

Continuous-Continuous:

Pearson's Correlation if normally distributed:

$$\rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y} \quad r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Spearman correlation (works with ranks) if not:

$$r_s = \rho_{R(X),R(Y)} = \frac{\text{cov}(R(X),R(Y))}{\sigma_{R(X)} \sigma_{R(Y)}} \quad \rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

Categorical-Categorical:

Goodman Kruskal Lambda:

$$\lambda_{B|A} = \frac{P_E - P_{E|A}}{P_E} \quad P_E \text{ is the error in predicting B without knowing A}$$

$P_{E|A}$ is the error in predicting B with knowledge of A

- without knowing A:

$$P_E = 1 - \frac{\max_c N_{c+}}{N_{++}} = 1 - \frac{700}{300+700}$$

- knowing A:

$$P_{E|A} = 1 - \frac{\sum_r \max_c N_{rc}}{N_{++}} = 1 - \frac{118+628}{300+700}$$

Chi Square:

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} \quad O_i \text{ is Observed frequency for the } i^{\text{th}} \text{ category.}$$

E_i is Expected frequency for the i^{th} category,

$$E_{ij} = \frac{(\text{Row Total}) * (\text{Column Total})}{\text{Grand Total}} \quad \text{df} = (n_{\text{Rows}} - 1) * (n_{\text{Columns}} - 1)$$

Continuous-Categorical:

2 classes: **t-Test** (normally distributed), **Mann-Whitney U Test**

(one is not normally distributed)

More than 2: **ANOVA**, **Kruskal-Wallis**

Color: [Red, Blue, Green, ...]
o **Nominal Scale**
Size: [Small, Medium, Large]
o **Ordinal Scale**
Price: [3.65, 2.98, 6.12, ...]
o **Continuous**
Petals: [5, 6, 4, 7, ...]
o **Discrete**

Fisher Score:

$$F(x) = \frac{\sum_{k=1}^K N_k (\mu_k - \mu)^2}{\sum_{k=1}^K N_k \sigma_k^2}$$

- K is the number of classes.
- N_k is the number of samples in class k .
- μ_k is the mean of the feature x in class k .
- μ is the overall mean of the feature x across all classes.
- σ_k^2 is the variance of the feature x in class k .

Entropy based binning:

Find the binary split boundary that minimizes the entropy function over all possible boundaries:

$$H(S_i) = -\sum_{j=1}^N p_{ij} \log_2 p_{ij} \quad H(S,T) = \sum_{i=1}^K \frac{|S_i|}{|S|} H(S_i)$$

Greater the gain: better the split $\text{Gain} = H(S) - H(S,T)$

Character Based Similarity:

Edit Distance: Insert, delete and replace – cost = 1

Affine Distance: Insert, delete, replace, open gap, extend gap

S1 = Chris R Lang
S2 = Christopher Richard Lang

Open gap – 1
Extend "topher" – 2.5 (5 * 0.5)
Open gap – 1
Extend "richard" – 2.5 (5 * 0.5)
Distance = 1 + 2.5 + 1 + 2.5 = 7

Smit-Waterman Distance: Beginning and end mismatches- < weights; ignored.

Token Based Similarity:

$$\text{Overlap}(S_1, S_2) = \frac{|\text{tok}(S_1) \cap \text{tok}(S_2)|}{\min(|\text{tok}(S_1)|, |\text{tok}(S_2)|)} \quad \text{Jaccard}(S_1, S_2) = \frac{|\text{tok}(S_1) \cap \text{tok}(S_2)|}{|\text{tok}(S_1) \cup \text{tok}(S_2)|}$$
$$\text{Dice}(S_1, S_2) = \frac{2|\text{tok}(S_1) \cap \text{tok}(S_2)|}{|\text{tok}(S_1)| + |\text{tok}(S_2)|}$$

Phonetic Similarity:

(1) - Keep the first letter of the surname and ignore all occurrences of W and H (Peter – P) (2) - Assign following codes: (B, F, P, V = 1), (C, G, J, K, Q, S, X, Z = 2), (D, T = 3), (L = 4), (M, N = 5), (R = 6) (Pe3e6) (3) - A, E, I, O, U and Y are not coded and serve as separators (0) (P0306) (4) - Consolidate sequences of identical codes by keeping only the first occurrence (P0306) (5) - Drop separators (P36) (6) - Keep the letter prefix and the three first codes, padding with zeros if there are fewer than 3 codes (P360)

Undersampling:

Random Undersampling: random elimination of majority samples

Tomek Links: For every sample E_i in the dataset, find its nearest neighbor E_j . - If E_i and E_j belong to **different classes**, proceed to the next step. - Verify that E_i and E_j are the **closest to each other** compared to other points in the dataset. - If the conditions above are met, remove one or both points in the pair.

