

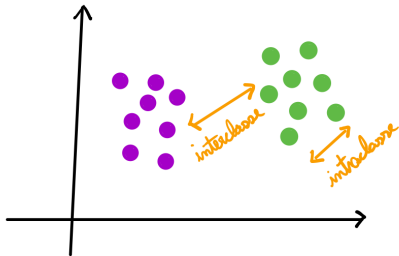
Análise Discriminante Linear

Advanced Institute for Artificial Intelligence – AI2

<https://advancedinstitute.ai>

Introdução

Análise Discriminante Linear, do inglês *Linear Discriminant Analysis* - LDA, é uma técnica **supervisionada** que objetiva **maximizar a separabilidade** entre as classes. A ideia é gerar agrupamentos em que a distância intraclasses (elementos de mesma classe) seja pequena, enquanto que a distância interclasses (elementos de classes distintas) seja grande. **Queremos, então, gerar agrupamentos compactos e distantes.**



Seja $\mathcal{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_z, y_z)\}$ um conjunto de dados rotulado tal que $\mathbf{x}_i \in \mathbb{R}^n$. Temos que $\mathcal{X}_1 \subset \mathcal{X}$ e $\mathcal{X}_2 \subset \mathcal{X}$ representam duas partições que denotam os conjuntos de treinamento e teste, respectivamente. Ademais, seja $\mathcal{Y} = \{\omega_1, \omega_2, \dots, \omega_c\}$ o conjunto de rótulos possíveis. Em nosso exemplo, vamos assumir que temos apenas duas classes, isto é, ω_1 e ω_2 .

Podemos calcular as médias de cada agrupamento da seguinte forma:

$$\boldsymbol{\mu}_1 = \frac{1}{m_1} \sum_{\mathbf{x}_i \in \omega_1} \mathbf{x}_i \quad \text{e} \quad \boldsymbol{\mu}_2 = \frac{1}{m_2} \sum_{\mathbf{x}_j \in \omega_2} \mathbf{x}_j,$$

em que m_1 e m_2 correspondem ao número de amostras de treinamento das classes 1 e 2, respectivamente. Já $\boldsymbol{\mu}_1$ e $\boldsymbol{\mu}_2$ denotam os centros dos agrupamentos das amostras das classes 1 e 2, respectivamente.

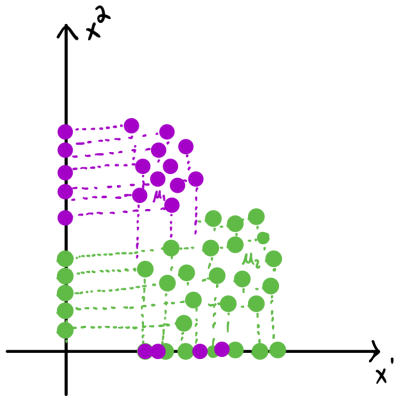
Objetivo: dado um conjunto de treinamento \mathcal{X}_1 , achar um vetor direção $\mathbf{w} \in \mathbb{R}^n$ de tal forma que, quando projetarmos nossas amostras nesta direção, maximizaremos a sua separabilidade.

Sejam $\hat{\mu}_1$ e $\hat{\mu}_2$ as médias dos dados das classes 1 e 2, respectivamente, projetados na direção \mathbf{w} :

$$\begin{aligned}\hat{\mu}_1 &= \frac{1}{m_1} \sum_{\mathbf{x}_i \in \omega_1} \mathbf{w}^T \mathbf{x}_i \\ &= \underbrace{\mathbf{w}^T}_{\text{constante}} \left(\frac{1}{m_1} \sum_{\mathbf{x}_i \in \omega_1} \mathbf{x}_i \right) \\ &= \mathbf{w}^T \mu_1.\end{aligned}\quad (1)$$

$$\begin{aligned}\hat{\mu}_2 &= \frac{1}{m_2} \sum_{\mathbf{x}_j \in \omega_2} \mathbf{w}^T \mathbf{x}_j \\ &= \underbrace{\mathbf{w}^T}_{\text{constante}} \left(\frac{1}{m_2} \sum_{\mathbf{x}_j \in \omega_2} \mathbf{x}_j \right) \\ &= \mathbf{w}^T \mu_2.\end{aligned}\quad (2)$$

Um critério interessante seria maximizar a diferença entre as médias projetadas, ou seja, $|\hat{\mu}_1 - \hat{\mu}_2|$:



Neste exemplo, a distância entre as médias é maior no eixo x^1 , mas a separabilidade é maior no eixo x^2 .

Por que isto ocorre? A abordagem acima não leva em consideração a variância, ou seja, o **espalhamento** entre as classe.

