

Pattern Discrimination: (1) **feature extraction**, mapping between raw data space and the feature space; (2) **classification**, mapping between feature space and interpretation space. We do classification based on similarity (distance between feature vectors). **Stages of RP system** include measurement/sensing, raw data preprocessing (optional), feature extraction, feature preprocessing (optional), feature selection/reduction (optional), classification; postprocessing; performance evaluation (testing.) **Main concerns:** difficult to know if we have a sufficiently large and representative dataset; features should be discriminative and not redundant, there should be an adequate parameter search for any method chosen; evaluate model with different feature/parameter sets, for robustness; consider complexity vs performance trade-off. **Precision:** $TP/TP+FP$, $Acc = TP+TN/(TP+TN+FP+FN)$, $Sensitivity = TP/(TP+FN)$, $Specificity = TN/(TN+FP)$, $F1 = 2 * Acc * Sensitivity / (Acc + Sensitivity)$

Decision functions: maps between feature space and interpretation space, splits IS into decision regions, if we apply a monotonic function to it we get exact same classification; also decision function can be a linear combination of transformations to features. This is helpful to turn complex non-linear problems into linearly separable problems in higher dimension or transformed spaces. Polynomial $d(x)$ requires setting $(k+d)/l$ ($k \neq l$) parameters, with $d=n_features$ and $k=max_pol_degree$.

Multi-class classification: divide into binary problems. Can be **absolute separation** (one-vs-all, we develop C classifiers, classifier i tells us whether sample belongs to i or not, and then we make hierarchical decisions) or **pairwise separation** (one-vs-one, we develop $C(C-1)/2$ ($n-2$) classifiers, each tells us what class in that pair the point most likely belongs to, final decision through majority voting).

Covariance Matrix: indicates how variables varies and co-variety. It's symmetric. It's related to the **Mahalanobis distance**: adjusts to many cluster shapes, ie, data distribution, it's **scale invariant** (e. circle -> ellipsis). The mahalanobis decision surface between class 1 and 2 is linear if $C_1=C_2$ and quadratic otherwise. **PCA:** find a reduced set of new variables with less redundancy, minimizing loss. Project data into directions (eigenvectors of C) that maximize variance. The eigenvalues represent the variance along those directions. **Steps:** standardize data ($\|f\|_2=1$), compute C ; find eigenvalues λ such that $|I - C| = 0$ and then eigenvalue matrix such that $C = V\lambda V^T$. Project data according to important PCs. **Kaiser Criterion:** keep eigenvectors with eigenvalues > 1 ; **Scree test:** discard eigenvalues after plot stabilizes (subjective). PCA is **unsupervised** (doesn't need labels) which can be good but also means we can cut features because they don't contribute to variance but they actually contribute to overall discrimination. Also PCs are linear transformations, so they struggle with non-linear relationships. PCs have no physical/semantic meaning.

Visualization: graphical inspection of features helps see their discriminative power. Histograms let us know where data is located (mean/median), how spread it is (std/quartiles), whether it's symmetric or skewed and whether there are outliers. We should do a distribution model assessment to know whether we can apply parametric tests to it. **Kolmogorov-Smirnov (KS):** quantifies the difference between the empirical cumulative distribution function (CDF) of the data and the theoretical CDF of a normal standard distribution. The null hypothesis states that the data is drawn from a normal distribution. **Shapiro-Wilk:** is the same but for $n < 25$. **Techniques to avoid overfitting:** Cross-validation, regularization to constrain complexity, pruning (in decision trees, ensemble methods (bagging/random forest)) **Kruskal-Wallis:** non-parametric statistical test to assess features discriminative power. Formula: (1) sort feature values; (2) assign ordinal rank; (3) compute the H statistics: as high as H is, the more discriminative the feature. **H0:** samples from different classes are drawn from the same population (no discriminative power.) **Feature Redundancy:** use correlation matrix. **Dimensionality Ratio problem:** $dr = n/d$, where n is the number of samples and d number of dimensions. If we have low n , we have poor representation of problem, however, with too many dimensions we can overfit data. This is the **curse of dimensionality**: there exists a critical feature dimension (CFD) above which performance degrades (we can't just add more dimensions while keeping the number of training samples the same). In geometric pov, if we have d features each divided in m buckets, then we have m^d hypercubes dividing space – grows exponentially w d and can lead to very sparse space, representing the problem poorly. CFD not theoretically ascertainable. **MDC:** consider class means as class prototypes (representatives) and assign to each new sample the class of the nearest prototype. **Euclidean MDC:** the decision function forms hyperplanes perpendicular to the segment linking the means and passing at half-distance.