Z-Score:

$$Z = \frac{X - \mu}{\sigma} \quad \sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}} \quad |Z| > k = \text{outlier}$$

SMOTE – Oversampling Technique:

New Point = Starting Point + Fraction × (Target Point – Starting Point)

$$x = x_i + \text{rand}(0 - 1)(\hat{x}_i - x_i)$$

Confusion Matrix:

	Predicted		
	P	N	
Real	P	TP	FN
	N	FP	TN

$$P = TP + FN$$
$$N = TN + FP$$

Correct classification:

- TP = true positives (P classified as P)
 - Malignant tumors classified as malignant
- TN = true negative (N classified as N)
 - Benign tumors classified as benign

Misclassification:

- FN = false negatives (P classified as N)
 - Malignant tumors classified as benign
- FP = false positives (N classified as P)
 - Benign tumors classified as malignant

K-Nearest Neighbors:

- Calcular **distância** entre um ponto e todos os outros. – Selecionar os **k pontos** mais **próximos**. – Para **classificar**, entre os k, ver a classe que existe em **maioria**.

PCA:

- Centralizar usando **z-score**. – Achar a **matriz de co-variância**, para conseguir achar os **eigenvalues** e os **eigenvectors**. – Ordenar eigenvalues e consoante essa ordenação, ordenar os eigenvectors. – **Projetar** na matriz original (**matriz O x eigenvectors ordenados**)

Explained Variance:

$$\beta_i = \frac{\lambda_i}{\sum \lambda_i}$$

Nominal – Tags
Ordinal – Rank

Error Rate (tudo menos diagonal):

$$ER = \frac{FP + FN}{TP + TN + FP + FN} = \frac{FP + FN}{S} \quad C: \text{number of classes}$$

N_i : number of samples of class i

Accuracy (diagonal):

$$Acc = 1 - ER = \frac{TP + TN}{S}$$

Balanced Accuracy:

$$B_{Acc} = \frac{1}{C} \sum_{i=1}^C \frac{TP_i}{TP_i + FN_i} = \frac{1}{C} \sum_{i=1}^C \frac{TP_i}{N_i}$$

Class True Positive Rate:

$$TPR_i = \frac{TP_i}{P_i} = \frac{TP_i}{TP_i + FN_i}$$

P_i : positives for class i
 C : Nº of classes
 S : Nº of test instances

Global TP Rate weighted = Acc

Global TP Rate macro = BAcc

Global TP Rate micro = Acc

Class False Positive Rate:

$$FPR_i = \frac{FP_i}{N_i} = \frac{FP_i}{FP_i + TN_i}$$

Global FP Rate w:

$$FPR_w = \frac{1}{S} \sum_{i=1}^C P_i \cdot FPR_i$$

Global FP Rate ma:

$$FPR_m = \frac{1}{C} \sum_{i=1}^C FPR_i$$

Global FP Rate mi:

$$FPR_m = \frac{\sum_{i=1}^C FPR_i}{\sum_{i=1}^C N_i} = \frac{1}{2S} \sum_{i=1}^C FPR_i$$

Class Precision:

$$Pr_i = \frac{TP_i}{TP_i + FP_i}$$

Global Precision w:

$$Pr_w = \frac{1}{S} \sum_{i=1}^C P_i \cdot Pr_i$$

Global Precision ma:

$$Pr_m = \frac{1}{C} \sum_{i=1}^C Pr_i$$

Global Precision mi = Acc

Class Recall:

$$Re_i = \frac{TP_i}{TP_i + FN_i}$$

All Global Recall = Acc

Class F1-Score:

$$F1_i = 2 \frac{Pr_i \cdot Re_i}{Pr_i + Re_i}$$

Global F1 w:

$$F1_w = \frac{1}{S} \sum_{i=1}^C P_i \cdot F1_i$$

Global F1 ma:

$$F1_m = \frac{1}{C} \sum_{i=1}^C F1_i$$

Global F1 mi = Acc

Discard Data:

Complete case – eliminar linha

Pairwise – eliminar o valor

Dropping variables – eliminar coluna

Undersampling:

Condensed Nearest Neighbor Rule: Eliminate the samples of the majority class that are **far away from the decision border**.

