

Projeto AM 2018-1

Francisco de A. T. de Carvalho¹

1 Centro de Informatica-CIn/UFPE
Av. Prof. Luiz Freire, s/n -Cidade Universitaria, CEP 50740-540, Recife-PE, Brasil,
fatc@cin.ufpe.br

- 1) No conjunto de dados "Image Segmentation" do site uci machine learning repository considere a tabela de dados segmentation.test (<http://archive.ics.uci.edu/ml/machine-learning-databases/image>). Essa tabela de dados contém 2100 objetos e 7 classes. Os objetos são descritos por 19 variáveis que podem ser divididas em 2 views:

- Shape view: as primeiras 9 variáveis
- RGB view: as 10 últimas variáveis

Execute a variante KCM-K-GH do algoritmo KCM-K-H descrito na seção 3.1 do artigo "FAT de Carvalho, EC Simões, LVC Santana, MRP Ferreira, Gaussian Kernel C-Means Hard Clustering Algorithms with Automated Computation of the Width Hyper-Parameters, Pattern Recognition, 79, 370-386, 2018" na tabela de dados completa (complet view, 2100 objetos e 19 variáveis), na tabela shape view (2100 objetos e 9 variáveis) e na tabela RGB view (2100 objetos e 10 variáveis), 100 vezes para obter uma partição em 7 grupos. Em cada caso selecione o melhor resultado segundo a função objetivo. Em cada caso, calcule o índice de Rand corrigido.

Observações:

- No algoritmo 1, página 374 da seção 3.1, os representantes dos grupos são calculados segundo a equação (14), o vetor de hiperparâmetros é calculado com a equação (16), a afetação dos objetos aos grupos é realizada segundo a equação (18);
- Parâmetros: número de grupos $c = 7$; parâmetro $\gamma = (\frac{1}{\sigma^2})^p$ onde, p é o número de variáveis e σ^2 é a média entre o 0.1 e o 0.9 quantil de $\|\mathbf{x}_l - \mathbf{x}_k\| \neq k$;
- Para o melhor resultado obtido para cada conjunto de dados imprimir: i) o representante de cada grupo, ii) o número de objetos de cada grupo, iii) o vetor de hiperparâmetros, iv) a partição (para cada grupo, a lista de objetos), v) o índice de Rand corrigido.

- 2) Considere novamente a tabela de dados "Image Segmentation". Os exemplos são rotulados segundo as classes "brickface", "sky", "foliage", "cement", "window", "path", "grass".
- Use validação cruzada estratificada "30 times ten fold" para avaliar e comparar os classificadores descritos abaixo. Se necessário, retire do conjunto de aprendizagem, um conjunto de validação para fazer ajuste de parâmetros e depois treine o modelo novamente com os conjuntos aprendizagem + validação.
 - Obtenha uma estimativa pontual e um intervalo de confiança para a taxa de acerto de cada classificador;
 - Usar Friedman test (teste não paramétrico) para comparar os classificadores. Se necessário, usar também o Nemenyi test (pos teste);

Considere os seguintes classificadores:

- Classificador bayesiano gaussiano. Considere a seguinte regra de decisão: afetar o exemplo \mathbf{x}_k à classe ω_i se $P(\omega_i|\mathbf{x}_k) = \max_{i=1}^7 P(\omega_i|\mathbf{x}_k)$ com $P(\omega_i|\mathbf{x}_k) = \frac{p(\mathbf{x}_k|\omega_i)P(\omega_i)}{\sum_{r=1}^C p(\mathbf{x}_k|\omega_r)P(\omega_r)}$
 - Estime $P(\omega_i)$ pelo método de máxima verossimilhança.
 - Para cada classe ω_i ($i = 1, 2, 3$) estime $p(\mathbf{x}_k|\omega_i) = p(\mathbf{x}_k|\omega_i, \theta_i)$ pelo método da máxima verossimilhança, supondo uma normal multivariada, onde:
 - $\theta_i = \begin{pmatrix} \mu_i \\ \Sigma \end{pmatrix}$
 - $p(\mathbf{x}_k|\omega_i, \theta_i) = (2\pi)^{-\frac{d}{2}} (|\Sigma^{-1}|)^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_k - \mu_i)^{tr} \Sigma^{-1} (\mathbf{x}_k - \mu_i) \right\}$
 - $\mu_i = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$,
 - $\Sigma = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \mu_i)(\mathbf{x}_k - \mu_i)^{tr} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k(\mathbf{x}_k)^{tr} - n \mu_i(\mu_i)^{tr}$

- ii) Usar um classificador bayesiano baseado em janela de Parzen para fazer a classificação dos dados. Treine três classificadores bayesianos baseados em janela de Parzen , um para cada view. Use a função de kernel multivariada produto com um mesmo h para todas as dimensões e a função de kernel unidimensional. Use conjunto de validação para fixar o parametro h .

- iii) Regra da soma: afetar o exemplo \mathbf{x}_k a classe ω_j se

$$(1 - L)P(\omega_j) + P_{GAUSS, VIEW1}(\omega_j|\mathbf{x}_k) + P_{GAUSS, VIEW2}(\omega_j|\mathbf{x}_k) + P_{GAUSS, VIEW3}(\omega_j|\mathbf{x}_k) + \\ P_{PARZEN, VIEW1}(\omega_j|\mathbf{x}_k) + P_{PARZEN, VIEW2}(\omega_j|\mathbf{x}_k) + P_{PARZEN, VIEW3}(\omega_j|\mathbf{x}_k) = \\ \max_{r=1}^7 ((1 - L)P(\omega_r) + P_{GAUSS, VIEW1}(\omega_r|\mathbf{x}_k) + P_{GAUSS, VIEW2}(\omega_r|\mathbf{x}_k) + P_{GAUSS, VIEW3}(\omega_r|\mathbf{x}_k) + \\ P_{PARZEN, VIEW1}(\omega_r|\mathbf{x}_k) + P_{PARZEN, VIEW2}(\omega_r|\mathbf{x}_k) + P_{PARZEN, VIEW3}(\omega_r|\mathbf{x}_k))$$

com $L = 3$ (três views: complete view, shape view, RGB view)

Observações Finais

- No Relatório e na saída da ferramenta devem estar bem claros:
 - a) como foram organizados os experimentos de tal forma a realizar corretamente a avaliação dos modelos e a comparação entre os mesmos.
Fornecer também uma descrição dos dados.
- Data de apresentação e entrega do projeto: TERÇA-FEIRA 12/06/2018
- Enviar por email : o programa fonte, o executável (se houver), os dados e o relatório do projeto
- Tempo de apresentação: 10 minutos (rigoroso).
- PASSAR NA MINHA SALA PARA ASSINAR A ATA DE ENTREGA DO TRABALHO EM 12/06/2018
- O PROJETO DEVE SER REALIZADO COM 3 ALUNOS.