• The H statistics is:

$$\begin{aligned} H &= \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} r_{ij} - \bar{R}^2 \\ R &= \frac{\sum_{i=1}^n \lambda_i^2}{\sum_{i=1}^M \lambda_i^2} \times 100 \end{aligned}$$

$$R_i \rightarrow \text{Average of ranks for the samples belonging to class } i$$

$$\bar{R} \rightarrow \text{Average of all ranks for all classes}$$

$$r_{ij} \rightarrow \text{Number of samples belonging to class } i$$

$$n \rightarrow \text{Total number of samples}$$

For **Mahalanobis MDC**: we can use as C the pooled covariance matrix (mean of Cs of classes.) This MDC equals the Euclidean MDC when C is proportional to I . **Fisher LDA:** goal is to separate samples of diff classes by placing them in a space maximizing between-class separability and minimizing within-class variability. We can use this for dimensionality reduction and linear classifier (Fisher LDA). LDA assumes implicitly that the C of each class is equal, because the same within-class scatter matrix is used for all the classes considered. We find vector w and project data according to that direction, then use euclidean MDC. LDA reduces data to at most C-1 dimensions, and it has

• $w = S_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$; $S_w = S_1 + S_2 = (N_1 - 1)C_1 + (N_2 - 1)C_2 = 2(N-1)C$; $d(x) = w^T x - \frac{w^T \mathbf{m}_1 + w^T \mathbf{m}_2}{2}$

$S_w^{-1} = \frac{1}{2(N-1)}C^{-1}$ thus $w = \frac{1}{2(N-1)}C^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$

trouble with some layouts, for example when centroids of classes are the same (ex. Concentric) or when there's a big overlap between classes in any direction.

Bayesian Classifier: takes a probabilistic approach, makes decisions based on probabilities and their associated costs. Training involves getting a priori probabilities and class conditional probabilities. **A priori probability** $P(w_i)$: probability of class w_i without knowing any specific measurements or data distribution. **Class conditional probability density function** $p(x|w_i)$: likelihood of measurement x to occur when the class is w_i , we can estimate this using training data. **A posterior probability**, $P(w_i|x)$: probability of class w_i when the measurement value is x . **Minimum Error Classification Starting Point:** we want to choose the most probable class, ie, choose w_1 if $P(w_1|x) > P(w_2|x)$.

Decision boundary: we can say decision boundary is $y(x) = P(x|w_1)P(w_1) - P(x|w_2)P(w_2)$, or $d(x) = g_1(x) - g_2(x)$, with $g_i(x) = \ln(P(x|w_i)P(w_i)) = \ln(P(x|w_i)) + \ln(P(w_i))$. This is helpful when we equal $P(x|w_i)$ to the normal distribution. If we develop the function $d(x) = g_1(x) - g_2(x)$ using the formula to the basically as the same as Mahalanobis

$$\Delta(x) = \frac{p(x|\omega_1)}{p(x|\omega_2)} > \frac{p(\omega_2)}{p(\omega_1)}$$
 then $x \in \omega_1$ else $x \in \omega_2$ or $\Delta(x) = \frac{p(x|\omega_1)}{p(x|\omega_2)} \leq \frac{p(\omega_2)}{p(\omega_1)}$

Bayesian behaves exactly like Likelihood ratio Prevalence Threshold

for both classes, identical to MDC. **If** reduce number of incorrect classifications for the higher prevalence class.

• $g_i(x) = \ln(p(x|\omega_i)P(w_i)) = \ln(p(x|\omega_i)) + \ln(P(w_i)) = -\frac{1}{2}(x - \mu_i)^T C_i^{-1}(x - \mu_i) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(|C_i|) + \ln(P(w_i))$

• $d_j(x) = g_i(x) - g_j(x) = w^T x + w_0$,

Where,

$w = C^{-1}(\mu_i - \mu_j)$; $w_0 = -\frac{1}{2}(\mu_i^T C^{-1} \mu_i - \mu_j^T C^{-1} \mu_j) + \ln\left(\frac{P(\omega_i)}{P(\omega_j)}\right)$

For **Minimum Risk Classification Starting Point**: decide w_1 if $I_1 < I_2$, I_i is the loss function for choosing w_i , aka, describes how much expected loss is associated with that decision. Se desenvolvemos $I_1 < I_2$ vai dar algo semelhante ao LRT. **If equal loss** then the righthand side is simply the prevalence ratio, aka it would equal the **minimum error classification**. **Loss function**: $L(w_i, w_j)$ cost of taking action i when the true state is j . **Risk/Total Risk:** expected loss for a specific decision given an input x (loss over all possible true states weighted by probabilities): $R(w_i) = \sum_j P(j|w_i) L(w_i, j)$

Risk: Minimum possible risk achievable for a given problem (depends on the data itself, not the classifier, obtained when bayes rule is perfectly followed). **Bayes Rule:** Select action a with lowest Total Risk $R(a|w)$. **Strategies 4 obtaining decision rule:** **Generative:** Estimate $P(w_i|x)$ and $P(w_i)$ then apply Bayes. **Discriminative:** Learn decision boundary or the posterior prob, directly w/o underlying distributions. **Generalization:** Good in unseen data, **Memorization:** Overfit.

• $\Delta(x) = \frac{p(x|\omega_1)}{p(x|\omega_2)} > \frac{P(\omega_2)}{P(\omega_1)}$

• $c = 2$ (two class problem):

• $R = \lambda_{12} \int_{\mathcal{R}_1} p(x|\omega_2)P(\omega_2)dx + \lambda_{21} \int_{\mathcal{R}_2} p(x|\omega_1)P(\omega_1)dx$

Resubstitution: whole set X used for training and testing. Corresponds to setting the error estimate on the training set (apparent error). This is an optimistic and low variance error estimate. **Holdout:** total dataset randomly or not divided into two disjoint sets, one for training and one for testing. Improve testing error estimate by doing different runs.

