

Classification

Ruwen Qin
Stony Brook University

October 2, 2025

Data used in this module:

- ai4i2020.csv

Python notebook used in this module

- Classification.ipynb

1 Introduction

Considering an entity characterized by an N dimensional feature representation: $\mathbf{X} = [X_1, \dots, X_N]^T$. We would like to infer a target variable of interest, Y , which takes values from a finite set of K categories, $\{c_1, \dots, c_K\}$. Classification is about developing a classifier: $h: \mathbf{X} \rightarrow Y$, which can assign an observation with feature values $\mathbf{x} = [x_1, \dots, x_N]^T$ to a class label y .

There are many classification problems in the real world. For example, given an image captured by a robot in infrastructure inspection, a classifier can determine if the image contains defect. Given measures of machine attributes in operations, a classifier can predict the type of failure the machine will confront. This learning module is developed based on related chapters in [1, 2, 3, 4, 5].

1.1 Bayes Classifier

Define the probability that the class of the response is c_k conditional on the observed feature values is:

$$p_k = P(c_k|\mathbf{x}), \quad (1)$$

and $\sum_{k=1}^K p_k = 1$.

According to Bayes' theorem,

$$p_k = \frac{f_k(\mathbf{x})\pi_k}{\sum_{r=1}^K f_r(\mathbf{x})\pi_r} \quad (2)$$

where π_k is the prior probability:

$$\pi_k = P(c_k), \quad (3)$$

and $f_k(\mathbf{x})$ is the likelihood probability,

$$f_k(\mathbf{x}) = P(\mathbf{x}|c_k). \quad (4)$$

The Bayes classification rule says that the optimal rule is one that picks the class with the maximum posterior probability:

$$h^*(\mathbf{x}) = \arg \max_k p_k = \arg \max_k f_k(\mathbf{x})\pi_k. \quad (5)$$

1.2 Three Approaches to Classification

The Bayes rule depends on unknown quantities $f_k(\mathbf{x})$ and π_k for $k = 1, \dots, K$. Therefore, we need to use the data to find some approximation to the Bayes rule.

Generally speaking, there are three major approaches:

- The first approach is density estimation. We use data to estimate the prior and likelihood probabilities and then compute the posterior probability using the Bayes's theorem.
- The second approach is to model the conditional probability p_k directly, for example, as a parametric model.
- The last approach involves constructing a discriminant function that directly assign each vector to a specific class.

Classifiers can be developed in one of these approaches or a mix of them. In this learning module, we will introduce several classification methods.

2 Quadratic and Linear Discriminant Analysis

One approach to classification is to use the density estimation strategy and assume a parametric model for the densities.

Assume $f_k(\mathbf{x})$ are multivariate Gaussians:

$$f_k(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^N |\Sigma_k|}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)\Sigma_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)^T \right\} \quad (6)$$

for $k = 1, \dots, K$, where

$$\mathbf{x} = [x_1, \dots, x_N]^T \quad (7)$$

is the feature vector of a data point,

$$\boldsymbol{\mu}_k = \mathbb{E}[\mathbf{x}|c_k] = [\mu_{k,1}, \dots, \mu_{k,N}]^T \quad (8)$$

is the mean feature vector for class k , and

$$\Sigma_k = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu}_k)^T(\mathbf{x} - \boldsymbol{\mu}_k)|c_k] \quad (9)$$

is the class- k feature covariance matrix. $|\cdot|$ means the determinant of a matrix.

2.1 Quadratic Discriminant Analysis (QDA)

According to (6),

$$\log [f_k(\mathbf{x})\pi_k] = -\frac{N}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)\Sigma_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)^T + \log \pi_k, \quad (10)$$

where the first term on the right hand side, $-\frac{N}{2} \log 2\pi$, is a constant. Therefore, we drop it and define the discriminant function:

$$\delta_k(\mathbf{x}) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)\Sigma_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)^T + \log \pi_k. \quad (11)$$

$(\mathbf{x} - \boldsymbol{\mu}_k)\Sigma_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)^T$ in (11) is called Mahalanobis distance, measuring a point relative to a distribution.

The Bayes rule assigns \mathbf{x} to a specific class:

$$h^*(\mathbf{x}) = \arg \max_k \delta_k(\mathbf{x}). \quad (12)$$

The discriminant function in Eq. (11) shows that, in the feature space,

- the smaller the distance from a data point to the mean of class k , the larger the discriminant function value, if all else are the same.
- the smaller the determinant of class- k co-variance matrix, the tighter the distribution, the larger the discriminant function, if all else are the same.
- the larger the prior probability of class k , the larger the discriminant function value, if all else is the same.

The discriminant function $\delta_k(\mathbf{x})$ in (11) is quadratic. Therefore, the classification procedure in (11) and (12) is called Quadratic Discriminant Analysis (QDA).

Given a training sample with M data points, we can estimate the prior probability π_k , and parameters of the likelihood probability, Σ_k and $\boldsymbol{\mu}_k$, which are required for determining the discriminant function $\delta_k(\mathbf{x})$ in (11).

Denote by $\mathbf{X} \in \mathbb{R}^{N \times M}$ the feature values of training dataset. We partition \mathbf{X} by classes: $\mathbf{X} = \cup_{k=1}^K \mathbf{X}_k$, where $\mathbf{X}_k = \{\mathbf{x}_i|c_k\}$ represents the attributes of data points in class c_k and M_k is the number of such data points:

$$M_k = \sum_{i=1}^M 1\{y_i = c_k\}. \quad (13)$$

Then, for each class k , the prior probability is estimated as:

$$\hat{\pi}_k = \frac{M_k}{M}, \quad (14)$$

and the parameters of the likelihood function are estimated as:

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{M_k} \sum_{i=1}^M 1\{y_i = c_k\} \mathbf{x}_i \quad (15)$$

$$\hat{\Sigma}_k = \frac{1}{M_k - 1} (\mathbf{X}_k - \hat{\boldsymbol{\mu}}_k)^T (\mathbf{X}_k - \hat{\boldsymbol{\mu}}_k). \quad (16)$$

2.2 Linear Discriminant Analysis (LDA)

If $\Sigma_1 = \dots = \Sigma_K = \Sigma$, that is, features have the same covariance matrix across all classes, the discriminant function becomes linear. Let's substitute Σ for Σ_k in the log function of the posterior probability, leading to

$$\begin{aligned}\log f_k(\mathbf{x})\pi_k &= -\frac{N}{2}\log 2\pi - \frac{1}{2}\log |\Sigma| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) + \log \pi_k \\ &= -\frac{N}{2}\log 2\pi - \frac{1}{2}\log |\Sigma| - \frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x} + \mathbf{x}^\top \Sigma^{-1}\boldsymbol{\mu}_k - \frac{1}{2}\boldsymbol{\mu}_k^\top \Sigma^{-1}\boldsymbol{\mu}_k + \log \pi_k.\end{aligned}\quad (17)$$

We drop the first three terms on the right hand side, which are class-independent, to obtain the updated discriminant function:

$$\delta_k(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1}\boldsymbol{\mu}_k - \frac{1}{2}\boldsymbol{\mu}_k^\top \Sigma^{-1}\boldsymbol{\mu}_k + \log \pi_k, \quad (18)$$

which is linear in \mathbf{x} .

