

Linear Discriminant Analysis - LDA

Gustavo Aparecido de Souza Viana

Abstract—A área de machine learning vem crescendo exponencialmente, com intuito de desenvolver algoritmos ou aplicações para analisar grande bases de dados e assim tirar conclusões. Partindo desses principio, o LDA foi criado, depois que o PCA, com objetivo de classificar dados. Atualmente, nos deparamos com alguns trabalhos que utilizam o PCA como redutor de dimensões antes de aplicar o LDA, ou seja, PCA como otimizador do LDA. O objetivo deste trabalho é utilizar o PCA como otimizador do LDA. Os resultados mostraram que o é possível aplicar o PCA para reduzir dimensões e em seguida aplicar o LDA.

I. INTRODUÇÃO

A necessidade de classificações está sendo bem comum na atualidade, podemos dar o exemplo de uma empresa que necessita prever qual estado da máquina (bom, normal, ou ruim) e assim com base no histórico dela e com base nos dados atuais, conseguimos classificar a mesma.

Conseguimos prever alguma informação baseada em uma base conhecida com machine learning ou com métodos estatísticos. É de extrema importância conhecer o ramo de atividade da aplicação pois métodos de machine learning geralmente estão relacionados a rede neural, consequentemente necessitam de mais dados para que as previsões tenham uma maior acurácia. De contra partida, métodos estatísticos não necessitam.

Neste trabalho foi abordado a implementação do *Linear Discriminant Analysis* com objetivo de classificar a base de dados "Iris" mas aplicando o PCA previamente com intuito de deixar somente as dimensões necessárias.

II. CONCEITOS FUNDAMENTAIS

Nesta seção, será apresentado os conceitos básicos utilizados para o *Linear Discriminant Analysis* e para o *Principal Components Analysis* que inclui a Covariância, e Autovalores e Autovetores.

A. Covariância

A covariância na estatística, ou variância conjunta, é uma medida do grau de interdependência numérica entre duas variáveis aleatórias [2]. Ela pode ser positiva e negativa ou zero, positiva quando a variação tem o mesmo sentido, negativa quando são sentidos opostos e zero quando não há variação. A Equação 1 representa o cálculo entre x e y , onde \bar{x} representa a média da sua própria população, e N a quantidade de observações.

$$\sqrt{\sum \frac{(x_i - \bar{x})(y_i - \bar{y})}{N - 1}} \quad (1)$$

B. Autovalor e Autovetor

Autovetor e autovalor, são vetores e valores respectivamente de uma matriz onde suas multiplicações não alteram a direção, apenas a magnitude [4]. Sendo assim, dado um λ , sendo um valor escalar, multiplicado por um vetor X de uma matriz quadrada A , podemos dizer que λ é autovalor da matriz A caso exista um vetor $\neq 0$ tal que $Ax = \lambda x$, sendo assim o vetor x é chamado de autovetor. Abaixo a Equação 2 representa o conceito do autovetor e valor, e a Equação 3 representa um exemplo de autovalor e autovetor.

$$Matriz_a Vetor_i = \lambda Vetor_i \quad (2)$$

$$\begin{vmatrix} 2 & 3 \\ 2 & 1 \end{vmatrix} \begin{vmatrix} 3 \\ 2 \end{vmatrix} = 4 \begin{vmatrix} 3 \\ 2 \end{vmatrix} \quad (3)$$

C. Linear Discriminant Analysis

A ideia do LDA é maximizar a distância entre as classes de uma base de dados, para assim facilitar a classificação de uma predição [3].

O processo de classificação inicia pelo cálculo do *between class* e *with-in class*. Após isso calculamos os autovetores e autovalores. Os autovalores demonstram o quanto aquela componente representa ou explica do dado em si, quanto maior o autovalor mais importante e consequentemente, explica majoritariamente o dado. O autovetor é a componente em si, uma matriz de transformação, ou seja, se multiplicarmos os dados pelo autovetor teremos uma nova base maximizada entre as classes e minimizada dentro das classes.

D. Principal Component Analysis

A ideia do PCA é reduzir a dimensionalidade dos dados para assim representar o mesmo dado com menos dimensões que são novas variáveis linearmente independentes chamadas de componentes [1].

O processo de conversão inicia pelo cálculo da covariância, e após isso conseguimos calcular os autovetores e autovalores. Os autovalores demonstram o quanto aquela componente representa ou explica do dado em si, quanto maior o autovalor mais importante e consequentemente, explica majoritariamente o dado. O autovetor é a componente em si, uma matriz de transformação, ou seja, se multiplicarmos os dados pelo autovetor teremos uma nova base reduzida. Importante ressaltar que cada componente é ortogonal a anterior e assim sucessivamente.

III. METODOLOGIA

Nesta seção será apresentado a metodologia utilizada para implementar o *Linear Discriminant Analysis* juntamente com o *Principal Analysis Component* para reduzir as dimensões desnecessárias, implementados em Python. O código fonte pode ser encontrado em <https://github.com/apparecido/master-special-learning-topic>.