Stratified cross-validation: ensure class prevalence along folds to handle imbalanced data. **Bootstrap:** generate artificial data samples by randomly selecting existing samples with uniform distribution. **Three set validation:** we use training and validation subsets to choose the best model (ex. Set of parameters) and then test that on a never-before-used testing subset. **ROC:** plots the proportion of correct responses (TP or sensitivity) against the FP rate as a threshold changes. Represents trade-off between sensitivity and specificity. A perfect classifier would have $AUC=1$. Random classifier would have a diagonal ROC (45° degree). **Parzen Window:** kernel density estimation. Creates histograms. **Non-parametric:** KNN: consider k points that are about data distribution) and a **lazy classifier** because training only implies storing the feature vectors and class labels of training samples. This is good because it requires less computation

for training. Also it's instance based therefore adapts locally to data structure (good for non-linearly separable problems).

Filters: pre-processing methods independent from classifier; simple features assume independence between features and only look at relevance (discriminative power.) Ideally it would be combined with a redundancy analysis (eg. minimum redundancy maximum relevancy algorithm (mRMR). **Advantages:** fast and computationally efficient, may provide a generic selection independent of classifier. **Disadvantages:** assume independence between features. Based just on data, not in classification performance. **Wrappers:** based on classifier performance: we train a model with different subsets of features and choose the one that yields best performance. **Advantages:** high accuracy, overfitting control associated with classifier. **Disadvantages:** computationally expensive and lacks generality. **Embedded:** use all features and define their relevance internally during learning. More important features affect the internal parameters of the classifier more.

VC dimension, h: largest number of samples that can be correctly separated (shattered) by a classifier during training, considering all possible ways to label these samples. **Structural risk:** Structural risk refers to minimizing a bound on the true (generalization) error by balancing empirical risk with model complexity, selecting among nested function classes so as to control overfitting and achieve the best trade-off between fit to data and capacity.

SVM: defining an optimal hyperplane that minimizes the training error and maximizes the separation margin among classes (minimizing structural risk.) It can extend this definition to non-linearly separable problems by having a penalty for misclassifications (soft-margin.) It can also use kernel trick to map data to high dimensional space where it's easier to classify with linear decision surfaces.

Margin, r : r is usually 1. Support vectors are those closest to decision boundary, aka $w^T x + b = \pm r$. **O SVM resolve o problema convexo que equivale a maximizar a margem entre as duas classes sob a restrição de classificar corretamente com margem pelo menos 1.** Resolve-se com o Lagrange multipliers method: apenas support vectors ficam com lagrange multipliers não nulos.

SVM soft-margin: this method handles non linearly separable data by accounting the deviations from the ideal separable problem, it allows patterns to fall within the margin and in the opposite region but penalizes them. It adds a **free parameter, C** , that influences margin: **higher C means smaller margins**, aka hard penalization of wrong classifications, but can overfit. **Smaller C means larger margins**, soft penalization of the wrong classification.

Kernel Trick: we can use functions to transform data to higher dimensional spaces to make problems linearly separable. But this can be computationally expensive to compute exhaustively. Kernel functions allow us to use SVMs without explicitly mapping the data; in the svm formulas we simply substitute inner products $x_i^T x_j$ for the kernel function $K(x_i, x_j) = f(x_i) \cdot f(x_j)$. **Gaussian kernel parameter:** this γ parameter controls the locality of influence of support vectors: a large γ makes each support vector affect only a small neighborhood (complex, wiggly boundary) smoother decision surface. **Multi-class SVM:** decompose problem in binary problems using one vs all or one vs one. We can use the **Error-Correcting Output**

design matrix to determine the classes that the binary classifiers will try to discriminate and set a base against which will decide the final design. To decode, the most **distance** between classifiers output and each line of the coding design matrix and attribute the class associated to the minimum distance. Hamming distance is count the number of different cells. **Get non-linear SVM w and w0 given support vectors without lagrange multipliers:** Identify SVs x_1 and x_2 in the original plane (on top of boundaries) -> map them to z_1 and z_2 with kernel 2 -> bias = $(z_1 - z_2)/2$ -> weight vector according to the slope of the line in the middle of them. That will give us $wz_1 + wz_2 + b$, change to $w(x_1^2 + x_2^2) + b$ if the kernel transformation is $f(x) = (x_1^2 + x_2^2)^{1/2}$. **Hybrid Classifier:** combines two or more different classification methods to exploit their complementary strengths and improve overall performance.

• $\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \varphi\left(\frac{x-x_i}{h}\right)$

Occam's Razor: when you have two competing theories which make exactly the same predictions, the simpler is the better. Aka, we should select the simplest model because it minimizes the number of incorrect assumptions. **No Free Lunch Theorem:** no algorithm is superior to any other over all possible classification problems. **Bias:** error rate due to the average performance of the classifier, how far is the average performance from optimal?. **Variance:** error rate due to training set changes, how does it vary for different training sets? **Bagging:** draw n' samples with replacement from D several times, use each set to train a classifier (parallel), combine results of each classifier to the final classification usually by majority voting. This decreases variance (use case: random forests.)