$\boldsymbol{\mu}_k$ in (18) can be estimated using (15) and Σ can be estimated with

$$\hat{\Sigma} = \frac{1}{M-1}(\mathbf{X} - \hat{\boldsymbol{\mu}})^\top (\mathbf{X} - \hat{\boldsymbol{\mu}}), \quad (19)$$

where $\hat{\boldsymbol{\mu}}$ is the estimator for the mean attribute vector using the entire training dataset:

$$\hat{\boldsymbol{\mu}} = \frac{1}{M} \sum_{i=1}^M \mathbf{x}_i = \frac{1}{M} \sum_{k=1}^K M_k \hat{\boldsymbol{\mu}}_k. \quad (20)$$

In the learning module “Feature Extraction”, we also introduced LDA and derived it in a different approach called Fisher’s Linear Discriminant.

3 Naive Bayes (NB)

3.1 Naive Bayes Classifier

Estimating π_k is simple, but estimating the likelihood probability $f_k(\mathbf{x})$ is not. In QDA and LDA, $f_k(\mathbf{x})$ is assumed a multivariate Gaussian, leading to the requirement of estimating its parameters $\boldsymbol{\mu}_k$ and Σ_k . A more general approach is to estimate $f_k(\mathbf{x})$ with some non-parametric density estimator. However, if \mathbf{x} is in high dimension, non-parametric density estimation is not reliable.

If we assume feature variables are independent, we can just estimate N 1-dimensional density functions. The joint distribution is simply the product of these 1-dimensional density functions. Although these are apparently over-simplified assumptions, Naive Bayes (NB) classifiers have worked quite well in many real-world situations.

NB classifiers are implemented below.

1. Estimate the density for the j th attribute pertaining to class k , $\hat{f}_{k,j}(x_j) = P(x_j|c_k)$, with a density estimator
2. The estimator of the likelihood probability is the product of the N 1-dimensional density estimators: $\hat{f}_k(\mathbf{x}) = \prod_{j=1}^N \hat{f}_{k,j}(x_j)$
3. Estimate the prior probability $\hat{\pi}_k = \sum_{i=1}^M 1\{y_i = c_k\}/M$.
4. The Bayes rule is

$$h^*(\mathbf{x}) = \arg \max_k \hat{\pi}_k \hat{f}_k(\mathbf{x}). \quad (21)$$

3.2 Likelihood Density Estimation

There are various NB classifiers. They differ mainly by the data type of features and so the assumption they make regarding the distribution $f_{k,j}$.

- Gaussian NB: assumes that $f_{k,j}$ is a Gaussian.

$$f_{k,j}(x_j) = \frac{1}{\sqrt{2\pi\sigma_{k,j}^2}} \exp \left\{ -\frac{(x_j - \mu_{k,j})^2}{2\sigma_{k,j}^2} \right\}. \quad (22)$$

where $\mu_{k,j}$ and $\sigma_{k,j}$ are parameters of the Gaussian for feature j in class k , estimated as

$$\begin{aligned} \hat{\mu}_{k,j} &= \frac{\sum_{i=1}^M x_{i,j} 1\{y_i = c_k\}}{M_k} \\ \hat{\sigma}_{k,j} &= \frac{\sum_{i=1}^M (x_{i,j} 1\{y_i = c_k\} - \hat{\mu}_{k,j})^2}{M_k - 1}. \end{aligned}$$

- Bernoulli NB: implemented for features that are distributed according to multivariate Bernoulli distributions; i.e., there may be multiple features but each one is assumed to be a binary-valued variable.

$$f_{k,j}(x_j) = \mu_{k,j}^{x_j} (1 - \mu_{k,j})^{(1-x_j)} \quad (23)$$

where $\mu_{k,j}$ is the probability that attribute j appears in a data point pertaining to class k , estimated as

$$\hat{\mu}_{k,j} = \frac{\sum_{i=1}^M x_{i,j} 1\{y_i = c_k\}}{M_k}.$$

- Multinomial NB: implemented for multinomially distributed data. That is, there may be multiple features but each is assumed to be a binomial variable.

$$f_{k,j}(x_j) = \frac{M_k!}{x_j!(M_k - x_j)!} \mu_{k,j}^{x_j} (1 - \mu_{k,j})^{(M_k - x_j)}. \quad (24)$$

where $\mu_{k,j}$ is the probability that attribute j occurs in a sample pertaining to class k , estimated as

$$\hat{\mu}_{k,j} = \frac{\sum_{i=1}^M 1\{y_i = c_k\} x_{i,j} + \alpha}{\sum_{i=1}^M \sum_{j=1}^N 1\{y_i = c_k\} x_{i,j} + \alpha N}, \quad (25)$$

where α is a smoothing factor that accounts for features not present in the training dataset and prevents zero probabilities in inference.

- Categorical NB: implemented for categorically distributed data. It assumes that each feature has its own categorical distribution. Let $s_{j,l}$ denote class l of attribute j , then the density $f_{k,j}(x_j)$ is estimated as

$$\hat{f}_{k,j}(x_j = s_{j,l}) = \frac{\sum_{i=1}^M 1\{x_{i,j} = s_{j,l}, y_i = c_k\} + \alpha}{M_k + \alpha N_j} \quad (26)$$

where α is a smoothing factor that accounts for features not present in the training dataset and prevents zero probabilities in inference, and n_j is the number of categories of attribute j .

4 Logistic Regression

We studied logistic regression in the learning module “Regression”, which estimates the conditional probability $p_k = P(c_k|\mathbf{x})$ by fitting a parametric model to data. We will not repeat it in this learning module. Logistic regression and LDA are almost the same in that they both lead to a linear decision rule. In LDA, we estimate the joint probability $f(\mathbf{x}|y)f(y)$, which equals $f(y|\mathbf{x})f(\mathbf{x})$. In logistic regression, we estimate the conditional probability $f(y|\mathbf{x})$ only and ignores $f(\mathbf{x})$.

