# PNEUMONIA DETECTION USING NON-LINEAR SVM

# SUPERVISED CLASSIFICATION OF X-RAY CHEST IMAGES

**01** SUPPORT VECTOR MACHINE

Linear and nonlinear

**02** DATASET

X ray chest images

**03** DATA PROCESSING

Image resize
Image Enhacement
Image Denoising

**04** PERFORMANCE MEASURES

Accuracy - Precision
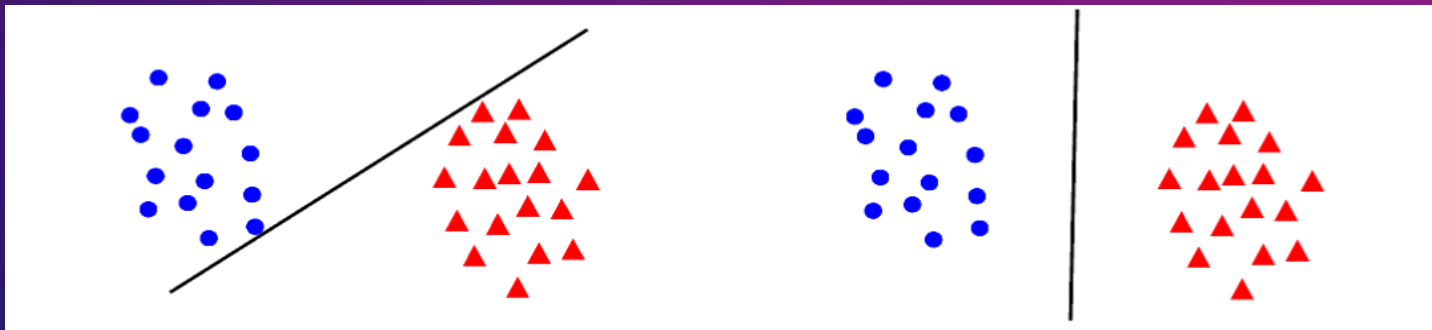Recall -F1- Score

# 01

# SUPPORT VECTOR MACHINE

*Linear and nonlinear*

# SVM

Let's assume that A and B are linearly separable sets, meaning there exists a hyperplane H = { x ∈ R^n : w^T x + b = 0 } such that

$$\begin{cases} w^T x + b > 0. & y_i = +1 \\ w^T x + b < 0. & y_i = -1 \end{cases} \longrightarrow y_i(w^T x + b) > 0$$
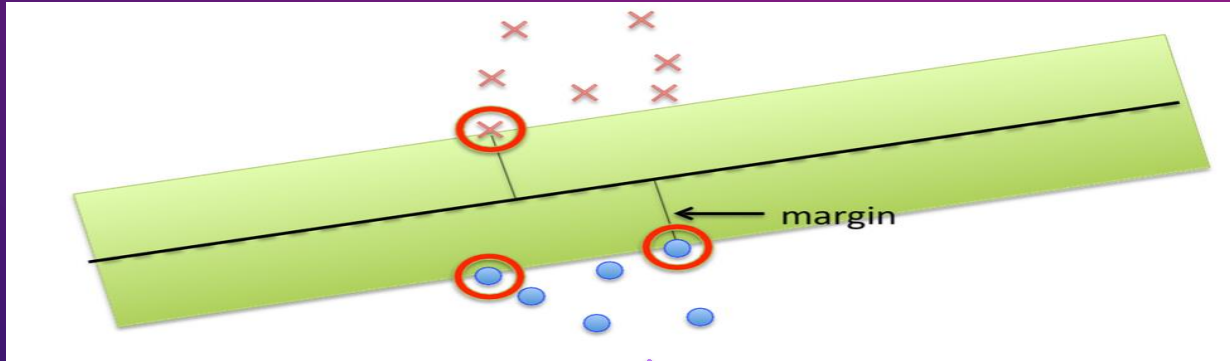
The separating hyperplane is not unique.

# SVM GOAL

Find the separating hyperplane that maximizes the distance between the closest training set and the separating hyperplane.



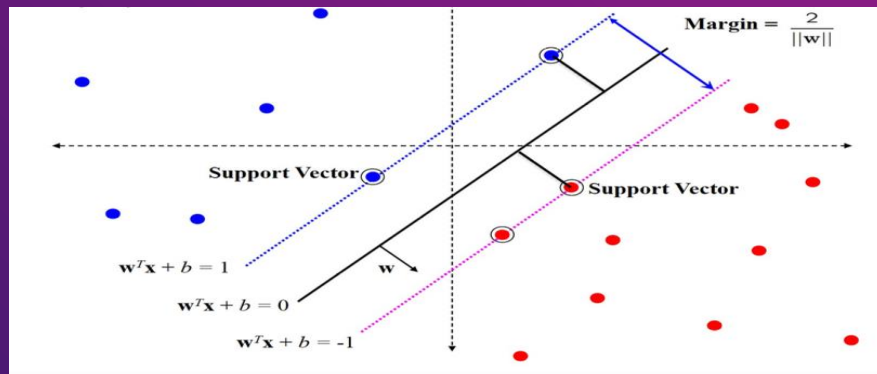The distance from the closest training point is called the margin.

