

Aula12

DATA SCIENCE IPT

TURMA 02

Aprendizado Supervisionado :

Há as amostras (X) e os resultados corretos para elas (Y). Queremos descobrir uma função que mapeie X a Y e possa prever, para novos X' , novos Y' . O nome decorre do fato que, após o aprendizado, podemos verificar se o algoritmos aprendeu (supervisionar o aprendizado) usando novas amostras (X) para as quais conhecemos as repostas (Y).

Exemplos : Regressão e Classificação

Aprendizado Não Supervisionado :

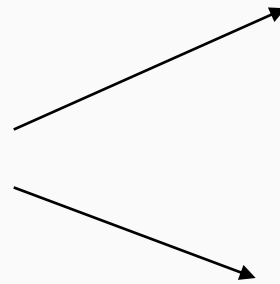
Há as amostras (X) .Queremos descobrir alguma estrutura ou distribuição de dados entre elas. Não há resposta certa ou errada verificável, como no caso dos algoritmos de aprendizado supervisionado.

Clustering : o objetivo é descobrir possíveis agrupamentos entre os dados. Exemplo : K-Means

Association : o objetivo é descobrir tendências do tipo : comprou A então comprou B (para usar em Cross –Sell, por exemplo).
Exemplo : Apriori

Conglomerados (Clusters)

Métodos Hierárquicos



Aglomerativos

Divisivo

Métodos não hierárquicos : Ex: **K-Means**

Algoritmos Hierárquicos

Criam uma hierarquia de conjuntos de classes por fusão de classes menores em classes maiores (ascendente) ou por divisão de classes maiores em classes menores (descendente).

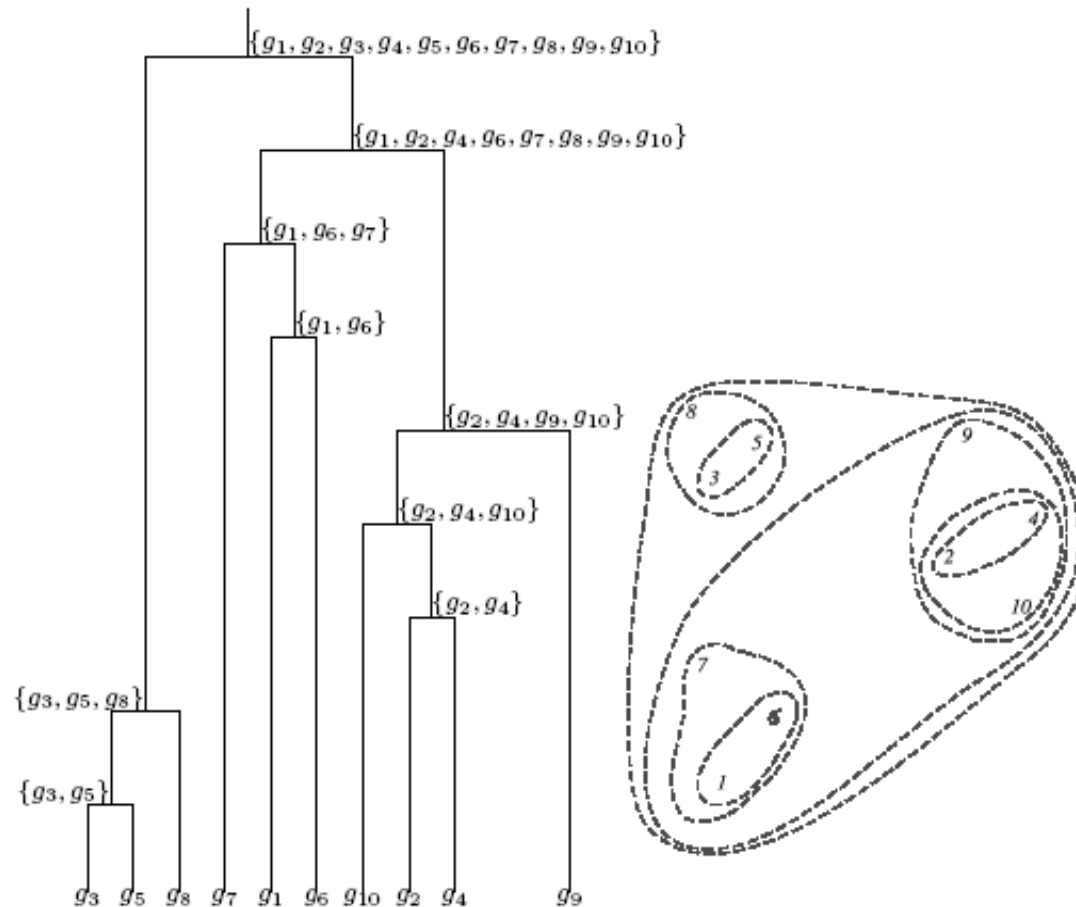
O resultado de um algoritmo hierárquico é uma árvore ou dendograma.