5 Support Vector Machine (SVM)

Let's consider a non-probabilistic classifier in the form of a decision boundary

$$f(\mathbf{x}) = \sum_{i=1}^M \lambda_i \mathcal{K}(\mathbf{x}, \mathbf{x}_i). \quad (27)$$

where $\mathcal{K}(\mathbf{x}, \mathbf{x}_i)$ is a kernel function measuring the similarity of a new data point, \mathbf{x} , to be tested with a training data point i , \mathbf{x}_i , in the feature space. By adding additional suitable constraints, many of the coefficients λ_i 's are zeros so that the prediction in test only depends on a subset of the training data points called *support vectors*. A resulting model is called support vector machine (SVM).

5.1 Large Margin Classifiers

Let's consider a linear decision boundary for a binary classification task:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (28)$$

where \mathbf{x} is the N -dimensional feature vector, \mathbf{w} is the N -dimensional coefficient vector, and w_0 is the bias. If one assume data points are in two classes $y \in \{1, -1\}$ and they are linearly separable, $f(\mathbf{x}) = 0$ split the feature space into two regions: $f(\mathbf{x}) > 0$ is the region for $y = 1$ and $f(\mathbf{x}) < 0$ is the region for $y = -1$. That is, $yf(\mathbf{x}) \geq 0$.

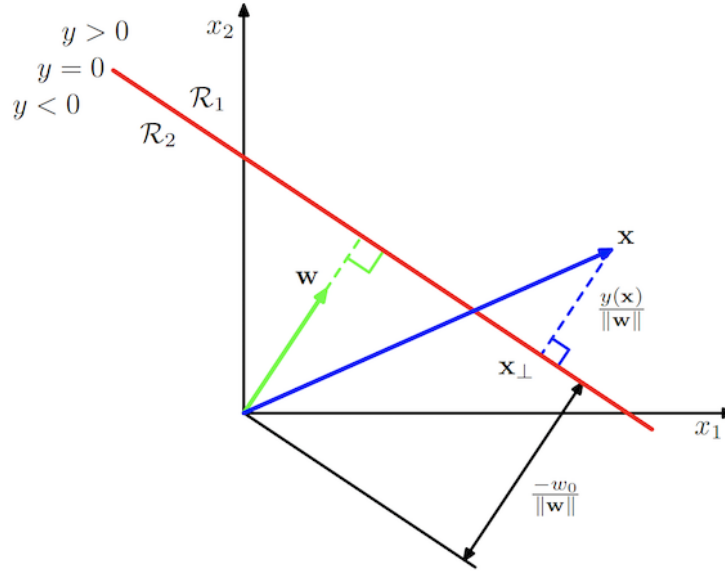


Figure 1: A linear decision boundary[4]

Ideally, we would like to pick a decision boundary that has maximum *margin*. The margin is the distance of the closet point to the decision boundary. Figure 1 illustrates a decision boundary on a 2D feature space. \mathbf{w} is perpendicular to the decision boundary¹, indicating that \mathbf{w} decides the direction of the boundary. w_0 is the distance from the boundary to the origin. A point in the feature space, \mathbf{x} , can be decomposed

$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \quad (29)$$

¹ $\mathbf{w}(\mathbf{x}_1 - \mathbf{x}_2) = 0$ for any two points, \mathbf{x}_1 and \mathbf{x}_2 , on the decision boundary. $\mathbf{x}_1 - \mathbf{x}_2$ is along the direction of $f(\mathbf{x}) = 0$. Therefore, \mathbf{w} must be perpendicular to $f(\mathbf{x}) = 0$.

where \mathbf{x}_\perp is the orthogonal projection of \mathbf{x} onto the boundary. $r \frac{\mathbf{w}}{\|\mathbf{w}\|}$ is the projection of \mathbf{x} along the direction of \mathbf{w} , where r is the length of \mathbf{x} along \mathbf{w} . Since

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + w_0 \\ &= \mathbf{w}^T \left(\mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + w_0 \\ &= \mathbf{w}^T \mathbf{x}_\perp + r \|\mathbf{w}\| + w_0, \end{aligned} \quad (30)$$

and $\mathbf{w}^T \mathbf{x}_\perp = 0$, we have

$$r = \frac{f(\mathbf{x})}{\|\mathbf{w}\|} + \frac{-w_0}{\|\mathbf{w}\|}. \quad (31)$$

Figure 1 indicates the distance from \mathbf{x} to the decision boundary is $\frac{f(\mathbf{x})}{\|\mathbf{w}\|}$. Therefore, to maximize the length of this distance is about solving the following optimization problem:

$$\max_{\mathbf{w}, w_0} \frac{1}{\|\mathbf{w}\|} \min_{i \in \{1, \dots, M\}} (y_i (\mathbf{w}^T \mathbf{x}_i + w_0)). \quad (32)$$

Scaling \mathbf{w} and w_0 in (32) does not change the value of the objective function. Therefore, one can force

$$\min_{i \in \{1, \dots, M\}} (y_i (\mathbf{w}^T \mathbf{x}_i + w_0)) = 1 \quad (33)$$

so that (32) becomes a quadratic programming (QP) problem:

$$\begin{aligned} &\min_{\mathbf{w}, w_0} \|\mathbf{w}\|^2 \\ &\text{s.t.} \\ &y_i (\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1, i = 1, \dots, M. \end{aligned} \quad (34)$$

To solve the QP problem in (34), one can construct a Lagrangian by introducing Lagrangian multipliers $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_M]^T \geq 0$,

$$L(\mathbf{w}, w_0, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^M \lambda_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1), \quad (35)$$

which is the lower boundary for the objective function in (34). We can maximize the Lagrangian function with respect to decision variables $\mathbf{w}, w_0, \boldsymbol{\lambda}$ to find the optimal solution to (34). First, the \mathbf{w} and w_0 are optimized out by solving the following linear system:

$$\nabla_{\mathbf{w}} L(\mathbf{w}, w_0, \boldsymbol{\lambda}) = \mathbf{w} - \sum_{i=1}^M \lambda_i y_i \mathbf{x}_i = 0, \quad (36)$$

$$\frac{\partial L(\mathbf{w}, w_0, \boldsymbol{\lambda})}{\partial w_0} = \sum_{i=1}^M \lambda_i y_i = 0. \quad (37)$$

Therefore,

$$\hat{\mathbf{w}} = \sum_{i=1}^M \lambda_i y_i \mathbf{x}_i, \quad (38)$$

$$\sum_{i=1}^M \lambda_i y_i = 0. \quad (39)$$

Bringing them into the Lagrangian in (35) leads to the Lagrangian dual function:

$$L(\boldsymbol{\lambda}) = -\frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^M \lambda_i. \quad (40)$$