One Sided Selection: Tomek Links + CNN

Neighborhood Cleaning Rule: Clean the data while reducing it. - Get the 3-NN of a data sample: If data sample of majority class and 3-NN contradict classification eliminate data sample else if data sample of minority class and 3-NN contradict classification eliminate data samples of 3-NN that belong to majority class

Oversampling:

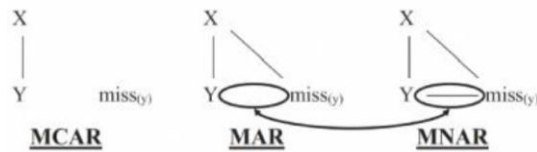
Naïve Random

Oversampling:

Randomly duplicate (with replacement) minority data points

SMOTE – já falado antes

MCAR, MAR, MNAR:



Imbalance Ratio:

$$IR = \frac{\text{\#Positive Class}}{\text{\#Negative Class}}$$

Linear Regression:

$$\underset{\beta}{\operatorname{argmin}} \sum_{i=1}^p (y_i - \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_p x_{i,p})^2 = \underset{\beta}{\operatorname{argmin}} \|Y - X\beta\|^2$$

$$\frac{\partial}{\partial \beta} \|Y - X\beta\|^2 = 0 \Leftrightarrow \beta = \underbrace{(X^T X)^{-1} X^T Y}_{X^+} \quad \text{Pseudo-inverse}$$

LDA e FDA:

Detect direction that maximizes the separation between classes
Maximum separation between means of projected classes and
Minimum variance within each projected class

$$\max \left\{ \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2} \right\}$$

Find a vector w that maximizes:

$$\underset{w}{\operatorname{argmax}} \left(\frac{(\mu_1^T w - \mu_2^T w)^2}{w^T \Sigma_1 w + w^T \Sigma_2 w} \right) \Leftrightarrow \underset{w}{\operatorname{argmax}} \left(\frac{\left(\frac{m^T}{(\mu_1^T - \mu_2^T) w} \right)^2}{w^T (\Sigma_1 + \Sigma_2) w} \right) \Leftrightarrow \underset{w}{\operatorname{argmax}} \left(\frac{(m^T w)^2}{w^T S w} \right)$$

- Standard deviation is the same for all classes

$$\underset{w}{\operatorname{argmax}} \left(\frac{((\mu_1^T - \mu_2^T) w)^2}{w^T S w} \right) \Leftrightarrow \underset{w}{\operatorname{argmax}} \left(\frac{w^T \overbrace{(\mu_1 - \mu_2)}^{S_b} (\mu_1 - \mu_2)^T w}{w^T S w} \right) \Leftrightarrow \underset{w}{\operatorname{argmax}} \left(\frac{w^T S_b w}{w^T S w} \right)$$

ReliefF: (ver código)

For each sample xi update feature score W[f] for feature f based on feature distance from point xi to nearest hit Hi (same class) and nearest miss Mi (other class)

$$W[f] = W[f] - \text{diff}(f, \mathbf{x}_i, \mathbf{H}_i) + \text{diff}(f, \mathbf{x}_i, \mathbf{M}_i)$$

Find the closest element of the same class to be the **Hit**, and find the closest element from an-other class to be the **Miss**. When there are more than two classes, find all the misses, sum them all up, and

Sum Squared Error:

$$SSE = \sum_{i=1}^N (y_i - \hat{y}(x_i, w))^2$$

 y_i : real values
 \hat{y} : predicted values
 N : size of the test set

Mean Squared Error:

$$MSE = \frac{1}{N} \cdot SSE$$

Root Mean Square Error:

$$RMSE = \sqrt{MSE}$$

Coefficient of Determination:

$$R^2 = 1 - \frac{SSE}{SST}, \text{ where } SST = \sum_{i=1}^N (y_i - \bar{y})^2$$

