo o gráfico das duas primeiras componentes
plot(Z1,Z2,'-')
%fazendo o gráfico das tres classes separadas
plot(L1(:,1),L1(:,2),'L1';L2(:,1),L2(:,2),'L2';L3(:,1),L3(:,2),'L3')

```

PCA

Roseli Ap. F. Romero

- Implementar a técnica PCA (use o Matlab ou algum pacote estatístico, ou 1 ling. de progr.) e verifique as 4 comp. Principais obtidas para o conj. iris

PCA

Roseli Ap. F. Romero

- Dada a matriz:

$$X = \begin{pmatrix} 3 & 0 & 1 & 5 \\ 4 & 7 & 4 & 3 \\ 3 & 2 & 6 & 2 \\ 5 & 3 & 6 & 4 \end{pmatrix}$$

- transformar essa matriz:
dados centrados pela média
escalamento pela variância

PCA

Roseli Ap. F. Romero