Dado o nosso conjunto de dados $\mathcal{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_z, y_z)\}$, temos que sua média amostral é dada por:

$$\boldsymbol{\mu} = \frac{1}{z} \sum_{i=1}^z \mathbf{x}_i, \quad (3)$$

e o seu espalhamento (*scatter*) é calculado como segue:

$$s^2 = \sum_{i=1}^z (\mathbf{x}_i - \boldsymbol{\mu})^2. \quad (4)$$

Ademais, seja $\hat{\mathbf{x}}_i = \mathbf{w}^T \mathbf{x}_i$ a projeção da amostra \mathbf{x}_i na direção do vetor \mathbf{w} .

Temos que os espalhamentos no espaço projetado para as classes ω_1 e ω_2 são definidos como segue:

$$\hat{s}_1^2 = \sum_{\hat{\mathbf{x}}_i \in \omega_1} (\hat{\mathbf{x}}_i - \hat{\mu}_1)^2, \quad (5)$$

e

$$\hat{s}_2^2 = \sum_{\hat{\mathbf{x}}_j \in \omega_2} (\hat{\mathbf{x}}_j - \hat{\mu}_2)^2. \quad (6)$$

Assim, desejamos **maximizar a distância entre as médias e minimizar o espalhamento entre as classes**. O Critério de Fisher atende à esta nossa exigência.

Temos, então, que maximizar o Critério de Fisher, dado por:

$$J(\mathbf{w}) = \frac{(\hat{\mu}_1 - \hat{\mu}_2)^2}{\hat{s}_1^2 + \hat{s}_2^2}. \quad (7)$$

Desta forma, queremos encontrar o vetor \mathbf{w} que maximiza o critério, acima, ou seja:

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} J(\mathbf{w}). \quad (8)$$

No entanto, precisamos reescrever tanto o numerador quanto o denominador em termos de \mathbf{w} .

Antes da projeção, podemos obter as matrizes de espalhamento (*scatter*) das classes ω_1 e ω_2 como segue:

$$S_1 = \sum_{\mathbf{x}_i \in \omega_1} (\mathbf{x}_i - \boldsymbol{\mu}_1)(\mathbf{x}_i - \boldsymbol{\mu}_1)^T \quad (9)$$

e

$$S_2 = \sum_{\mathbf{x}_j \in \omega_2} (\mathbf{x}_j - \boldsymbol{\mu}_2)(\mathbf{x}_j - \boldsymbol{\mu}_2)^T. \quad (10)$$

Seja $S_A = S_1 + S_2$ a matriz de espalhamento **intraclasse**. Podemos reescrever a Equação 5 como segue:

$$\hat{s}_1^2 = \sum_{\hat{\mathbf{x}}_i \in \omega_1} (\hat{\mathbf{x}}_i - \hat{\boldsymbol{\mu}}_1)^2 = \sum_{\mathbf{x}_i \in \omega_1} (\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \boldsymbol{\mu}_1)^2 = \sum_{\mathbf{x}_i \in \omega_1} [\mathbf{w}^T (\mathbf{x}_i - \boldsymbol{\mu}_1)]^T [\mathbf{w}^T (\mathbf{x}_i - \boldsymbol{\mu}_1)]. \quad (11)$$

Rearranjando os termos, temos que:

$$\begin{aligned}\hat{s}_1^2 &= \sum_{\mathbf{x}_i \in \omega_1} [\mathbf{w}^T(\mathbf{x}_i - \boldsymbol{\mu}_1)]^T [\mathbf{w}^T(\mathbf{x}_i - \boldsymbol{\mu}_1)] \\ &= \sum_{\mathbf{x}_i \in \omega_1} [(\mathbf{x}_i - \boldsymbol{\mu}_1)^T \mathbf{w}]^T [(\mathbf{x}_i - \boldsymbol{\mu}_1)^T \mathbf{w}] \\ &= \sum_{\mathbf{x}_i \in \omega_1} \mathbf{w}^T (\mathbf{x}_i - \boldsymbol{\mu}_1) (\mathbf{x}_i - \boldsymbol{\mu}_1)^T \mathbf{w} \\ &= \mathbf{w}^T \underbrace{\sum_{\mathbf{x}_i \in \omega_1} (\mathbf{x}_i - \boldsymbol{\mu}_1) (\mathbf{x}_i - \boldsymbol{\mu}_1)^T}_{S_1} \mathbf{w} \\ &= \mathbf{w}^T S_1 \mathbf{w}.\end{aligned}\tag{12}$$

De maneira análoga, podemos escrever:

$$\hat{s}_2^2 = \mathbf{w}^T \mathbf{S}_2 \mathbf{w}. \quad (13)$$

Substituindo-se as Equações 12 e 13 na Equação 7, temos que:

$$J(\mathbf{w}) = \frac{(\hat{\mu}_1 - \hat{\mu}_2)^2}{\hat{s}_1^2 + \hat{s}_2^2} = \frac{(\hat{\mu}_1 - \hat{\mu}_2)^2}{\mathbf{w}^T \mathbf{S}_1 \mathbf{w} + \mathbf{w}^T \mathbf{S}_2 \mathbf{w}} = \frac{(\hat{\mu}_1 - \hat{\mu}_2)^2}{\mathbf{w}^T \underbrace{(\mathbf{S}_1 + \mathbf{S}_2)}_{\mathbf{S}_A} \mathbf{w}} = \frac{(\hat{\mu}_1 - \hat{\mu}_2)^2}{\mathbf{w}^T \mathbf{S}_A \mathbf{w}}, \quad (14)$$

a qual é uma forma quadrática em \mathbf{w} .