Cortando a árvore num determinado nível é obtida uma partição dos indivíduos em k classes.

Hierárquicos aglomerativos: Partem de n indivíduos agrupados em n classes, cada classe com 1 indivíduo. Agrupam as classes sucessivamente até se obter uma única classe.

Hierárquicos divisivos: Partem de uma única classe que inclui os n indivíduos. As classes são sucessivamente divididas em classes menores até se obterem n classes, cada uma com um indivíduo.

Exemplo : Cluster Hierárquico Aglomerativo



Fonte : An Introduction to Bioinformatics Algorithms

Como é o K-Means em “pseudocódigo”

Defina o número de clusters (k)

Defina os centróides dos k clusters

Faça

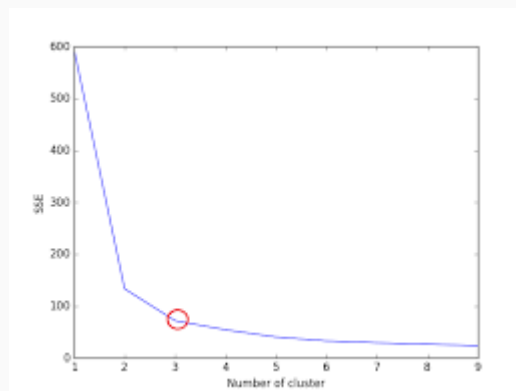
- Forme os k clusters associando cada objeto a seu centróide mais próximo

- Recompute o centróide de cada cluster

Enquanto mudarem os objetos dos clusters

Partindo de Cluster-carro-custo.ipynb, obtenha o custo de kmeans para 2 a 20 clusters (desenvolva a função loss).

Tente identificar o “cotovelo” (número de clusters onde o custo tende a “estabilizar”). É um dos critérios para se determinar o número ideal de clusters...



- Há várias métricas para avaliar a “qualidade” de uma partição do dataset em clusters: Dunn, Silhouette etc...

O índice Dunn é dado pelo quociente entre a menor *distância entre pontos de clusters diferentes e a maior distância entre pontos do mesmo cluster.

$$\text{Dunn} = \frac{\min \text{dist pontos entre diferentes clusters}}{\max \text{dist pontos do mesmo cluster}}$$

Quanto maior o índice, melhor a partição (mais compacto e separado é o cluster)

Para um ponto “i” de um cluster (com k pontos), a média das distâncias dele a cada um dos (k-1) outros pontos do cluster é dada por **a(i)**.

Para um ponto “i” de um cluster, a (menor) média das distâncias dele a cada um dos pontos de um outro cluster é dada por **b(i)**.

Silhouettte para um ponto “i” é $S(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$

Podemos calcular Silhouette para o cluster todo pela média dos pontos e para a partição toda, pela média dos clusters...

Quanto mais próximo de 1 o índice, melhor a partição.

Atividade :

1) Partindo de silhouette_res.ipynb com $k=2$.

Obter (na mão) o índice Silhouette para o pontos 0 e compará-lo com o gerado pelo Scikit

2) Rodar com $k=3$. Comparar índices com $k=3$..qual é a melhor partição?
Discussão.

Silhouette para ponto 0:

$$a=1.414$$

$$b = 5.65+6.4+4+5.09/4=5.27$$

$$s=(5.27-1.414)/5.27=0.73$$

Na classificação, uma métrica usual é a **acurácia**, dada pelo quociente entre o número de amostras corretamente classificadas e o número total de amostras classificadas.