A metodologia consiste em 5 passos. Iniciando com a leitura da base de dados citada na introdução, em seguida é aplicado o PCA para reduzir a dimensionalidade. No terceiro passo inicia o processo do LDA, calculando o *between e within class* baseada na nova base de dados retornada pelo PCA. No quarto passo é calculado o auto valor e auto vetor do resultado de $Matrix_{within}^{-1} * Matrix_{between}$, respectivamente. Por fim no último passo as componentes são aplicadas na base original à fim de transformá-la.

IV. EXPERIMENTOS E RESULTADOS

Nesta seção será apresentado os experimentos e seus respectivos resultados.

O experimento consiste na aplicação do *Principal Component Analysis* e em seguida a aplicação do LDA na base de dados "Iris".

A Figura 1 representa o gráfico da base Iris sem a aplicação do LDA, 2 representa o gráfico da base Iris com a aplicação do LDA e 3 representa o gráfico da base Iris com a aplicação do PCA e em seguida LDA.

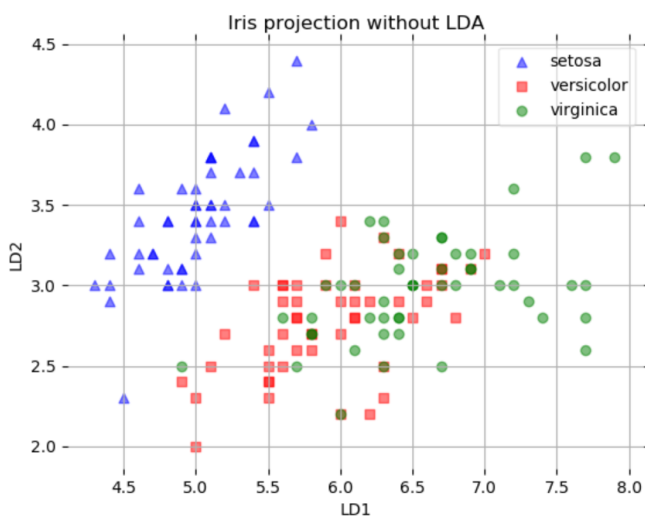


Fig. 1: Representação da base de dados Iris sem LDA

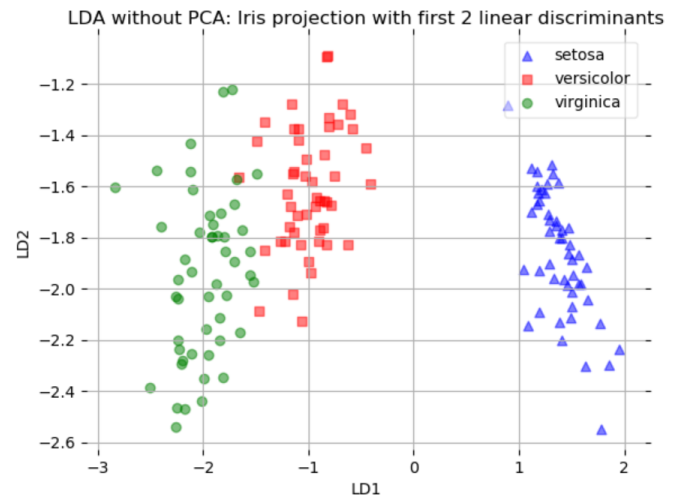


Fig. 2: Representação da base de dados Iris com LDA

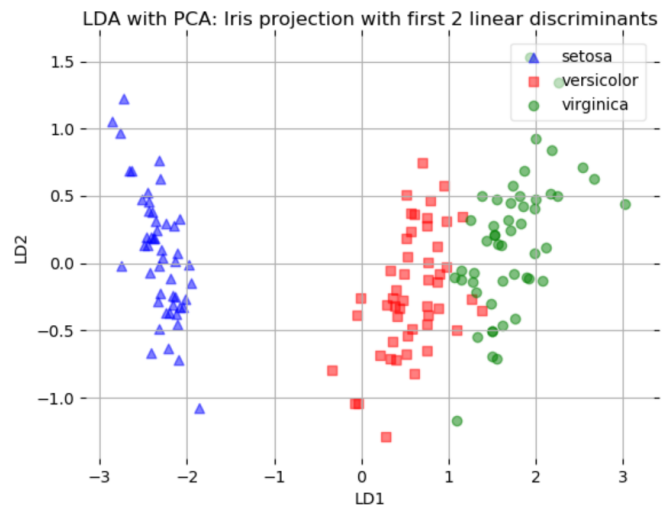


Fig. 3: Representação da base de dados Iris com PCA e LDA

V. CONCLUSÃO

Após a análise dos conceitos e resultados obtidos após a aplicação dos métodos *PCA* e *LDA*, e comparando os mesmos, podemos concluir que o *PCA* consegue reduzir a dimensionalidade, e além disso podemos representar o mesmo dado da base com menos características. E nesse caso a aplicação do PCA antes do LDA não surtiu tanta diferença se analisarmos os gráficos.

REFERENCES

- [1] M. Baxter and M. Heyworth. Principal components analysis of compositional data in archaeology. 1989.
- [2] P. J. Bickel and E. Levina. Covariance regularization by thresholding. 2008.
- [3] C. Eckart and G. Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1:211–218, 1936.
- [4] E. Oja and J. Karhunen. On stochastic approximation of the eigenvectors and eigenvalues of the expectation of a random matrix. 1985.