Extreme Value Analysis:

$$[\mu - k\sigma, \mu + k\sigma]$$

K-Means:

- **k** initial cluster **centers**. - each point **allocated** to the **closest center**. - calculate the **mean** of x and y for each cluster's values (new values are the new centers). - stop when nothing changes

DBSCAN:

A point *p* is a **core point** if at least **minPts** points are **within distance ε** of it (including *p*). - A point *q* is **directly reachable** from *p* if point *q* is **within distance ε** from core point *p*. Points are only said to be directly reachable from core points. - A point *q* is **reachable** from *p* if there is a **path** *p*1, ..., *p*_n with *p*1 = *p* and *p*_n = *q*, where each *p*_i+1 is directly reachable from *p*_i. Note that this implies that the initial point and all points on the path must be core points, with the possible exception of *q*. - All points **not reachable** from any other point are **outliers** or **noise points**.

Interquartile range:

$$[Q_1 - k \times IQR, Q_3 + k \times IQR]$$

$$IQR = Q3 - Q1$$

Normalization:

If it's not normally distributed

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Standardization:

It's normally distributed

$$Z = \frac{X - \mu}{\sigma}$$

Mutual Information:

$$I(X, Y) = I(Y, X)$$

$$= H(Y) - H(Y|X)$$

$$= H(X) - H(X|Y)$$

$$I(X; Y) = D_{KL}(P(X, Y), P(X)P(Y))$$

$$= \sum_{x \in \mathcal{X}_x} \sum_{y \in \mathcal{X}_y} P(x, y) \log_2 \frac{P(x, y)}{P(x)P(y)}$$

$$= \sum_{x \in \mathcal{X}_x} \sum_{y \in \mathcal{X}_y} P(x, y) \log_2 \frac{P(y|x)}{P(y)}$$

Entre F e C quero apresentar o ranking usando IM: - calcular probabilidades. - calcular IM para (F=S,R=B), (F=S,R=A), (F=N,R=B), (F=N,R=A). - somar esses IMs, o que vai dar IM(F,R). - fazer o mesmo para C. - se I(F,R) maior que I(C,R), então ranking vai ser 1º:F, 2º:C

ANOVA:

Determine degrees of freedom

$$df_{between} = k - 1, k = \text{nGroups}$$

$$df_{within} = n - k, n = \text{nPoints}$$

$$SS_{within} = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \mu_i)^2$$

$$Var_{within} = \frac{SS_{within}}{df_{within}}$$

$$SS_{Total} = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{\mu})^2$$

$$SS_{Total} = SS_{between} + SS_{within}$$

$$Var_{between} = \frac{SS_{between}}{df_{between}}$$

$$F = \frac{Var_{between}}{Var_{within}}$$

Overfitting: occurs when the model learns specific patterns in the training data that do not generalize to new unseen data – learns noise (ajuste exagerado dos dados de treino)

Underfitting: it happens when the model can't capture the underlying patterns in the data, resulting in poor performance on both training and test datasets

Training Set: no partition (training=test)

Train-Test: 70/30

Stratified Hold Out: TT but train and test have balanced classes

Repeated Stratified TT: repeat TT several times with different random samples

Repeated, Stratified, TTV: 40/30/30

Stratified K-Fold Cross Validation: randomly separates the dataset into K groups. Only 1 group for testing. Repeat K times

Repeated Stratified K-Fold: repeats k-fold several times after reshuffling and re-stratifying the data

Nested Repeated Stratified K-Fold: repeated stratified k-fold using TTV

Leave One Out Cross-Validation: each fold contains only 1 sample

ZeroRule (ZeroR):

Algorithm that classifies all samples in the test set into one class: the majority class in the train set.
Random Classifier: Randomly assigns a class to each instance
OneRule (OneR): Learn rules from a single feature (not a single rule).