Essa métrica pode ser inadequada quando há grande desbalanceamento entre o número de positivos e negativos. Exemplo : doenças raras.

É importante considerar as métricas precision, recall e F1-Score...

A acurácia pode ser uma métrica inadequada se há grande desbalanceamento entre o número de elementos em cada classe. Suponha uma doença muito rara, que afeta 0.001% da população. Se o algoritmo sempre classificar uma pessoa como negativa (sem a doença) em um milhão de pessoas, por exemplo, teremos :

Positivos ≈ 10

negativos ≈ 999990

O algoritmo que considerar todos negativos terá Acurácia : $999990/1000000=99.999\%$ por outro lado, não encaminhou ninguém para o tratamento...

Foi fraco em identificar os positivos.

CONFUSION MATRIX

TP = True Positives
TN = True Negatives
FP = False Positives
FN = False Negatives

	p' (Predicted)	n' (Predicted)
p (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{F1-score} = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

Precision mostra a parcela de TP entre os que a classificação considerou P

Recall (ou TP rate) mostra a parcela de TP identificados pela classificação entre os que são realmente P

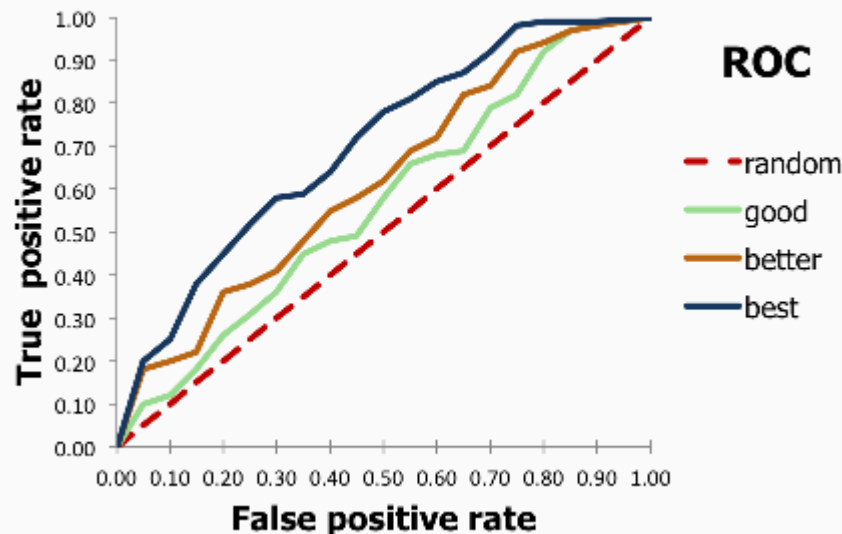
F1-Score combina precision e recall em métrica única, pois normalmente há um trade-off entre eles. (quando se varia o threshold, por exemplo).

Atividade :

Partindo de metrics.ipynb

1) Montar a Confusion Matrix e obter Acurácia, Precision, Recall e F1-Score para a lista de thresholds=[0.10,0.16,0.20,0.50]
Monte também os gráficos das predições

A Curva ROC (Receiver Operating Characteristic) é útil para comparar classificadores. Ela é traçada com FP rate em x e TP rate (Recall) em y para diferentes thresholds. O ponto FP rate=0 e TP rate=1 é o ideal...quando possível. Quanto maior a área sob a curva (mais próxima de 1.0), melhor o classificador. A diagonal (FP rate = TP rate) indica que o classificador é apenas “randômico”.



True Positive Rate (TPR) = Recall = Sensitivity (Prob. de detecção) = $TP/(TP+FN)$

False Positive Rate(FPR) = False alarm ratio = $FP/(FP+TN)$

True Negative Rate(TNR) = Specificity = $TN/(FP+TN) = 1-FPR$

Atividade : partindo de ROC_res.ipynb

Observar código da geração das curvas ROC de Logistic Regression e KNN

2) Comparar as áreas sobre as curvas ROC de logistic regression e KNN

Seja M uma matriz, v um vetor e λ um escalar.