The circled points are called support vectors.

# GEOMETRIC MARGIN

We define the optimal hyperplane (or maximum margin hyperplane) as:
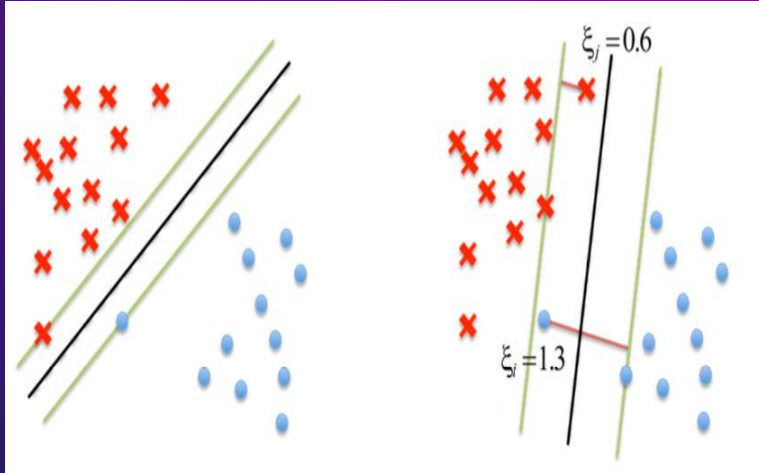$$(w*, b*) = argmax \ \rho(w, b)$$



The geometric margin between the two classes is $\frac{2}{\|w\|}$

The larger the margin (separation), the higher the expected generalization

# NON LINEARLY SEPARABLE SETS

Linear separability is generally too strong an assumption. However, the concept of an optimal separating hyperplane still makes sense, based on slack variables to handle outliers.



The ξi (Slack Variables) account for the non-separability of the data.

$ξi = 0 ⇒ x\_i$ correctly classified (outside the margin)

$ξi ∈ (0, 1) ⇒ x\_i$ correctly classified (but within the margin)

$ξi > 1 ⇒ x\_i$ misclassified (on the wrong side of the separating hyperplane))

# C-SVM

Maximization of the (Soft) Margin and Minimization of the Number of Misclassified Samples

$$\min \frac{1}{2} \parallel w \parallel_2^2 \; + c \sum_{i=1}^{n} \xi_i$$

$$\text{s.t. } y_i(w^T x + b) \geq \; 1 - \xi_i$$
$$\xi_i \geq 0$$
$$\text{Where C} > 0$$

The regularization parameter C takes into account the penalty for misclassified data.

Larger C⇒ Fewer exceptions (smaller margin, potential overfitting).

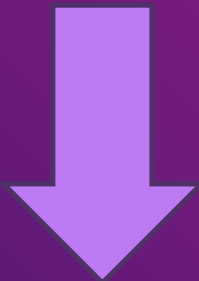Smaller C⇒ more eccezioni exceptions (larger margin, potential underfitting).

# Non linear SVM

A classification problem of complex patterns thrown into a high-dimensional nonlinear space is more likely to be linearly separable than in a low-dimensional space. (Thomas Cover, 1965)
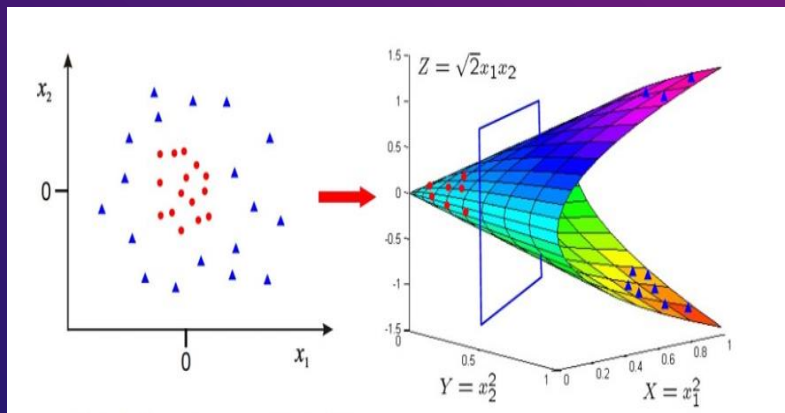The goal is to project into high-dimensional space and solve with a linear model