Seja, agora, S_B a matriz de espalhamento interclasses, que pode ser definida como segue:

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T, \quad (15)$$

a qual mede a separação entre os vetores média antes da projeção. Note que o numerador da Equação 14 pode ser expressado como:

$$(\hat{\mu}_1 - \hat{\mu}_2)^2 = (w^T \mu_1 - w^T \mu_2)^2 = [w^T (\mu_1 - \mu_2)]^T [w^T (\mu_1 - \mu_2)]. \quad (16)$$

Rearranjando os produtos internos, temos que:

$$(\hat{\mu}_1 - \hat{\mu}_2)^2 = [(\mu_1 - \mu_2)^T w]^T [(\mu_1 - \mu_2)^T w] = w^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T w = w^T S_B w. \quad (17)$$

Podemos reescrever o numerador da Equação 14 utilizando o resultado da Equação 17, obtendo uma expressão final para o Critério de Fisher:

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_A \mathbf{w}}. \quad (18)$$

Precisamos, agora, maximizar a Equação 18 com relação à \mathbf{w} , ou seja, calculamos $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$, resultando em uma equação fechada para o cálculo de \mathbf{w} , dada por:

$$\mathbf{w}^* = \mathbf{S}_A^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad (19)$$

que é, basicamente, a diferença entre as médias modulada pela matriz de espalhamento.

Suponha, agora, que tenhamos um problema de classificação por múltiplas classes, isto é, $\mathcal{Y} = \{\omega_1, \omega_2, \dots, \omega_c\}$. este caso, podemos reduzir a dimensionalidade do espaço original para, no máximo, $c-1$ classes, utilizando a Análise Discriminante Múltipla, do inglês *Multiple Discriminant Analysis* (MDA). Temos que a operação de transformação é dada por:

$$\hat{\mathbf{x}}_i = \mathbf{W}^T \mathbf{x}_i, \quad (20)$$

em que, novamente, $\mathbf{x}_i \in \mathbb{R}^n$ e $\mathbf{W} \in \mathbb{R}^{n \times d}$ é uma matriz de projeção para um espaço com $d < c$ dimensões. Temos que cada coluna de \mathbf{W} é uma direção ortogonal \mathbf{w}_j .

Relembrando que temos as seguintes definições:

- m_i : número de elementos do conjunto de treinamento da classe ω_i .
- μ_i : média dos elementos da classe classe ω_i .
- μ : média global, isto é, de todo o conjunto de treinamento.

A função objetivo (Critério de Fisher) a ser maximizada é dada por:

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_A \mathbf{W}|}, \quad (21)$$

em que $|\cdot|$ calcula o determinante de uma matriz.

Já as matrizes S_A e S_B são calculadas como segue:

$$S_A = \sum_{i=1}^c S_i = \sum_{i=1}^c \sum_{\mathbf{x}_j \in \omega_i} (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T, \quad (22)$$

em que S_i representa a matriz de espalhamento da classe ω_i , e

$$S_B = \sum_{i=1}^c m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T. \quad (23)$$

Pode-se mostrar que o posto (rank) da matriz S_B é $c - 1$, ou seja, o número de linhas ou colunas linearmente independentes é limitado superiormente à $c - 1$ (número de direções discriminantes do método).

A condição para maximização de nosso Critério de Fisher é dada por:

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = 0.$$

Resolvendo analiticamente a equação acima, chegamos na seguinte formulação:

$$(\mathbf{S}_B \mathbf{W}) = \lambda \mathbf{S}_A \mathbf{W}. \quad (24)$$

Caso \mathbf{S}_A admita inversa, temos que:

$$(\mathbf{S}_A^{-1} \mathbf{S}_B) \mathbf{W} = \lambda \mathbf{W}. \quad (25)$$

Para resolver o problema, basta selecionarmos, no máximo, $c - 1$ autovetores de $(\mathbf{S}_A^{-1} \mathbf{S}_B)$ associados aos $c - 1$ maiores autovalores.

Algumas limitações do LDA/MDA:

- Não é interessante para problemas com poucas classes e muitas características.
- Quanto a média das classes são muito próximas, o numerador da Equação 7 tende a 0.
- Situações em que a função de custo (objetivo) assume valores muito altos (classes com formas não lineares).

Como a LDA/MDA se assemelha ao PCA?

Dados os conceitos e modelagem, **por que** o PCA é mais utilizado que a LDA/MDA?