Dizemos que v é um autovetor de M e λ é seu autovalor associado quando:

$$Mv = \lambda v$$

Assim, a transformação linear causada por M em um vetor é equivalente ao produto de um escalar por esse vetor...

$$Mv = \lambda v$$
$$(M - \lambda I)v = 0$$

Essa equação (sistema de equações lineares homogêneas) tem uma solução trivial (vetor $v=0$). Para que existam outras soluções além da trivial, o determinante de $M - \lambda I$ deve ser nulo. Dessa imposição, chegamos em uma equação que nos permite encontrar os autovalores e, depois, os autovetores (obviamente).

Exemplo :

$$M = \begin{pmatrix} 1 & -2 \\ 3 & 9 \end{pmatrix}$$

$$M - \lambda I = \begin{pmatrix} 1 - \lambda & -2 \\ 3 & 9 - \lambda \end{pmatrix}$$

$$\text{Det}(M - \lambda I) = 0 \quad (1 - \lambda)(9 - \lambda) + 6 = 0$$

$$9 - \lambda - 9\lambda + \lambda^2 + 6 = 0$$

$$\lambda^2 - 10\lambda + 15 = 0$$

Raízes (aprox): 8.16 e 1.84....são os autovalores

Para $\lambda=8.16$, o autovetor correspondente é :

$$(M - \lambda I)v = 0$$

$$\begin{pmatrix} -7.16 & -2 \\ 3 & 0.84 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$-7.16x - 2y = 0$$

Exemplo: para $x=1$ $v = \{ 1, -3.58 \}$

$$\text{Teste: } \begin{pmatrix} 1 & -2 \\ 3 & 0.84 \end{pmatrix} \begin{pmatrix} 1 \\ -3.58 \end{pmatrix} = 8.16 \begin{pmatrix} 1 \\ -3.58 \end{pmatrix} \quad ? \quad \text{Vamos ao Python!}$$

Analogamente, para $\lambda=1.84$, o autovetor correspondente é :

$$(M - \lambda I)v = 0$$

$$\begin{pmatrix} -0.84 & -2 \\ 3 & 7.16 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\begin{aligned} -0.84x - 2y &= 0 \\ y &= -0.42x \end{aligned}$$

Exemplo: para $x=1$ $v = \{1, -0.42\}$

$$\text{Teste: } \begin{pmatrix} 1 & -2 \\ 3 & 9 \end{pmatrix} \begin{pmatrix} 1 \\ -0.42 \end{pmatrix} = 1.84 \begin{pmatrix} 1 \\ -0.42 \end{pmatrix} \quad ? \quad \text{Sim! A conferir no Python}$$

É interessante deixar os autovetores com módulo 1. Para isso, basta obter a norma e dividir seus componentes por essa norma...

$\lambda=8.16$ o vetor $v=\{1, -3.58\}$ tem norma: 3.72
e fica $v=\{0.27, -0.96\}$ com norma 1

$\lambda=1.84$ o vetor $v=\{1, -0.42\}$ tem norma: 1.08
e fica $v=\{0.92, -0.39\}$ com norma 1

Vamos conferir com numpy! `autovalores_res.ipynb`

Propriedade Importante:

O produto dos autovalores é o determinante da matriz:

$$M = \begin{pmatrix} 1 & -2 \\ 3 & 9 \end{pmatrix} \quad \det(M) = 9 + 6 = 15 = 8.16 * 1.84$$

A matriz M é ‘Hermitiana’ se a transposta da conjugada é a própria M (?)

Exemplos:

$$\begin{bmatrix} 3 & 2+i & -i \\ 2-i & 5 & 5-i \\ i & 5+i & -1 \end{bmatrix} \quad \begin{bmatrix} 0 & 2 & -i & 10+i \\ 2 & 1 & 0 & -5i \\ i & 0 & -2 & 0 \\ 10-i & 5i & 0 & 5 \end{bmatrix} \quad \begin{bmatrix} 0 & i & i \\ -i & 0 & i \\ -i & -i & 0 \end{bmatrix}$$

É fácil identificá-las: Na diagonal principal só há reais e os elementos simétricos são conjugados ($a+bi$, $a-bi$)

Toda matriz real simétrica é Hermitiana.