Ultimately, we will solve the following standard QP problem:

$$\begin{aligned}
& \max_{\boldsymbol{\lambda}} -\frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^M \lambda_i \\
& \text{s.t.} \\
& \sum_{i=1}^M \lambda_i y_i = 0 \\
& \lambda_i \geq 0, \quad i = 1, \dots, M.
\end{aligned} \tag{41}$$

Solution to (41), denoted by $\hat{\boldsymbol{\lambda}}$, must also satisfy the Karush–Kuhn–Tucker (KKT) conditions:

$$\begin{aligned}
& \lambda_i \geq 0 \\
& y_i f(\mathbf{x}_i) - 1 \geq 0 \\
& \lambda_i (y_i f(\mathbf{x}_i) - 1) = 0
\end{aligned} \tag{42}$$

for $i = 1, \dots, M$. The KKT conditions in (42) says that either $\hat{\lambda}_i = 0$ or the constraint $y_i f(\mathbf{x}_i) = 1$ is active. $y_i f(\mathbf{x}_i) = 1$ means the data point i is on the maximum margin of SVM. Those data points are support vectors. \mathcal{S} denotes the set of support vectors.

Finally, the SVM classifier is

$$\begin{aligned}
f^*(\mathbf{x}) &= f(\mathbf{x} | \hat{\mathbf{w}}, \hat{w}_0) \\
&= \hat{\mathbf{w}}^T \mathbf{x} + \hat{w}_0 \\
&= \sum_{i \in \mathcal{S}} \hat{\lambda}_i y_i \mathbf{x}_i^T \mathbf{x} + \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (y_i - \hat{\lambda}_i y_i \mathbf{x}_i^T \mathbf{x}_i),
\end{aligned} \tag{43}$$

which shows that only a portion of the training data points (i.e., support vectors) are used in predicting the class label given the feature vector \mathbf{x} .

5.2 Soft Margin Classifiers

If data are not linearly separable, we introduce a slack variable ξ_i (≥ 0) to replace the hard constraint $y_i f(\mathbf{x}_i) \geq 1$ with the soft constraint $y_i f(\mathbf{x}_i) \geq 1 - \xi_i$. Consequently, the optimization problem in (34) becomes

$$\begin{aligned}
& \min_{\mathbf{w}, w_0, \boldsymbol{\xi}} \|\mathbf{w}\|^2 + C \sum_{i=1}^M \xi_i \\
& \text{s.t.} \\
& y_i (\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 - \xi_i, \quad i = 1, \dots, M \\
& \xi_i \geq 0, \quad i = 1, \dots, M
\end{aligned} \tag{44}$$

where $C (\geq 0)$ is a hyperparameter that controls how many points are allowed to violate the margin constraint. If $C \rightarrow \infty$, it becomes the hard margin classifier.

The corresponding Lagrangian is

$$L(\mathbf{w}, w_0, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\eta}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^M \xi_i - \sum_{i=1}^M \lambda_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i) - \sum_{i=1}^M \eta_i \xi_i \tag{45}$$

where $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_M]^T \geq 0$ and $\boldsymbol{\eta} = [\eta_1, \dots, \eta_M]^T \geq 0$ are Lagrangian multipliers. By solving the linear system: $\nabla_{\mathbf{w}} L = 0$, $\frac{\partial L}{\partial w_0} = 0$, and $\nabla_{\boldsymbol{\xi}} L = 0$, \mathbf{w} , w_0 , and $\boldsymbol{\xi}$ in the Lagrangian are optimized out, leading to the

same Lagrangian dual function $L(\boldsymbol{\lambda})$ as in (40). Finally, we solve the following QP dual problem:

$$\begin{aligned}
& \max_{\boldsymbol{\lambda}, \boldsymbol{\eta}} -\frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^M \lambda_i \\
& \text{s.t.} \\
& \sum_{i=1}^M \lambda_i y_i = 0 \\
& C - \lambda_i - \eta_i = 0, \quad i = 1, \dots, M \\
& \lambda_i, \eta_i \geq 0, \quad i = 1, \dots, M
\end{aligned} \tag{46}$$

While the objective function in the soft margin case is identical to that in the hard margin case, the constraint on λ_i is different.

The KKT conditions for the soft margin case are:

$$\begin{aligned}
& \lambda_i \geq 0 \\
& y_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i \geq 0 \\
& \lambda_i(y_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i) = 0 \\
& \eta_i \geq 0 \\
& \xi_i \geq 0 \\
& \eta_i \xi_i = 0
\end{aligned} \tag{47}$$

for all data point in the training set (i.e., $i = 1, \dots, M$).

Moreover,

$$C - \lambda_i - \eta_i = 0 \tag{48}$$

The value of λ_i tells the position of the corresponding training data point in the feature space.

- $\lambda_i = 0$: the data point is outside the margin
On one hand, $\lambda_i = 0$ indicates $y_i f(\mathbf{x}_i) > 1 - \xi_i$. On the other hand, $\lambda_i = 0$ indicates $\eta_i = C > 0$ and so ξ_i must be zero. Therefore, the data point is outside the margin and it can be ignored.
- $0 < \lambda_i < C$: the data point is on the margin
This indicates that $y_i f(\mathbf{x}_i) - 1 = \xi_i$. And, $\lambda_i < C$ indicates that $\eta_i > 0$ and so $\xi_i = 0$. Therefore, the data point is on the margin.
- $\lambda_i = C$: the data point can be inside the margin
This indicates that $y_i f(\mathbf{x}_i) - 1 = \xi_i$. Moreover, $\lambda_i = C$ indicates $\eta_i = 0$ and so $\xi_i > 0$. Therefore, the data point is allowed to be inside the margin.

6 Classification and Regression Trees (CART)

CART, also called decision trees, are defined by recursively partitioning the feature space, and defining a local model in each resulting region of feature space. The overall model can be represented by a tree, with one leaf per region.

6.1 Definition of a Classification Tree

Let's focus on the classification trees in this learning module. A tree consists of a set of nested decision rules. At node l of the tree, $\mathcal{D}_l = \{(\mathbf{x}_i, y_i) \in \mathcal{T}\}$ is the subset of the data \mathcal{T} reaching that node. One can choose a feature to split \mathcal{D}_l and pass data points to the left and right branches, depending on the relationship between the feature value and a selected threshold value t_l . The tree ends with Q leaves: $\mathcal{V}_1, \dots, \mathcal{V}_Q$. Each leaf contains a distribution of data points over class labels:

$$p_{q,k} = P(c_k | \mathcal{V}_q), \tag{49}$$

Denote the predicted class of leaf \mathcal{V}_q as \hat{y}_q . One can determine \hat{y}_q according to $p_{q,k}$.