OneRule – Example:

Feature: **Age**

1. Discretise: one approach is to find the more accurate binary split. Ex: 39 or 43 years.

2. Rules

- 39: ≤ 39 : 5 Yes, 1 No: If Age $\leq 39 \rightarrow$ yes
- > 39 : 5 No, 4 Yes: If Age $> 39 \rightarrow$ no
- 43: ≤ 43 : 9 Yes, 3 No: If Age $\leq 43 \rightarrow$ yes
- > 43 : 3 No, 0 Yes: If Age $> 43 \rightarrow$ no

3. Error rate

- ER39 = $\frac{1+4}{15} = 33\%$
- ER43 = $\frac{3+0}{15} = 20\%$

Feature: **Income Range**

1. Discrete feature

2. Rules

- 50-60: 2 Yes, 0 No: If IR 50-60 \rightarrow yes
- 40-50: 3 No, 1 Yes: If IR 40-50 \rightarrow no
- 30-40: 4 Yes, 1 No: If IR 30-40 \rightarrow yes
- 20-30: 2 Yes, 2 No: If IR 20-30 \rightarrow no (because no is the less represented class)

3. Error Rate

ER = $\frac{0+1+1+2}{15} = 27\%$

Choose the rule with the **smallest ER**

Support Vector Machines: Find the **hyperplane** that best divides a dataset into **two classes**.

Artificial Neural Network:

Input, hidden and output variables are represented by **nodes**. Nodes are organized in **layers**. The weights are the **links** between nodes. The network can be **fully connected** or **sparse**, i.e. not all connections are present
Feed Forward NN: Data flows forward between layers, with no loops or cycles. (**Gradient Descent:** the **w** update comprise a small step in the direction of negative gradients: at each step, **w** is moved in the direction of the greatest rate of decrease of **E(w)**),
Backpropagation: the derivatives of the errors are propagated back through the network, to adjust **w**)
Convolutional NN: Nearby pixels are more strongly correlated than more distant pixels.

Naive Bayes:

Bayes Theorem:

$$P(Y | X) = \frac{P(X | Y) \cdot P(Y)}{P(X)}$$

For classification problems:

$$P(C_k | x) = \frac{P(x | C_k) \cdot P(C_k)}{P(x)}$$

x: features vector
C_k: output classes

Because features are assumed **independent**:

$$P(x | C_k) \cdot P(C_k) = P(C_k, x) \quad M = P(x)$$

Probability of class C_k:

$$P(C_k | x) = \frac{P(x | C_k) \cdot P(C_k)}{P(x)} = \frac{1}{M} \cdot P(C_k) \prod_{i=1:n} P(x_i | C_k)$$

Assigned class to a new sample will be:

$$\hat{C} = \operatorname{argmax}_k P(C_k) \prod_{i=1:n} P(x_{\text{new}} | C_k)$$

ID3:

- Select the most important attribute to place at the root node, i.e. attribute with the highest **Information Gain**. - Make a branch for each possible value. - Repeat the process (attribute selection) for each branch using only the instances in the branch. **Any attribute can appear at most once along any branch through the tree.**

Entropy:

$$H(S) = - \sum_{k=1}^K p_k \cdot \log_2(p_k)$$

Information Gain: measures how well a given attribute separates the training data according to their target classification

$$IG(S, A) = H(S) - \sum_{v \in A} \frac{S_v}{S} \cdot H(S_v) \quad \begin{matrix} A: \text{attribute with possible values } v \\ S_v: \text{samples with value } v \text{ for attribute } A \end{matrix}$$

Entropy after partitioning S using attribute A

$$IG(S, A) = H(S) - H(S|A)$$

Naïve Bayes example:

▪ New instance:

Outlook	Temperature	Humidity	Windy	Play
sunny	cool	high	true	?

▪ $P(C_k | \text{Sunny, Cool, High, Windy})$: $\hat{C} = \operatorname{argmax}_k P(C_k) \prod_{i=1:n} P(x_{\text{new}} | C_k)$