Toda Matriz hermitiana tem autovalores reais

Os autovetores associados a autovalores diferentes são ortogonais

Covariância:

$$\sigma(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Matriz de covariância
(para 3 features)

$$\begin{bmatrix} \text{Var}_1 & \text{Cov}_{1,2} & \text{Cov}_{1,3} \\ \text{Cov}_{1,2} & \text{Var}_2 & \text{Cov}_{2,3} \\ \text{Cov}_{1,3} & \text{Cov}_{2,3} & \text{Var}_3 \end{bmatrix}$$

Matrizes de Covariância são reais e simétricas,
portanto Hermitianas

Todas as matrizes de covariância são positivas semi definidas.
Isso implica que todos os autovalores são positivos.

Uma matriz positiva semi definida $A \Rightarrow v_t^T A v_t \geq 0$ para qualquer v diferente de zero

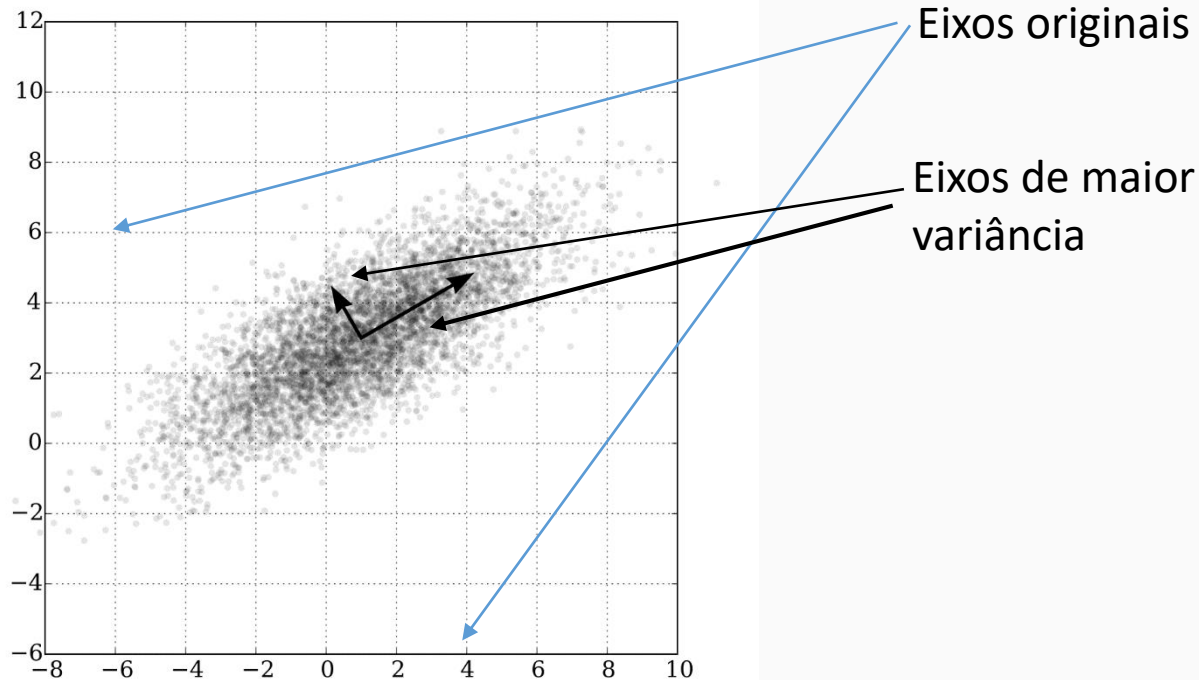
A Análise de Componentes Principais (PCA) é um ferramenta de redução de dimensionalidade que pode ser usada para reduzir um grande conjunto de variáveis para um pequeno conjunto que ainda contém a maioria das informações em o conjunto grande

Basicamente, partimos da matriz de covariância entre as features e extraímos seus autovalores e autovetores correspondentes.

Os autovetores associados a maiores autovalores indicarão maior “variância naquele eixo” => mais importância...mais variância, mais informação

Podemos representar o dataset com menos eixos (dimensões) que o inicial desprezando os eixos com autovalores baixos.