Leaves of a tree are a result of the nested decision rules. Formally, a decision tree is

$$f(\mathbf{x}, \boldsymbol{\theta}) \quad (50)$$

where $\boldsymbol{\theta}$ denotes the nested thresholds chosen for partitioning the data.

6.2 Optimization of a Classification Tree

The loss function is the sum of losses attained on the entire training dataset:

$$L(\boldsymbol{\theta}) = \sum_{i=1}^M l(y_i, f(\mathbf{x}_i, \boldsymbol{\theta})) = \sum_{q=1}^Q \sum_{i: \mathbf{x}_i \in \mathcal{L}_q} l(y_i, \hat{y}_i). \quad (51)$$

One can optimize the decision tree by minimizing the loss function. The loss function is not differential. A greedy procedure is commonly used to iteratively grow the tree one node at a time.

Suppose we are at node l . If one chose feature j for partitioning the data \mathcal{D}_l , and this feature is a real-valued feature, a threshold value t is chosen for passing the data to the left and right branches, respectively:

$$\mathcal{D}_l^L(j, t) = \{(\mathbf{x}_i, y_i) \in \mathcal{D}_l | x_{i,j} \leq t\}, \quad (52)$$

$$\mathcal{D}_l^R(j, t) = \{(\mathbf{x}_i, y_i) \in \mathcal{D}_l | x_{i,j} > t\}, \quad (53)$$

where the threshold t can be chosen from the set of unique values of feature j , denoted by \mathcal{X}_j .

If the feature chosen for partitioning data at node l is a categorical feature whose space is \mathcal{X}_j . A class $t \in \mathcal{X}_j$ is chosen for passing the data to the left and right branches, respectively:

$$\mathcal{D}_l^L(j, t) = \{(\mathbf{x}_i, y_i) \in \mathcal{D}_l | x_{i,j} = t\}, \quad (54)$$

$$\mathcal{D}_l^R(j, t) = \{(\mathbf{x}_i, y_i) \in \mathcal{D}_l | x_{i,j} \neq t\}. \quad (55)$$

One optimizes the tree by choosing the best feature to slit on, and the best threshold for that feature:

$$(j_l, t_l) = \arg \min_{j \in \{1, \dots, N\}, t \in \mathcal{X}_j} \frac{|\mathcal{D}_l^L|}{|\mathcal{D}_l|} L(\mathcal{D}_l^L) + \frac{|\mathcal{D}_l^R|}{L\mathcal{D}_l|} l(\mathcal{D}_l^R). \quad (56)$$

where $L(\cdot)$ is the loss at a branch.

To calculate the loss at branch $s \in \{L, R\}$, one can first calculate the empirical distribution of \mathcal{D}_l^s over class labels:

$$p_{l,k}^s = \frac{\sum_{y_i \in \mathcal{D}_l^s} 1\{y_i = c_k\}}{|\mathcal{D}_l^s|}. \quad (57)$$

Then, the expected error rate, called Gini Index, is calculated for each branch:

$$G_l^s = \sum_{k=1}^K p_{l,k}^s (1 - p_{l,k}^s) = 1 - \sum_{k=1}^K (p_{l,k}^s)^2, \quad (58)$$

which can be used as the loss at that branch. Alternatively, one can use the entropy as the loss for that branch:

$$l(\mathcal{D}_l^s) = - \sum_{k=1}^K p_{l,k}^s \log p_{l,k}^s. \quad (59)$$

A node is *pure* if the loss is zero, meaning that all data points are of one class.

6.3 Regularization

If one let the tree become deep enough, it can achieve 0 error on the training set, by partitioning the input space into sufficiently small regions where the output is constant. However, this will typically result in overfitting. To prevent this, there are two main approaches.

- stop the tree growing process according to some heuristic, such as having too few examples at a node, or reaching a maximum depth.
- grow the tree to its maximum depth, where no more splits are possible, and then to prune it back, by merging split subtrees back into their parent

This can partially overcome the greedy nature of top-down tree growing.

7 Ensemble Learning

Classification trees are a high variance estimator. The prediction result could vary a lot if the training data is perturbed. *Ensemble learning* is an approach to reducing the variance. Ensemble learning ensembles the outputs from multiple estimators to reach the final prediction. The assemble methods could be averaging, majority voting, relearning from base model's outputs, or others. In this learning module, we introduce a few ensemble learning methods, including stacking, bagging, random forests, and boosting.

7.1 Stacking

Stacking involves training a set of base models and then using their outputs to train a meta model that gives the final prediction. Base models are usually complex and quite different in their approaches to solving the classification problem. It's preferred that they have different strengths. Candidates for base models include decision trees, support vector machines, neural networks, and more. The meta model usually is relatively simple, such as the logistic regression.

7.2 Bagging

Bagging involves training different base models on different randomly sampled versions of the data. The training data for a model is sampled with replacement (called bootstrap) from the original dataset so that a data point can appear multiple times in the training data. The training data for each model are in the same size as the original dataset. However, the training dataset for a model only contains about 63% of the original data. The remaining 37% out-of-sample (oos) can be used as the test data for that model.

The advantage of bootstrap is to prevent the ensemble from relying too much on any individual data points. Therefore, bagging is useful if the base models are unstable estimators. Bagging is simple. It retrains the same learning algorithm on different subsets of the data, and hopefully obtain diverse base models.

7.3 Random Forests

Random forests method involves learning different trees on various subsets of training data and use randomly chosen subset of features at each node of trees.

Ensembles of trees, fit by random forests or bagging, are in the form of

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{v=1}^V \beta_v F_v(\mathbf{x}; \boldsymbol{\theta}_v) \quad (60)$$

where F_v is the v th model.

7.4 Boosting

Boosting involves training base models in a sequential manner. The first base model F_1 is trained on the entire training dataset. In fitting the second base model F_2 , the training data is weighted according to the

loss on F_1 . This process continues until M (a hyper-parameter of boosting algorithms) base models have been developed. As long as F_v is more accurate than chance, the final ensemble has a higher accuracy than base models. Boosting is a process that turns a weaker learner (i.e., the accuracy is slightly better than chance) into a stronger learner.