Boosting: learn a strong classifier by combining weak classifiers (classifiers with low bias and variance, only slightly better than random guessing.) We learn weak classifiers iteratively (not parallel) and reweight data

10: Final classifier: $H(x) = \text{sign}\left(\sum_{j=1}^M \alpha_j H_j(x)\right)$. The final decision is a weighted sum of all the weak learners. (use case: AdaBoost.) **Stacking:** combine several heterogeneous models by training an additional classifier with the individual decision function values. Learning from diverse predictions may lead to improved overall classification. **Generalized Probabilistic Classifier Combination:** combine results of R classifiers where each classifier k outputs the probabilities of each class $P_k(w_j|x)$. We can apply the **Product Rule** (product of every likelihood, sensitive to very low values) or the **sum rule** (sum of every likelihood), sum rule found to perform well in practice. If $\prod(k=1 \rightarrow K)P(w_j|x) > \prod(k=1 \rightarrow K)P(w_i|x)$ assign w_j , else assign w_i . **Advantages vs Disadvantages LDA:** advantages: computationally efficient, noise reduction, enhance classifier performance by maximizing separability, can be easily interpretable. Disadvantages: assumes data has the same covariance matrix, works best when linearly separable, difficult with some layouts (ex. Above). can produce at most C-1 discriminants. **Clustering, Types:** Hierarchical/Tree: Use linkage rules to create hierarchical sequence of clustering solutions. **Centroid adjustment algorithms / partitioning:** Based on an initial proposal for the prototypes (centroids), the algo iteratively evaluates and adjusts centroid position. **Density-based:** Join areas of high pattern density into clusters. Good for arbitrarily-shaped clusters because it can be formed as long as dense areas can be connected. **Distribution-based:** Estimate parametric data distributions from data, e.g., gaussian. **Others:** **Hierarchical Clustering:** Agglomerative: Start w/ N clusters where each cluster is composed of a single sample -> Merge two nearest cluster -> repeat until all clusters are part of a single cluster -> output dendrogram. **Divisive:** Start with all in 1 cluster -> split into 2 ... **Linkage rules:** What will be the new distance between the join clusters and other clusters? **Single:** Dissimilarity given by the closest elements (good for clusters with filamentary shape, promote connectedness). **Complete:** Dissimilarity given by the furthest elements (good for compact clusters, original spherical-shaped clusters). **Average:** Pair-group using arithmetic average (effective for several types of clusters). **Centroid:** Euclidean distance between centroids of the two clusters (effective for several types of clusters). **Ward's method:** Sum of the squared within-cluster distances (minimize intra-cluster variability and consequently maximize cluster separability). **Lance-Williams formula:** update the distance matrix in agglomerative methods after merging clusters (the d in the formula can be the distance for any of the previously mentioned methods).

Centroid adjustment algorithms (K-MEANS): Assign data point to clusters -> move centroid describing the clusters. Number of clusters need to be defined a-priori (we can use hierarchical clustering in a small portion of data for initialisation. Minimize and overall within cluster distance from the patterns to the centroids (loss function to the right). K-means works if: Clusters are spherical -> clusters are of similar volumes -> clusters have similar number of points **K-means Alg:** assign each datapoint to the cluster represented by the nearest centroid -> compute, for the previous partition, the new centroids -> repeat until max num of iteration exceeded or there are no more changes on the centroids. **Density Algorithms (DBSCAN):** Clusters are dense areas separated by areas of lower density. Defined by two parameters: ϵ (radius of neighborhood with respect to some point) minPts (min number of points required to consider a region dense). Advantages: robust to outliers (points in very low populated areas); do not need an initial guess for the number of clusters; deal with arbitrary cluster shapes. **Point Types:** **Core Point:** If at least minPts points are within a distance ϵ of it (including the point itself). **Border Point:** a point that has at least one core point at a distance ϵ and if the number of points inside the volume defined by ϵ is less than minPoint. **Noise Point:** Not core nor border. **Reachability:** If a data point can be accessed from another datapoint (directly or indirectly). **Connectivity:** If two datapoints belong to the same cluster or not. **Two points in DBSCAN can be:** Directly density-reachable if $d(y, x) < \epsilon$ and y is core point, density-reachable if they can be reached indirectly, density connected if both x and y are density-reachable from O. **DBSCAN Alg:** Pick a point not assigned to a cluster or classified as outlier, implement neighborhood jumps to find all points in its neighborhood, repeat until all points are labeled as cluster members.

