

UNIVERSIDADE FEDERAL DE MINAS GERAIS

Graduação em Engenharia de Sistemas Introdução ao Reconhecimento de Padrões - Aula 15 Nikolas Dias Magalhães Fantoni - 2018019400

Análise de Componentes Principais (PCA)

Exercício

No exercício proposto, foram obtidos dados contendo 400 fotos dos rostos de 40 pessoas diferentes, sendo 10 de cada, em posições ligeiramente diferentes. Cada foto possui 4096 pixels (64x64), sendo que cada pixel pode assumir um valor de 0 a 255, correspondendo a uma escala de preto e branco. Assim, cada pixel possui 1 byte de tamanho, e cada foto aproximadamente 4KB, sendo a base inteira composta por 800KB de dados. A biblioteca que gerou estas imagens é a RnavGraphImageData do R.

As fotos foram consideradas dados de entrada e foram divididas em 4 folds, seguindo então a proporção percentual de 75/25 em relação aos dados para treinamento e teste, respectivamente. Cada pixel tornouse então uma variável diferente do problema.

Devido ao alto número de variáveis envolvidas, a técnica PCA foi implementada para definir as variáveis que possuiam maior variância nos dados. Assim, mesmo que com uma certa perda de valores, foi possível analisar apenas uma quantidade considerada pequena de pixels para definir a quem os rostos do grupo de teste pertenciam. Para isso, cada pessoa foi considerada uma classe e, com isso, o problema tornou-se um classificador para definir em qual das 40 classes (pessoas) um certo dado de teste (foto) pertencia.

Para a técnica PCA, primeiro foram calculados os autovalores e os autovetores utilizando funções próprias do R. Na figura 1 são mostrados os maiores autovalores, em ordem decrescente, das variáveis. As 3 primeiras variáveis correspondem a aproximadamente 46% da variância total do problema, enquanto as 8 primeiras correspondem à aproximadamente 70%. Porém, para garantir uma boa acurácia da classificação, foram utilizados dados que correspondiam a mais de 98% da variância total do problema.

Foram feitas 10 simulações, dividindo-se os dados em 75/25 para cada iteração. O número médio de pixels analisados nas divisões foi de 109, uma redução de pouco menos de 36 vezes na quantidade de valores. Assim, foram analisados e testados apenas 43KB em média de dados, contra os 800KB totais.

Quanto à acurácia, a média dos testes foi de 75,96% e o desvio padrão dos resultados foi de 5,05%. A tabela das acurácias obtidas em cada iteração é mostrada na tabela 1. Já a matriz de confusão é mostrada na tabela 2.

Iteração	Acurácia
1	82,83%
2	72%
3	80,20%
4	72%
5	74%
6	67%
7	76,24%
8	$75,\!25\%$
9	82,83%
10	77,23%
Média	$75{,}96\%$

Tabela 1

Mesmo utilizando as variáveis correspondentes a 98% da variância total do problema, a acurácia média manteve-se baixa em termos relativos. Isso se explica pois ao dividir o conjunto de fotos em dados para o teste e para o treinamento, as fotos com alta variação (como quando uma pessoa tirou fotos com óculos e sem óculos, ou quando uma pessoa afastou-se demais da câmera ou mesmo inclinou demais a cabeça

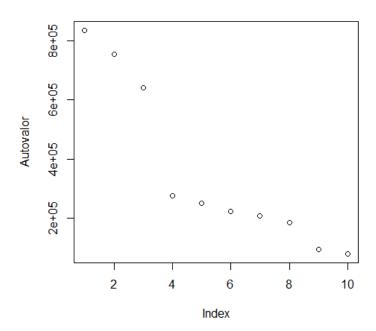


Figura 1: Autovalores dominantes dos dados de entrada.

			Previsto																																					
	Classes	1	2	3	4	5	6	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40
	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	2	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	3	0	0	3	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	4	0	0	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	5	0	0	0	0	3	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	6	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	8	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	9	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	11	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	12	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Obtido	13	0	0	0	0	0	0	0	0	0	0	0	2	1	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	14	0	0	0	0	0	0	0	0	0	0	0	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	15	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	16	0	0	0	0	0	0	0	0	0	0	0	0	0	1	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Obtido	21	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	4	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
	22	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	3	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
	23	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
	24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	25	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	26	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
	28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	3	0	0	0	0	0	0	0	0	0	0	0	0
	29	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	30	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0	0	0	0	0	0	1	0	0	0
	31	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0
	32	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0
	33	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0
	34	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
	35	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	3	0	0	0	0	0
	36	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	2	0	0	0	0
	37	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0
	38	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	3	0	0
	39	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	2	0
	40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	4

Tabela 2: Matriz Confusão

para os lados) podem ter influenciado no treinamento e, mais ainda, na classificação. Ou seja, como a distribuição foi aleatória, o treinamento pode ter sido feito com fotos mais uniformes enquanto as fotos separadas para o teste podem diferir muito das originais. O ideal para melhorar a acurácia do problema seria um conjunto muito maior de fotos, para que a classificação se desse de forma mais precisa ainda, ou uma seleção não estocástica (mas sim analítica) das fotos usadas para o treinamento do algoritmo classificador.

Por fim, a técnica do PCA mostrou-se muito efetiva em reduzir a dimensão do problema, economizando recursos computacionais e tempo ao fazer uma classificação, uma vez que só considera as variáveis de maior impacto na classificação e elimina as que possuem pouca variabilidade.