The forward stepwise additive model for boosting is

$$\boldsymbol{\theta}_v, \beta_v = \arg \min_{\boldsymbol{\theta}, \beta} \sum_{i=1}^n l(y_i, f_{v-1}(\mathbf{x}_i) + \beta F(\mathbf{x}_i, \boldsymbol{\theta})) \quad (61)$$

where

- y_i is the true label,
- $f_{v-1}(\mathbf{x}_i)$ is the prediction from the previous ensemble of models,
- $F(\mathbf{x}_i, \boldsymbol{\theta})$ is the new base model added in step v ,
- β is the scaling factor for the new model,
- l is the loss function.

The model is updated in each step by adding the new model F_v to the current ensemble prediction

$$f_v(\mathbf{x}) = f_{v-1}(\mathbf{x}) + \beta_v F(\mathbf{x}; \boldsymbol{\theta}_v) = f_{v-1}(\mathbf{x}) + \beta_v F_v(\mathbf{x}). \quad (62)$$

There are different versions of boosting algorithms, such as AdaBoosting, Gradient Boosting, Extreme Gradient Boosting (XGBoost). We encourage students to read Section 18.5 in [4].

8 Multiple Binary Classification Tasks

Some classifiers are inherently binary classifiers. However, they can be used for multiclass classification tasks too because one can decompose a multiclass classification task into a set of binary classification tasks.

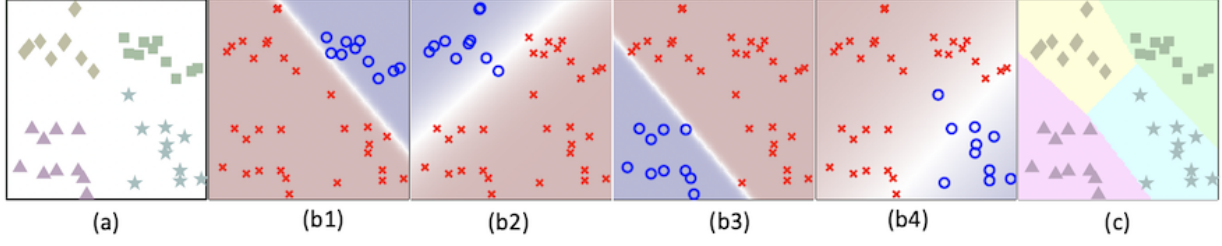
8.1 One-Versus-All Classifier

A straightforward reduction of the multiclass classification to binary classification is based on the so-called one-versus-all (OVA) or one-versus-the-rest (OVR) technique. Figure 2 illustrates the concept of this method. If the response has K classes, K binary classifiers are trained and each is for differentiating one class from the rest. The final classifier integrates the K classifiers together. The feature vector of a new data point, \mathbf{x} , is evaluated by all of the K classifiers, which respectively gives the probability that data point belongs to class k for $k = 1, \dots, K$. The predicted class for \mathbf{x} is the one associated with the top-scorer.

The one-vs-all classifier is simple because the number of binary classifiers required is equal to the number of classes. But it has a limitation of unbalance. That is, usually the “rest” class has a larger size than the “one” class.

8.2 One-Versus-One Classifier

An alternative approach is to train a binary classifier for each pair of classes. If the response variable Y has K classes, $\binom{K}{2}$ binary classifiers are trained. Figure 3 below illustrates an example. The response Y has four classes. Therefore, $\binom{4}{2} = 6$ binary classifiers are trained. The final classifier is a tournament champion of the six classifiers.



(a). A training dataset: $\{(x_i, y_i) | i = 1, \dots, m\}$, $y_i \in \{\blacksquare, \blacklozenge, \blacktriangle, \star\}$

(b1). Binary classifier: $h_{\blacksquare} = P(x_i; \beta_{\blacksquare})$ trained on $\{(x_i, y_i^{\blacksquare}) | i = 1, \dots, m\}$, where $y_i^{\blacksquare} = \begin{cases} 1, & \text{if } y_i = \blacksquare \\ 0, & \text{if } y_i = \blacklozenge, \blacktriangle, \star \end{cases}$

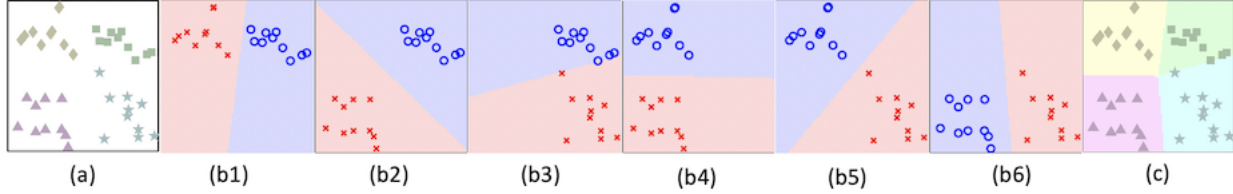
(b2). A binary classifier: $h_{\blacklozenge} = P(x_i; \beta_{\blacklozenge})$ trained on $\{(x_i, y_i^{\blacklozenge}) | i = 1, \dots, m\}$, where $y_i^{\blacklozenge} = \begin{cases} 1, & \text{if } y_i = \blacklozenge \\ 0, & \text{if } y_i = \blacksquare, \blacktriangle, \star \end{cases}$

(b3). A binary classifier: $h_{\blacktriangle} = P(x_i; \beta_{\blacktriangle})$ trained on $\{(x_i, y_i^{\blacktriangle}) | i = 1, \dots, m\}$, where $y_i^{\blacktriangle} = \begin{cases} 1, & \text{if } y_i = \blacktriangle \\ 0, & \text{if } y_i = \blacksquare, \blacklozenge, \star \end{cases}$

(b4). A binary classifier: $h_{\star} = P(x_i; \beta_{\star})$ trained on $\{(x_i, y_i^{\star}) | i = 1, \dots, m\}$, where $y_i^{\star} = \begin{cases} 1, & \text{if } y_i = \star \\ 0, & \text{if } y_i = \blacksquare, \blacklozenge, \blacktriangle \end{cases}$

(c). The one-vs-all classifier: $g = \underset{k \in \{\blacksquare, \blacklozenge, \blacktriangle, \star\}}{\operatorname{argmax}} \{P(x_i; \beta_k)\}$

Figure 2: The concept of the one-versus-all classifier



(a). A training dataset: $\{(x_i, y_i) | i = 1, \dots, m\}$, $y_i \in \{\blacksquare, \blacklozenge, \blacktriangle, \star\}$

(b1). Binary classifier: $h_1 = P(x_i; \beta_{\blacksquare, \blacklozenge})$ trained on the subset $\{(x_i, y_i) | y_i = \blacksquare \text{ or } \blacklozenge\}$

(b2). Binary classifier: $h_2 = P(x_i; \beta_{\blacksquare, \blacktriangle})$ trained on the subset $\{(x_i, y_i) | y_i = \blacksquare \text{ or } \blacktriangle\}$