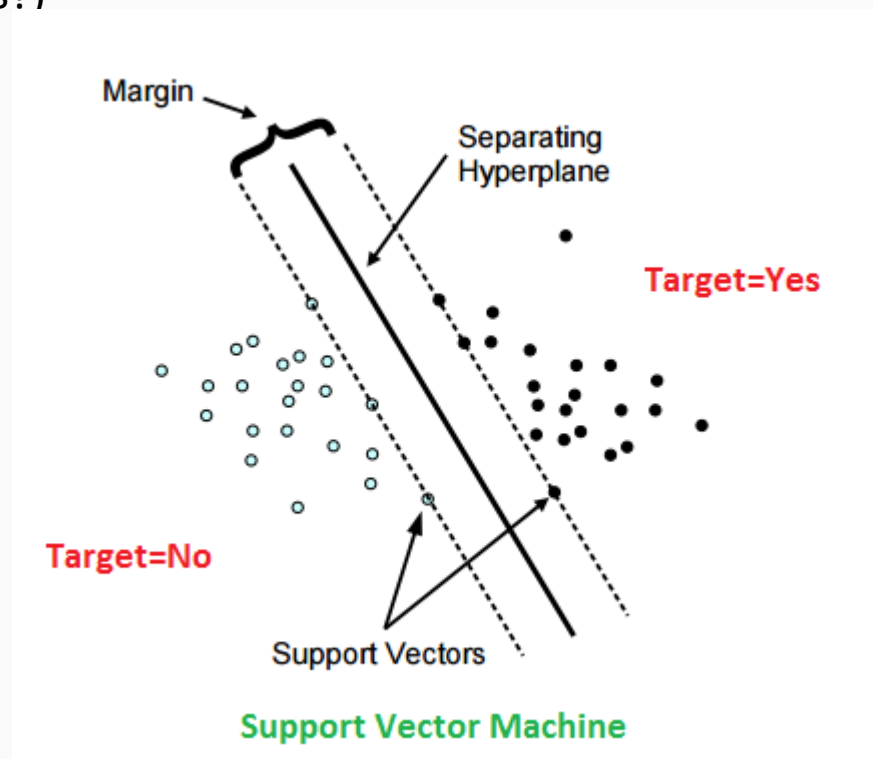
É possível voltar ao dataset inicial, mas com perdas....