Single Link $a_i = a_j = 0.5; b = c = -0.5$
Complete Link $a_i = a_j = 0.5; b = 0; c = 0.5$
Average Link $a_i = \frac{n_j}{n_i+n_j}; a_j = \frac{n_i}{n_i+n_j}, b = \frac{n_i n_j}{n_i+n_j}, c = 0$
Centroid $a_i = \frac{n_j}{n_i+n_j}; a_j = \frac{n_i}{n_i+n_j}, b = -\frac{n_i n_j}{n_i+n_j}, c = 0$
Ward's Method $a_i = \frac{n_j+n_k}{n_i+n_j+n_k}; a_j = \frac{n_j+n_k}{n_i+n_j+n_k}; b = -\frac{n_i}{n_i+n_j+n_k}; c = 0$

repeat until all points are labeled as cluster members. If we connected, cutting the dendrogram at height calculate distance to the k-th nearest part of a cluster). **Cluster evaluation:** (large between clusters distances) (e.g. Davies-Bouldin idx, Silhouette, Dunn Index, ...). **Supervised/external:** Need labels to evaluate. (e.g. Purity / Accuracy, Precision, Recall, ...).

$$\diamond d(w_{i+j}, w_k) = a_i d(w_i, w_k) + a_j d(w_j, w_k) + bd(w_i, w_j) + c |d(w_i, w_k) - d(w_j, w_k)|$$

Performance Evaluation. Confusion Matrix with TP, TN, FP, FN. Accuracy = $(TP+TN)/(TP+TN+FP+FN)$. Sensitivity/Recall = $TP/(TP+FN)$. Specificity = $TN/(TN+FP)$. F1-score = $2 \cdot (Precision \cdot Recall) / (Precision + Recall)$. Macro averaging computes metrics per class and averages them (class-balanced); Micro averaging computes metrics globally and is dominated by majority classes (\approx accuracy). For class imbalance: Binary \rightarrow use Recall, Specificity, F1 or Balanced Accuracy; Multiclass \rightarrow use BACC/BACCw. **Decision Trees.** Classification through recursive if-then rules; each root-to-leaf path is a conjunction, the full tree is a disjunction. Advantages: interpretable, fast, handles mixed data. Disadvantages: greedy search, high variance, sensitive to data changes. **Entropy and Information Gain (ID3).** Entropy measures class heterogeneity: $H(S) = -\sum p_i \log(p_i)$. Information Gain: $Gain(S,A) = H(S) - \sum_i |S_i|/|S| H(S_i)$. ID3 selects the feature with highest IG until entropy is zero. IG is biased toward features with many values and sensitive to noise and class imbalance. **CART.** Uses Gini impurity $G = 1 - \sum p_i^2$, binary splits only, supports classification and regression. Overfitting controlled by pruning, depth limits, minimum samples per leaf, cross-validation or Random Forests. **Naive Bayes.** Assumes conditional independence between features given the class. Likelihood Ratio Test: $p(x|w_1)p(x|w_2) > P(w_1)P(w_2) \Rightarrow$ choose w_1 . Capital P denotes discrete probabilities, lowercase p denotes probability density functions; $p(x)$ is the evidence. **Bias–Variance Decomposition.** Total error = bias + variance. High bias \rightarrow underfitting; high variance \rightarrow overfitting. Overfitting corresponds to very low training error and high testing error, often due to high model complexity (high VC dimension). **Feature Selection Additions.** ROC/AUC can be used to assess individual feature relevance; higher AUC implies better discriminative power. Regularization as embedded selection: LASSO (L1) forces irrelevant feature weights to zero; Ridge (L2) shrinks weights to reduce variance. **Distance Metrics.** A valid metric satisfies non-negativity, identity, symmetry and triangle inequality. A norm also satisfies $d(ax,ay)=|a|d(x,y)$. **k-NN Theory.** As $N \rightarrow \infty$, $k \rightarrow \infty$ and $k/N \rightarrow 0$, k-NN approaches the Optimal Bayes classifier. Small k yields low bias and high variance (overfitting).

```
class Node:
    def __init__(self, f=None, t=None, l=None, r=None, v=None):
        self.f, self.t, self.l, self.r, self.v = f, t, l, r, v
    def gini(y):
        if len(y) == 0: return 0
        p = np.bincount(y) / len(y)
        return 1 - np.sum(p**2)
def _fit_tree(X, y, depth=0, max_d=5):
    if depth == max_d or len(np.unique(y)) == 1:
        return Node(v=Counter(y).most_common(1)[0][0])
    # Random Subspace: Select sqrt(D) features
    n_feat = X.shape[1]
    # feats = np.random.choice(n_feat, int(np.sqrt(n_feat)), replace=False)
    best = {'f': float('inf'), 't': None, 'l': None}
    for f in feats:
        thresholds = np.unique(X[:, f])
        for t in thresholds:
            mask = X[:, f] < t
            if np.sum(mask) == 0 or np.sum(~mask) == 0: continue
            # Weighted Gini Impurity
            g = (len(y[mask]) * gini(y[mask]) + len(y~mask)) * gini(y~mask)
            if g < best['g']:
                best = {'f': f, 't': t, 'l': t}
    if best['f'] is None: return Node(v=Counter(y).most_common(1)[0][0])
    mask = X[:, best['f']] <= best['t']
    return Node(best['f'], best['t'], _fit_tree(X[mask], y[mask], depth+1, max_d),
                _fit_tree(X~mask, y~mask, depth+1, max_d))
def _predict(n, x):
    return n.v if n.v is not None else _predict(n.l if x[n.f] < n.t else n.r, x)
return n.v if n.v is not None else _predict(n.l if x[n.f] < n.t else n.r, x),
```

```
def rf_training(Xtr, Ttr):
    # Xtr is D x P, Transpose to P x D for processing
    X, y = Xtr.T, Ttr.flatten().astype(int)
    trees = []
    # Bagging: Train 10 trees on bootstrap samples
    for _ in range(10):
        idx = np.random.choice(len(y), len(y), replace=True)
        trees.append(_fit_tree(X[idx], y[idx]))
    return trees
def rf_testing(Xte, Tte, model):
    X, y = Xte.T, Tte.flatten().astype(int)
    # Majority Voting
    preds = [_predict(t, x) for x in X for t in model['trees']]
    preds = np.array(preds).reshape(len(X), len(model['trees']))
    out = np.array([Counter(row).most_common(1)[0][0] for row in preds])
    # Metrics
    l1 = Pos, 2-Neg
    TP = np.sum((out == 1) & (y == 1))
    TN = np.sum((out == 2) & (y == 2))
    ss = TP / np.sum(y == 1) if np.sum(y == 1) else 0.0
    sp = TN / np.sum(y == 2) if np.sum(y == 2) else 0.0
    return out, ss, sp
```

```
def agg_clustering_avg(dist_matrix, labels) -> list[list[Any]]:
    D = np.array(dist_matrix, dtype=float)
    fill_diagonal(D, np.inf)
    clusters = labels[:]
    sizes = [1] * len(labels)
    history = [[list(clusters)]]
    while len(clusters) > 1:
        # 1. Find the closest pair of clusters
        i, j = np.unravel_index(np.argmin(D), D.shape)
        if i > j: i, j = j, i # Ensure i < j for consistent deletion
        # 2. Update Distances (Average Linkage: Weighted average)
        # Formula: d(i, k) = (n_i * d(i, k) + n_j * d(j, k)) / (n_i + n_j)
        n_i, n_j = sizes[i], sizes[j]
        D[i, :] = (n_i * D[i, :] + n_j * D[j, :]) / (n_i + n_j)
        D[:, j] = D[i, :]
        # Maintain symmetry
        D[i, j] = np.inf
        # Reset diagonal
        D[i, i] = 0
        # 3. Merge clusters and sizes
        clusters[i] += clusters[j] (clusters[i]) (clusters[j])
        sizes[i] += sizes[j]
        # 4. Remove cluster j (delete row and column)
        D = np.delete(np.delete(D, j, 0), j, 1)
        clusters.pop(j)
        sizes.pop(j)
        history.append(list(clusters))
    return history
# --- Test Case ---
dists = [
    {0, 754, 564, 138}, # (0, 754, 219, 869),
    {754, 0, 219, 869}, # (564, 219, 0, 669),
    {564, 219, 0, 669}, # (138, 869, 669, 0)
]
lbls = ['Milano', 'Napoli', 'Rome', 'Torino']
```