Outlook	Temperature	Humidity	Windy	Play
sunny 2/9 3/5	hot 2/9 2/5	high 1/3 4/5	false 2/3 1/5	9/14 5/14
overcast 4/9 0/5	mild 4/9 3/5	normal 2/3 1/5	true 1/3 3/5	
rainy 1/3 2/5	cool 1/3 1/5			

$$\left. \begin{aligned} P(C_{\text{yes}}) \prod_{i=1:n} P(x_i | C_{\text{yes}}) &= \frac{9}{14} \cdot \frac{2}{9} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} = \frac{18}{3402} = 0.0053 \\ P(C_{\text{no}}) \prod_{i=1:n} P(x_i | C_{\text{no}}) &= \frac{5}{14} \cdot \frac{3}{5} \cdot \frac{1}{5} \cdot \frac{4}{5} \cdot \frac{3}{5} = \frac{180}{8750} = 0.0206 \end{aligned} \right\} \begin{aligned} P(C_{\text{yes}} | x) &= \frac{0.0053}{0.0053 + 0.0206} = 0.205 \\ P(C_{\text{no}} | x) &= \frac{0.0206}{0.0053 + 0.0206} = 0.795 \end{aligned}$$

C4.5:

Gain Ratio: incorporates a term sensitive to how broadly and uniformly the attribute split the data

$$GR(C, A) = \frac{IG(C, A)}{H(A)}$$

Example:

$$IG(S, \text{Wind}) = H(S) - \sum_{v \in A} \frac{S_v}{S} \cdot H(S_v) = 0.05$$

Wind=weak $\rightarrow 8$
Wind=strong $\rightarrow 6$

$$H(A = \text{Wind}) = -\frac{8}{14} \cdot \log_2\left(\frac{8}{14}\right) - \frac{6}{14} \cdot \log_2\left(\frac{6}{14}\right) = 0.985$$

$$GR(S, A) = \frac{IG(S, \text{Wind})}{H(A = \text{Wind})} = \frac{0.05}{0.985} \approx 0.051$$

Missing values: assign the most common value within the instances of the same class.

ID3 – example:

$$\begin{aligned} IG(S, \text{Wind}) &= H(S) - \sum_{v \in A} \frac{S_v}{S} \cdot H(S_v) \\ \text{Wind=weak} &\rightarrow (6 \text{ Y}, 2 \text{ N}) \\ \text{Wind=strong} &\rightarrow (3 \text{ Y}, 3 \text{ N}) \\ H(S_{v=\text{weak}}) &= -\frac{6}{8} \cdot \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \cdot \log_2\left(\frac{2}{8}\right) = 0.81 \\ H(S_{v=\text{strong}}) &= -\frac{3}{6} \cdot \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \cdot \log_2\left(\frac{3}{6}\right) = 1 \\ IG(S, \text{Wind}) &= 0.94 - \frac{8}{14} \cdot 0.81 - \frac{6}{14} \cdot 1 \\ IG(S, \text{Wind}) &= 0.05 \quad \text{Low IG} \end{aligned}$$
$$\begin{aligned} H &= -p_{\text{No}} \cdot \log_2(p_{\text{No}}) - p_{\text{Yes}} \cdot \log_2(p_{\text{Yes}}) \\ p_{\text{No}} &= 5/14 = 0.36 \\ p_{\text{Yes}} &= 9/14 = 0.64 \\ H &= -0.36 \cdot \log_2(0.36) - 0.64 \cdot \log_2(0.64) \\ H(S) &= 0.94 \end{aligned}$$

Still ID3:

The training procedure looks to minimize the SD of the target on each decision node

Procedure: Compute the SD of the target: **SD(T)**. - Compute the SD for each target and feature:

$SD(T, x_i) = \sum_{v \in x_i} p(v) SD(T_{x_i=v})$ $v \rightarrow$ values of the feature x_i . - Compute the SD reduction (**SDR**) considering each feature as decision node and select the feature with higher SDR: $SDR(T, x_i) = SD(T) - SD(T, x_i)$

Another example:

Compute the SD(T)= 9.321

2. Compute the SD for each target and feature:

$$SD(T, x_i) = \sum_{v \in x_i} p(v) SD(T_{x_i=v})$$

Outlook	P(v)	SD(T _v)	Temp	P(v)	SD(T _v)	Humidity	P(v)	SD(T _v)	Windy	P(v)	SD(T _v)
Overcast	4/14	3.491	Hot	4/14	8.955	High	7/14	9.363	True	6/14	10.593
Rainy	5/14	7.782	Mild	6/14	7.652	Normal	7/14	8.734	False	8/14	7.873
Sunny	5/14	10.87	Cool	4/14	10.512	SD(T, Outlook) = 9.048		SD(T, Outlook) = 9.039			
SD(T, Outlook) = 7.659			SD(T, Temp) = 8.841								