Partindo de `pca.ipynb`

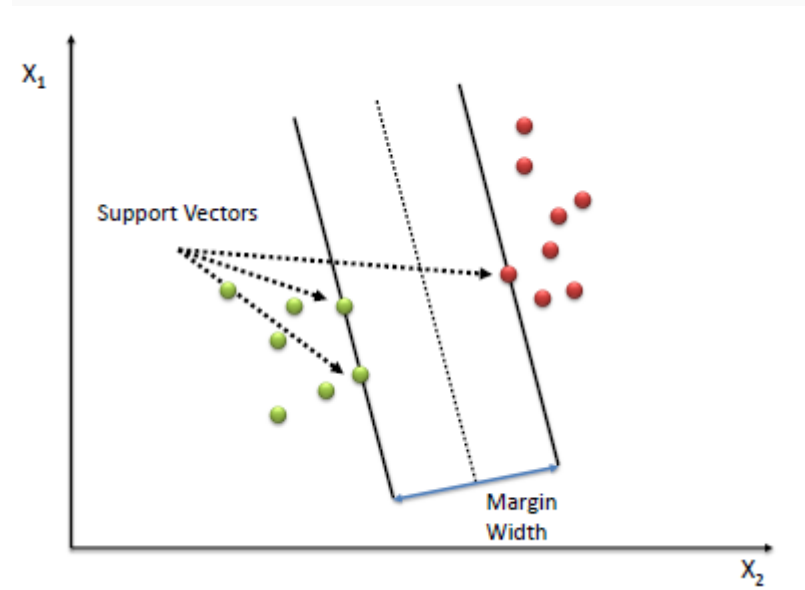
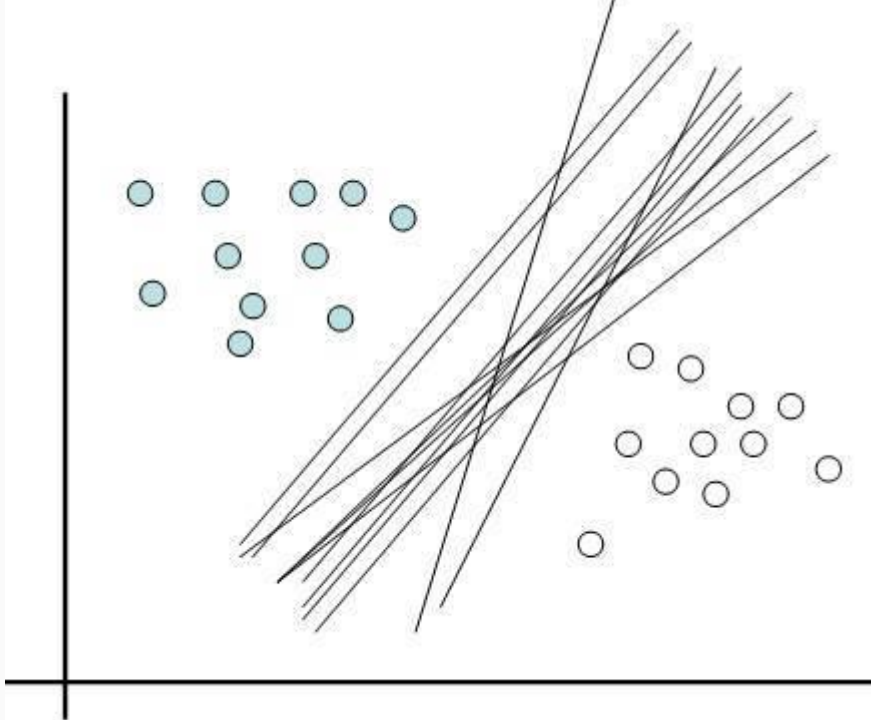
- 1) Analisar o código sem e com scikit
- 2) Verificar a matriz de covariância antes e depois das transformações
- 3) Calcular a acurácia (logistic regression) para (binário) versicolor só com dois componentes e com 4 componentes
- 4) Reconstruir o dataset (de 2 para 4 features) e observar perdas (é uma lossy compression)

Introdução (Linear SVM)...as classes na figura abaixo são claramente linearmente separáveis. Pra criar a “separação” (hiperplano) que deixe as classes o mais “distante” possível, procuramos uma reta que maximiza a distância de cada um dos pontos próximos à fronteira (os de mais difícil classificação), ou seja, maximiza a margem. Esses pontos são os “Support Vectors” (vectors?)