TEST LDA FISHER

```
def predict(self, X_te):
    X_te = np.asarray(X_te)
    X_te_proj = X_te.dot(self.projection_vector)
    # Classify based on minimum Euclidean distance in 1D
    predictions = np.where(np.abs(X_te[:, 0] - self.m0_proj) < abs(X_te[:, 1] - self.m1_proj)) else 1 for x in X_te_proj]
    return np.array(predictions)
```

Test MDCs

```
def min_dist_testing(Xte, Tte, model):
    Xte = np.asarray(Xte), Tte = np.asarray(Tte)
```

```
D, Nte = Xte.shape
```

```
mu1 = model['mu1'], mu2 = model['mu2']
```

```
metric = model['metric']
```

```
preds = np.empty(Nte, dtype=int)
```

```
if metric == 'euc':
    for j in range(Nte):
        x = Xte[:, j]
        d1 = np.sum((x - mu1)**2), d2 = np.sum((x - mu2)**2)
        preds[j] = 1 if d1 < d2 else 2
else: # 'manh'
    invC = model['inv_cov']
    for j in range(Nte):
        x = Xte[:, j], diff1 = x - mu1, diff2 = x - mu2
        d1 = diff1 @ invC @ diff1
        d2 = diff2 @ invC @ diff2
        preds[j] = 1 if d1 < d2 else 2
tp, tn, fp, fn = get_confusion(preds, Tte, pos_label=1)
```

Find C for SVM:

```
X = np.asarray(Xtr)
y = np.asarray(Ytr)
```

```
mean_f1 = []
std_f1 = []
skf = StratifiedKFold(n_splits=n_runs, shuffle=True)
```

```
for C in Cs:
    f1_scores = []
    for train_idx, test_idx in skf.split(X, y):
        X_train, X_val = X[train_idx], X[test_idx]
        y_train, y_val = y[train_idx], y[test_idx]
        clf = LinearSVC(C=C, max_iter=5000, dual=False)
        clf.fit(X_train, y_train)
        y_pred = clf.predict(X_val)
        f1 = f1_score(y_val, y_pred, pos_label=1)
        f1_scores.append(f1)
    mean_f1.append(np.mean(f1_scores))
    std_f1.append(np.std(f1_scores, ddof=1))
return np.array(mean_f1), np.array(std_f1)
```