3. Compute the SDR:

$$\begin{aligned} SDR(T, \text{Outlook}) &= 9.321 - 7.659 = 1.662 \\ SDR(T, \text{Temp}) &= 9.321 - 8.841 = 0.48 \\ SDR(T, \text{Humidity}) &= 9.321 - 9.048 = 0.273 \\ SDR(T, \text{Windy}) &= 9.321 - 9.039 = 0.282 \end{aligned} \quad \rightarrow \text{Maximum SD reduction: root node}$$

6.2 - Auxiliary method to calculate distances between sounds for the ReliefF algorithm

```
def smallerDistance(sound, classIndexes):
```

```
    classDistances = []
```

```
    for index in classIndexes:
```

```
        # Calculate the Euclidean distance between the given sound and the sound at the current index
```

```
        distance = np.sqrt(np.sum((sound - featureMatrix[index]) ** 2))
```

```
        # Append the distance and the corresponding index as a tuple to the classDistances list
```

```
        classDistances.append((distance, index))
```

```
# Sort the classDistances list based on the distance (first element of the tuple)
```

```
classDistances.sort(key=lambda x: x[0])
```

```
nearestIndex = classDistances[0][1] # Get the index of the nearest sound (smallest distance)
```

```
nearestSound = featureMatrix[nearestIndex] # Retrieve the nearest sound from the featureMatrix
```

```
return nearestSound
```

6.2 - Method to implement the ReliefF

```
def reliefF():
```

```
    reliefResults = np.zeros(featureMatrix.shape[1]) # Initialize an array to hold the ReliefF results for each feature
```

```
    nSamples = 100
```

```
    for i in range(nSamples):
```

```
        randomIndex = np.random.randint(0, featureMatrix.shape[0]) # Randomly select an index from the featureMatrix
```

```
        # Get the sound and its corresponding class for the randomly selected index
```

```
        randomSound = featureMatrix[randomIndex]
```

```
        randomClass = keyMatrix[randomIndex]
```

```
        # Find indexes of sounds that belong to the same class as the random sound
```

```
        sameClassIndexes = np.where(keyMatrix == randomClass)[0]
```

```
        # Exclude the randomly selected index from the same class indexes
```

```
        sameClassIndexes = sameClassIndexes[sameClassIndexes != randomIndex]
```

```
        nearestHit = smallerDistance(randomSound, sameClassIndexes) # Find the nearest sound (hit) in the same class
```

```
        nearestMisses = []
```

```
        for keyIndex, key in enumerate(keys):
```

```
            # Check if the current key index is not the same as the random class
```

```
            if keyIndex != randomClass:
```

```
                diffClassIndexes = np.where(keyMatrix == keyIndex)[0] # Find indexes of sounds that belong to a different class
```

```
                # Find the nearest sound (miss) from this different class
```

```
                nearestMiss = smallerDistance(randomSound, diffClassIndexes)
```

```
                nearestMisses.append(nearestMiss) # Append the nearest miss to the nearestMisses list
```

```
        for featureIndex in range(featureMatrix.shape[1]):
```

```
            hitDiff = abs(randomSound[featureIndex] - nearestHit[featureIndex]) # Calculate the difference for the hit
```

```
            reliefResults[featureIndex] -= hitDiff # Update the ReliefF results for this feature by subtracting the hit difference
```

```
            # Loop through each miss and accumulate the differences
```

```
            missDiff = 0
```

```
            for miss in nearestMisses:
```

```
                missDiff += abs(randomSound[featureIndex] - miss[featureIndex])
```

```
            missDiff /= (len(keys) - 1) # Average the miss differences over the number of classes (excluding the hit class)
```

```
            # Update the ReliefF results for this feature by adding the average miss difference
```

```
            reliefResults[featureIndex] += missDiff
```

```
    reliefResults /= nSamples # Average the ReliefF results over the number of samples
```

```
    finalResults = []
```

```
    for i, result in enumerate(reliefResults):
```

```
        finalResults.append((result, i + 1))
```

```
    finalResults.sort(key=lambda x: x[0], reverse = True) # Sort the final results in descending order based on the result value
```

```
    return finalResults
```