(b3). Binary classifier: $h_3 = P(x_i; \beta_{\blacksquare, \star})$ trained on the subset $\{(x_i, y_i) | y_i = \blacksquare \text{ or } \star\}$

(b4). Binary classifier: $h_4 = P(x_i; \beta_{\blacklozenge, \blacktriangle})$ trained on the subset $\{(x_i, y_i) | y_i = \blacklozenge \text{ or } \blacktriangle\}$

(b5). Binary classifier: $h_5 = P(x_i; \beta_{\blacklozenge, \star})$ trained on the subset $\{(x_i, y_i) | y_i = \blacklozenge \text{ or } \star\}$

(b6). Binary classifier: $h_6 = P(x_i; \beta_{\blacktriangle, \star})$ trained on the subset $\{(x_i, y_i) | y_i = \blacktriangle \text{ or } \star\}$

(c). The one-vs-one classifier: $g = \underset{k, l \in \{\blacksquare, \blacklozenge, \blacktriangle, \star\}}{\operatorname{argmax}} \{P(x_i; \beta_{k, l})\}$

Figure 3: The concept of the one-versus-one classifier

9 Classification Performance Metrics

9.1 Binary Classification Results

Assume that we have a test dataset to evaluate the performance of a classification model. There are four scenarios of classification result:

- *True Positive (TP)*: A true positive instance is correctly predicted as positive;
- *True Negative (TN)*: A true negative instance is correctly predicted as negative;
- *False Positive (FP)*: A true negative instance is mistakenly predicted as positive;
- *False Negative (FN)*: A true positive instance is mistakenly predicted as negative;

9.2 Confusion Matrix

The confusion matrix in Table 1 counts the numbers of instances according to their true label and predicted label.

Table 1: Confusion Matrix

		predicted label (\hat{y})		
		1	...	K
true label (y)	1			
	\vdots			
	K			

9.3 Error Rates

Overall Error Rate is the percentage of misclassified observations:

$$\text{Overall Error Rate} = \frac{\sum_{i=1}^M 1\{y_i \neq \hat{y}_i\}}{M}. \quad (63)$$

Overall error rate conveys an aggregate measure of misclassification, it counts a false positive the same as a false negative. To account for the asymmetric costs in misclassification, we define error rate with respect to the individual classes.

Class k error rate is the percent of class k instances that are misclassified:

$$\text{error rate}_k = \frac{\sum_{i=1}^M 1\{y_i = k\} 1\{\hat{y}_i \neq k\}}{\sum_{i=1}^M 1\{y_i = k\}}. \quad (64)$$

9.4 Correct Rates

Accuracy is the percentage of correctly classified observations:

$$\text{Accuracy} = \frac{\sum_{i=1}^M 1\{y_i = \hat{y}_i\}}{\sum_{i=1}^M 1\{y_i\}}. \quad (65)$$

Accuracy is 1-overall error rate.

Recall is the ability to correctly predict class k instances, and it is calculated as:

$$\text{Recall}_k = \frac{\sum_{i=1}^M 1\{\hat{y}_i = y_i = k\}}{\sum_{i=1}^M 1\{y_i = k\}}. \quad (66)$$

Precision is the proportion of instances predicted to be class k that are actually in class k :

$$\text{Precision}_k = \frac{\sum_{i=1}^M 1\{\hat{y}_i = y_i = k\}}{\sum_{i=1}^M 1\{\hat{y}_i = k\}}. \quad (67)$$

F1 Score combines precision and recall into a single measure and is defined as:

$$\text{F1}_k = \frac{2}{\frac{1}{\text{Recall}_k} + \frac{1}{\text{Precision}_k}}. \quad (68)$$

9.5 Precision-Recall Curve and Average Precision (AP)

Precision and recall for each class are dependent of the threshold for translating the probabilistic prediction into class prediction. By changing the threshold value, we can get different pairs of precision and recall. Accordingly, a *precision-recall curve* can be created for displaying the trade-off between them, where x -axis is the recall and the y -axis is the precision. The *Average Precision* (AP_k) is the area under that curve.

For binary classification, the precision-recall curve is created for the positive class (i.e., class 1). For multiclass classification, this curve is created for every class.

9.6 Receiver Operating Characteristic (ROC) and Area Under ROC (AUC)

Receiver Operating Characteristic (ROC) curve is a graphical approach for displaying the tradeoff between a classifier's ability to correctly identify class k instances and the error rate for misclassifying instances of other classes as class k instances. The x -axis is the false positive rate for class k and the y -axis is Recall $_k$. The curve is obtained by changing the threshold value for classification. The area under the ROC curve is called AUC_k and measures the quality of the classifier in classifying class k . The larger the AUC_k , the better the classifier performs.

For binary classification, the ROC curve is created for the positive class (i.e., class 1), with x -axis being the class 0 error (i.e., false positive rate) and y -axis being the recall for the positive class (i.e., the true positive rate). For multiclass classification, this curve is created for every class.

9.7 Macro-level Metrics

A macro-level average is an unweighted average of class-level metric values:

$$\text{mRecall} = \frac{1}{K} \sum_{k=1}^K \text{Recall}_k \quad (69)$$

$$\text{mPrecision} = \frac{1}{K} \sum_{k=1}^K \text{Precision}_k \quad (70)$$

$$\text{mF1} = \frac{1}{K} \sum_{k=1}^K \text{F1}_k \quad (71)$$

$$\text{mAP} = \frac{1}{K} \sum_{k=1}^K \text{AP}_k \quad (72)$$

$$\text{mAUC} = \frac{1}{K} \sum_{k=1}^K \text{AUC}_k \quad (73)$$

9.8 Weighed Metrics

A weighted average is about calculating the class-level metric values and then averaging them with the weights determined according to class support (i.e., the number of true instances for each class):

$$\text{mRecall} = \sum_{k=1}^K \frac{\sum_{i=1}^M 1\{y_i = k\}}{M} \text{Recall}_k \quad (74)$$

$$\text{mPrecision} = \sum_{k=1}^K \frac{\sum_{i=1}^M 1\{y_i = k\}}{M} \text{Precision}_k \quad (75)$$

$$\text{mF1} = \sum_{k=1}^K \frac{\sum_{i=1}^M 1\{y_i = k\}}{M} \text{F1}_k \quad (76)$$

$$\text{mAP} = \sum_{k=1}^K \frac{\sum_{i=1}^M 1\{y_i = k\}}{M} \text{AP}_k \quad (77)$$

$$\text{mAUC} = \sum_{k=1}^K \frac{\sum_{i=1}^M 1\{y_i = k\}}{M} \text{AUC}_k \quad (78)$$

9.9 Micro Average

To find a micro average, we calculate metrics globally. The accuracy in (65) and the overall error rate in (63) are both micro averages.