Train Bayes for 1 and 2 labels

```
X = np.asarray(Xtr)
y = np.asarray(Ytr)
# assume rótulos exatamente 1 e 2
mu1 = X[y==1].mean(axis=0)
mu2 = X[y==2].mean(axis=0)
cov1 = np.cov(X[y==1], rowvar=False) + eps*np.eye(X.shape[1])
cov2 = np.cov(X[y==2], rowvar=False) + eps*np.eye(X.shape[1])
prior1 = np.mean(y==1)
prior2 = np.mean(y==2)
# pre-calc inversas e determinantes para pdf
inv1 = np.linalg.inv(cov1); det1 = np.linalg.det(cov1)
inv2 = np.linalg.inv(cov2); det2 = np.linalg.det(cov2)
return {
    'mu1': mu1, 'inv1': inv1, 'det1': det1, 'prior1': prior1,
    'mu2': mu2, 'inv2': inv2, 'det2': det2, 'prior2': prior2,
    'dim': X.shape[1]
```

```
def svm_testing(Xte, Tte, model) -> tuple[NDArray[Any], NDArray[Any]]:
    X, y = Xte.T, Tte.flatten()
    # Decision function: w^T * x + b [Slide 15]
    decision = np.dot(X, model['w']) + model['b']
    out = np.where(decision >= 0, 1, 2)
    TP = np.sum((out == 1) & (y == 1))
    TN = np.sum((out == 2) & (y == 2))
    ss = TP / np.sum(y == 1) if np.sum(y == 1) else 0.0
    sp = TN / np.sum(y == 2) if np.sum(y == 2) else 0.0
    return out, ss, sp
```

Train MDC, can be eucl or mah

```
def min_dist_training(Xtr, Ttr, metric='euc', eps=1e-6):
    Xtr = np.asarray(Xtr)
    Ttr = np.asarray(Ttr)
    D, Ntr = Xtr.shape
    # separar as duas classes
    mask1 = (Ttr == 1)
    mask2 = (Ttr == 2)
    # calcular médias (vetores de dimensão D)
    mu1 = Xtr[:, mask1].mean(axis=1)
    mu2 = Xtr[:, mask2].mean(axis=1)
    model = {'mu1': mu1, 'mu2': mu2, 'metric': metric}
    if metric == 'mah':
        # covariâncias amostrais com regularização eps*I
        X1 = Xtr[:, mask1]
        X2 = Xtr[:, mask2]
        # np.cov com rowvar=True: linhas=variáveis, colunas=amostras
        cov1 = np.cov(X1, rowvar=True) + eps*np.eye(D)
        cov2 = np.cov(X2, rowvar=True) + eps*np.eye(D)
        n1 = X1.shape[1]
        n2 = X2.shape[1]
        # covariância pooled (estimador da máxima verossimilhança ou amostral)
        pooled = ((n1-1)*cov1 + (n2-1)*cov2) / (n1+n2-2)
        inv_pooled = np.linalg.inv(pooled)
        model['inv_cov'] = inv_pooled
    return model
```

TRAIN LDA FISHER

```
def train(self, X, train, Y_train):
    X_train = np.asarray(X_train)
    Y_train = np.asarray(Y_train)
    # Compute class means and scatter matrices
    class0 = X_train[Y_train == 0]
    class1 = X_train[Y_train == 1]
    self.m0 = np.mean(class0, axis=0)
    self.m1 = np.mean(class1, axis=0)
    # Compute within-class scatter matrix Sw
    S0 = np.cov(class0, rowvar=False)
    S1 = np.cov(class1, rowvar=False)
    Sw = S0 + S1
    # Compute LDA projection vector
    self.projection_vector = np.linalg.inv(Sw).dot(self.m1 - self.m0)
    self.m0_proj = np.dot(self.m0, self.projection_vector)
    self.m1_proj = np.dot(self.m1, self.projection_vector)
```

Test Bayes

```
y = np.asarray(Yte)
n, d = X.shape
m = model
# função inline de pdf multivariada Gaussiana
def pdf(x, mu, invC, detC):
    diff = x - mu
    return np.exp(-0.5*(diff @ invC @ diff)) / np.sqrt((2*np.pi)**d * detC)
# calcular posteriors para todas amostras de uma vez
# mas pdf é por ponto; fazemos em loop curto
out = np.empty(n, dtype=int)
for i in range(n):
    x = X[i]
    p1 = pdf(x, m['mu1'], m['inv1'], m['det1']) * m['prior1']
    p2 = pdf(x, m['mu2'], m['inv2'], m['det2']) * m['prior2']
    out[i] = 1 if p1 > p2 else 2
# sensibilidade e especificidadeide vetorializadas
ss = np.mean(out==1) == 1 if np.any(y==1) else 0.0
sp = np.mean(out==2) == 2 if np.any(y==2) else 0.0
def svm_training(Xtr, Ttr) -> dict[str, Any]:
    # Transpose to PxD; Labels: 1 -> 1, 2 -> -1
    X, y = Xtr.T, np.where(Ttr.flatten() == 1, 1, -1)
    n_samples, n_features = X.shape
    # Initialize Lagrange Multipliers (alpha)
    alpha = np.zeros(n_samples)
    b = 0
    C = 1.0 # Regularization parameter (Slide 27)
    # Simplified SMO (Coordinate Descent on Dual)
    # Iterates to satisfy KKT conditions (Slide 393)
    for _ in range(50):
        for i in range(n_samples):
            # Calculate margin: f(x) - y
            w = sum(alpha * y * X[i]) [Slide 21]
            error = np.dot(X[i], w) - y[i]
            # Update alpha if strict KKT conditions are violated
            if (y[i] * error < -0.01 and alpha[i] < C) or \
               (y[i] * error > 0.01 and alpha[i] > 0):
                # Gradient update for alpha (Simplified)
                step = 0.01
                alpha[i] = np.clip(alpha[i] + step * y[i] * error, -b, C)
    # Compute final w from alphas (Slide 21)
    w = np.dot(alpha * y, X)
    # Save model with Support Vectors (where alpha > 0)
    return {'w': w, 'b': b, 'alpha': alpha}
```