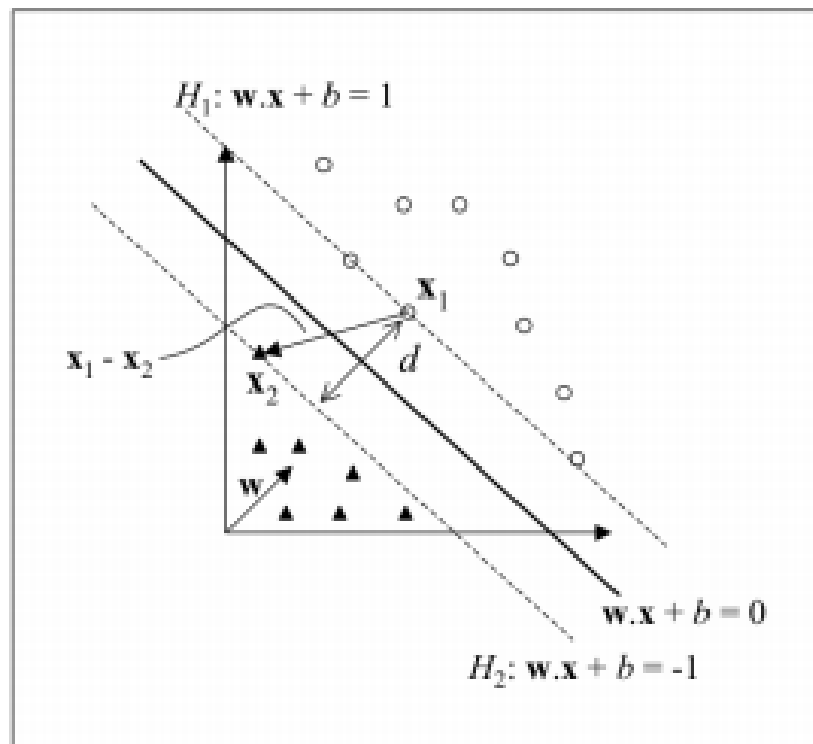


Support Vector Machine é um tipo de algoritmo bastante utilizado para classificação, tipicamente. O principal fundamento desse algoritmo é buscar um hiperplano que crie uma maior margem (distância) entre as classes. Os pontos considerados (matematicamente) para maximizar essa margem são os support vectors.

Como escolher o hiperplano?



Como escolher o *hiperplano?



d (margem) é a decomposição de $x_1 - x_2$ na direção do vetor normal w

A Margem é o Módulo da Projeção de $\mathbf{x1-x2}$ sobre \mathbf{w}

$$\frac{||\mathbf{x1} - \mathbf{x2}|| * \langle \mathbf{x1} - \mathbf{x2}, \mathbf{w} \rangle}{||\mathbf{x1} - \mathbf{x2}|| * ||\mathbf{w}||}$$

$\langle \mathbf{x1}, \mathbf{w} \rangle$ é $1-b$

$-\langle \mathbf{x2}, \mathbf{w} \rangle$ é $1+b$

Assim, a margem é : $\frac{2}{||\mathbf{w}||}$

Para maximizar a margem, devemos minimizar o módulo de \mathbf{w} ...

Minimizar $\| \mathbf{w} \|$ ou $\frac{1}{2} \| \mathbf{w} \|^2$ chega no mesmo resultado. Optamos pela forma quadrática por ser um problema matemático de otimização bem estudado.

O problema final é (margens rígidas):

$$\underset{\mathbf{w}, b}{\text{Minimizar}} \quad \frac{1}{2} \| \mathbf{w} \|^2$$

Com as restrições: $y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0, \quad \forall i = 1, \dots, n$

Matematicamente, criamos uma função Lagrangeana (?) que engloba a função a ser minimizada e as restrições.

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1)$$

Equação 1

Impondo:

$$\frac{\partial L}{\partial b} = 0 \quad \text{e} \quad \frac{\partial L}{\partial \mathbf{w}} = 0$$

Equação 2

Chegamos a:

$$\sum_{i=1}^n \alpha_i y_i = 0$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

Substituindo 3 em 1,
Chegamos ao problema
dual:

$$\text{Maximizar}_{\alpha} \quad \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\text{Com as restrições: } \begin{cases} \alpha_i \geq 0, \quad \forall i = 1, \dots, n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

Sendo α^* vindo do problema dual e os correspondentes w^* e b^* ...há as condições de Kuhn-Tucker para problemas de otimização:

$$\alpha_i^* (y_i (\mathbf{w}^* \cdot \mathbf{x}_i + b^*) - 1) = 0, \forall i = 1, \dots, n$$

Para que as condições sejam satisfeitas com $\alpha^* > 0$, \mathbf{x}_i deverá estar nas bordas....será um Support Vector!

A função decisora será o sinal da fórmula abaixo...

$$g(\mathbf{x}) = \text{sgn}(f(\mathbf{x})) = \text{sgn}\left(\sum_{\mathbf{x}_i \in \text{SV}} y_i \alpha_i^* \mathbf{x}_i \cdot \mathbf{x} + b^*\right)$$

Partindo de lousa-svm_res.ipynb

Atividade1 : analisar código

Atividade 2 : mostrar que a reta é $x_1=1.5$

Atividade 3 : apresentar vetores de suporte

Atividade 4: criar função decisora com base em support vectors e coeficientes da solução dual...

SVM continua na próxima aula

SVM Analisando o código da classe MKMeans (Prof. Leston)



Cursos com Alta Performance de
Aprendizado

© 2019 – Linked Education