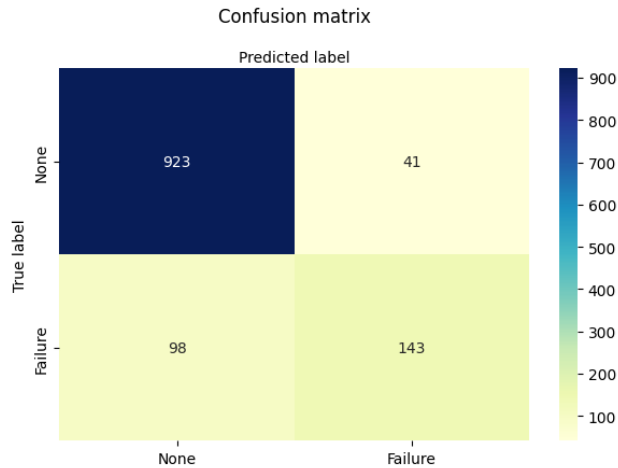
10 Examples

10.1 An Example of Binary Classification

Let's take a look at the binary classification problem in the Python Notebook. Results below is obtained using an LDA classifier.

Table 2 is the confusion matrix given the cutoff point is 0.5.

Table 2: Confusion matrix for the concrete defect class prediction



Below is the report of performance metrics:


```

1
2 Class-level report:
3     precision    recall  f1-score   support
4
5     None         0.90      0.96      0.93       964
6     Failure       0.78      0.59      0.67       241
7
8     accuracy          0.88          0.88          1205
9     macro avg         0.84          0.78          1205
10    weighted avg      0.88          0.88          1205
11
12 Task-level report:
13 ACC      0.885
14 AUC      0.897
15 F1       0.673
16 dtype: float64

```

Figure 4 shows the precision-recall curve. The average precision is 0.761.

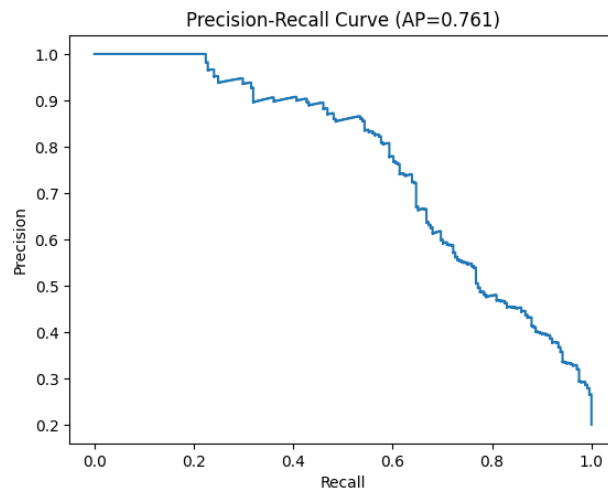


Figure 4: Precision Recall Curve

The orange colored curve in Figure 5 shows the ROC curve. The AUC is 0.90.

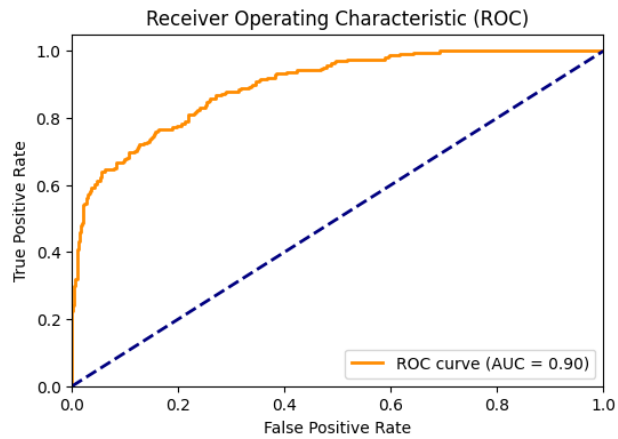


Figure 5: ROC and AUC

10.2 An Example of Multiclass Classification

Let's take a look at the multiclass classification problem in the Python Notebook. The following is the result using a QDA classifier.

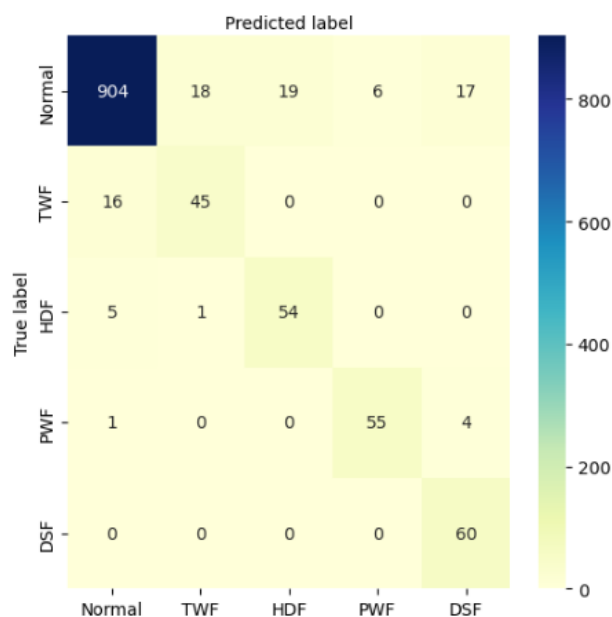
```

1
2 Class-level report:
3           precision    recall  f1-score   support
4
5   Normal      0.98      0.94      0.96      964
6   TWF         0.70      0.74      0.72       61
7   HDF         0.74      0.90      0.81       60
8   PWF         0.90      0.92      0.91       60
9   DSF         0.74      1.00      0.85       60
10
11   accuracy          0.93      1205
12   macro avg         0.81      0.90      0.85      1205
13   weighted avg      0.94      0.93      0.93      1205
14
15 Task-level report:
16 ACC      0.928
17 AUC      0.983
18 F1       0.930
19 dtype: float64

```

The confusion matrix is in Figure 6 .

Figure 6: Confusion matrix for the testing result on MNIST dataset



References

- [1] Yaser S Abu-Mostafa, Malik Magdon-Ismail, and Hsuan-Tien Lin. *Learning From Data*, volume 4. AMLBook New York, 2012.
- [2] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [3] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.

- [4] Kevin P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022.
- [5] Larry Wasserman. *All of statistics: a concise course in statistical inference*, volume 26. Springer, 2004.