adaboost_training

```
X, y = Xtr.T, np.where(Ttr.flatten() == 1, 1, -1)
n_samples, n_features = X.shape
# Initialize weights uniformly
w = np.full(n_samples, 1 / n_samples)
models = []
# Train 50 Decision Stumps (Weak Learners)
for _ in range(50):
    stump = {'err': float('inf'), 'alpha': 0}
    # Greedy search for best feature/threshold minimizing weighted error
    for f in range(n_features):
        thresh = np.unique(X[:, f])
        for t in thresh:
            p = 1 # Polarity
            preds = np.ones(n_samples)
            preds[X[:, f] < t] = -1
            # Calculate weighted error
            err = np.sum(w[preds != y]) # Invert polarity if error > 0.5
            if err > 0.5:
                p, err = -1, 1 - err
            if err < stump['err']:
                stump['err'] = err
                stump['f'] = f
                stump['t'] = t
                stump['p'] = p
    # Calculate Alpha (Vote Power)
    EPS = 1e-10
    stump['alpha'] = 0.5 * np.log((1 - stump['err']) + EPS) / (stump['err'] + EPS)
    # Update Weights: Increase weight of misclassified samples
    preds = stump['p'] * np.where(X[:, stump['f']] < stump['t'], -1, 1)
    w *= np.exp(stump['alpha'] * y * preds)
    w /= np.sum(w) # Normalize
    models.append(stump)
return models
```

adaboost_testing

```
X, y = Xte.T, Tte.flatten()
# Aggregate predictions: sign(sum(alpha * prediction))
final_preds = np.zeros(len(X))
for m in model['models']:
    pred = m['p'] * np.where(X[:, m['f']] < m['t'], -1, 1)
    final_preds += m['alpha'] * pred
out = np.where(np.sign(final_preds) >= 0, 1, 2)
TP = np.sum((out == 1) & (y == 1))
TN = np.sum((out == 2) & (y == 2))
ss = TP / np.sum(y == 1) if np.sum(y == 1) else 0.0
sp = TN / np.sum(y == 2) if np.sum(y == 2) else 0.0
return out, ss, sp
```

Head Sheets

```
description = pd.read_excel(file_name, sheet_name="Description")
data = pd.read_excel(file_name, sheet_name="Data")
```

Extract features and labels

Assume last two columns are Apgar1 and Apgar5

X = data.iloc[:, :-2].values

apgar1 = data["Apgar1"].values

Binary labels: bad prognosis if Apgar1 < 6

y = (apgar1 <= 6).astype(int)

(b) Feature ranking using AUC

auc_scores = []

for i in range(X.shape[1]):
 auc = roc_auc_score(y, X[:, i])
 auc_scores.append(auc)

auc_scores = np.array(auc_scores)

Sort features by descending AUC
sorted_idx = np.argsort(auc_scores)[-1:-1]

Select two most discriminatory features
best_features_idx = sorted_idx[:2]
X_best = X[:, best_features_idx]

print("Selected feature indices:", best_features_idx)
print("Corresponding AUCs:", auc_scores[best_features_idx])

(c) (d) (e) k-NN classification with 10 trials

errors = []

for trial in range(10):
 X_train, X_test, y_train, y_test = train_test_split(
 X, y, test_size=0.5, shuffle=True)

k-NN classifier with k = 8
knn = KNeighborsClassifier(n_neighbors=8)
knn.fit(X_train, y_train)

y_pred = knn.predict(X_test)

Classification error
error = np.mean(y_pred != y_test)

errors.append(error)

Results

errors = np.array(errors)

KNN PDF

```
def knn_pdf(data, k, eps=1e-8):
    x = np.asarray(data).ravel()
    n = x.size
    pdf = np.empty(n, dtype=float)
    for i in range(n):
        xi = x[i]
        # obtém índices e distâncias dos k vizinhos mais próximos a xi
        idxs, dists = knn1d.search(xi, k)
        # A distância ao k-ésimo vizinho é:
        if len(idxs) < k:
            # se houver menos pontos que k, usamos a maior distância disponível
            r = np.max(dists) if dists.size > 0 else eps
        else:
            # supondo dists ordenados ou não: garantimos pegar o k-ésimo menor
            r = np.partition(dists, k-1)[k-1]
        if r < eps:
            r = eps
        pdf[i] = k / (n * 2 * r)
    return pdf
```