Mapping the data (input space) into a feature space of higher dimensions using a non-linear transformation $\Phi(x) \in R^{\wedge}p$ (p > n)

# KERNEL

*Applying SVM to $\Phi(x\_i)$ rather than to $x\_i$, such that the formulation of SVM (in the feature space) is expressed only using the inner product $\langle\Phi(x\_i), \Phi(x\_j)\rangle$. To achieve this, it is sufficient to know how to compute the scalar product $k(x\_i, x\_j)$ in the feature space.*
*k is called a kernel, and the same kernel can correspond to different feature maps $\Phi$.*



*The idea of the kernel function is to perform operations in the input space rather than in the potentially high-dimensional feature space. Therefore, the inner product does not need to be evaluated in the feature space. We want the function to map observations from the input space to the feature space.*

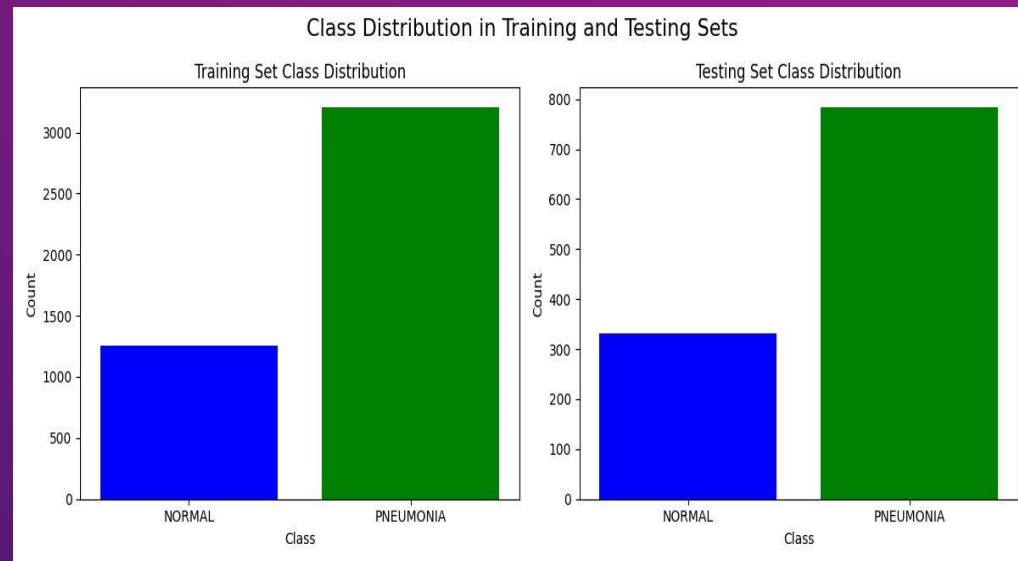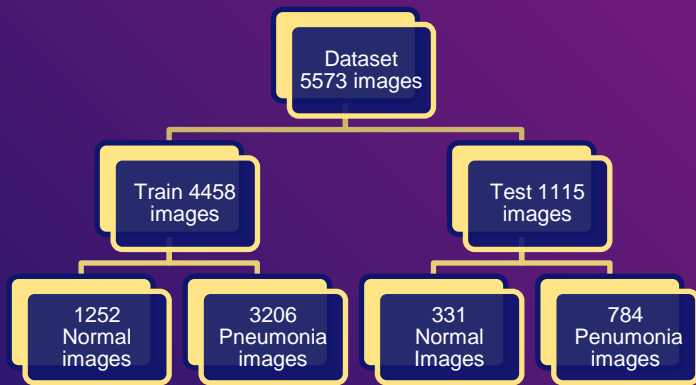# DATASET

**02**

*X ray chest images*

# DATASET



*The dataset is organized into 2 different types (train and test) and contains various categories of images (Pneumonia/Normal). There are 5,573 chest X-ray images (JPEG) and 2 categories (Pneumonia/Normal).*

*For the analysis of chest X-ray images, all chest X-rays initially underwent a quality check, eliminating all low-quality or unreadable scans. The diagnoses of the images were then classified by two expert physicians before being authorized for model training.*
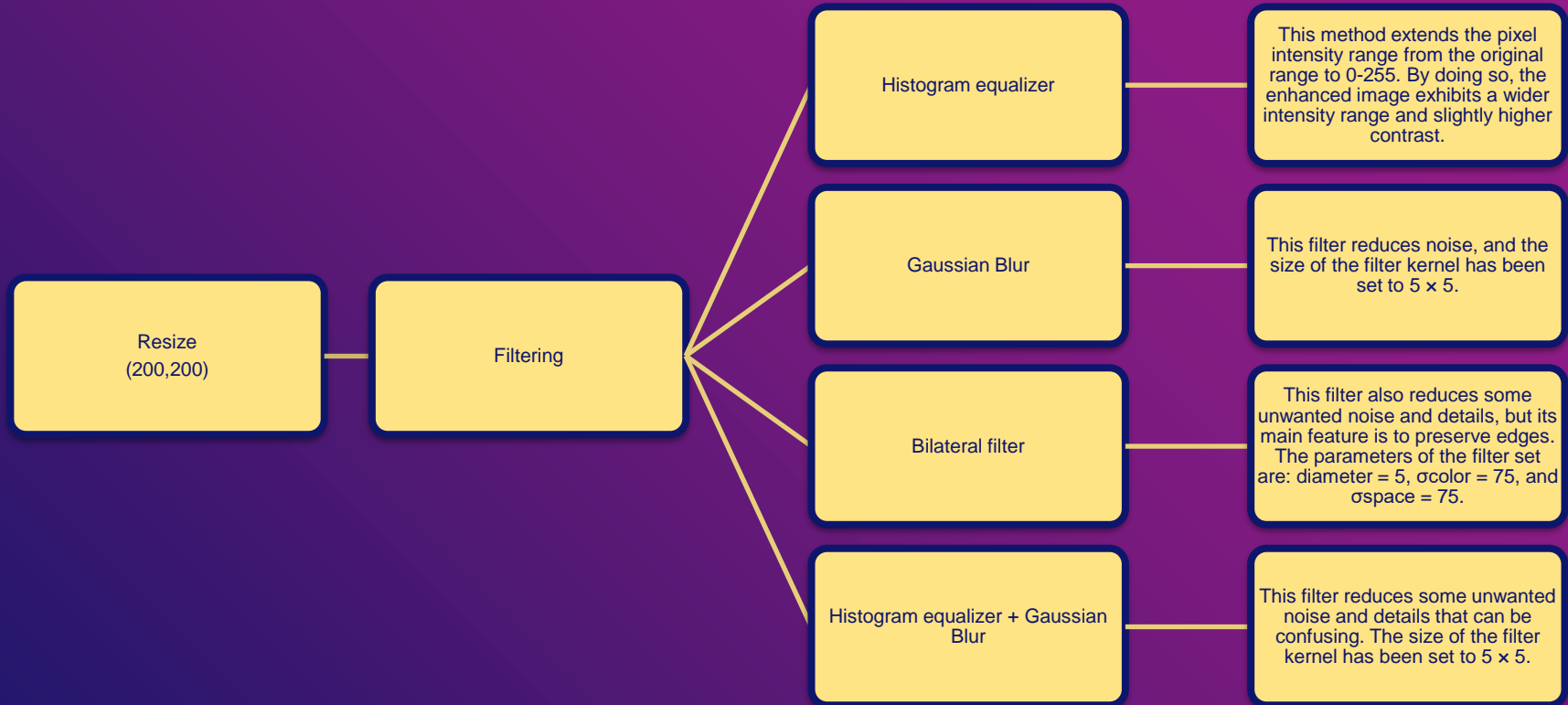
# DATASET SPLITTING

# DATA PROCESSING

## 03

Image resize
Image Enhacement
Image Denoising

# PREPROCESSING DATASET

Resize (200,200) — Filtering

**Histogram equalizer**

This method extends the pixel intensity range from the original range to 0-255. By doing so, the enhanced image exhibits a wider intensity range and slightly higher contrast.

**Gaussian Blur**

This filter reduces noise, and the size of the filter kernel has been set to 5 × 5.

**Bilateral filter**

This filter also reduces some unwanted noise and details, but its main feature is to preserve edges. The parameters of the filter set are: diameter = 5, σcolor = 75, and σspace = 75.

**Histogram equalizer + Gaussian Blur**

This filter reduces some unwanted noise and details that can be confusing. The size of the filter kernel has been set to 5 × 5.
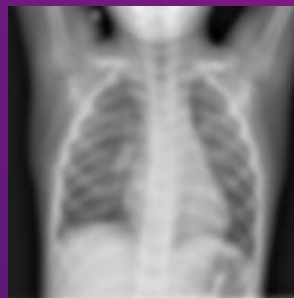
# RESIZING AND FILTERING
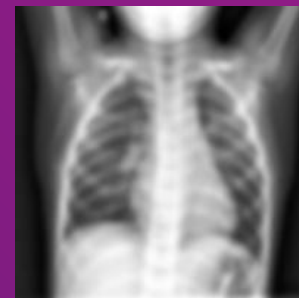
**Immagine Resized**

**Histogram equalizer**

**Gaussian Blur**

**Bilateral Filter**

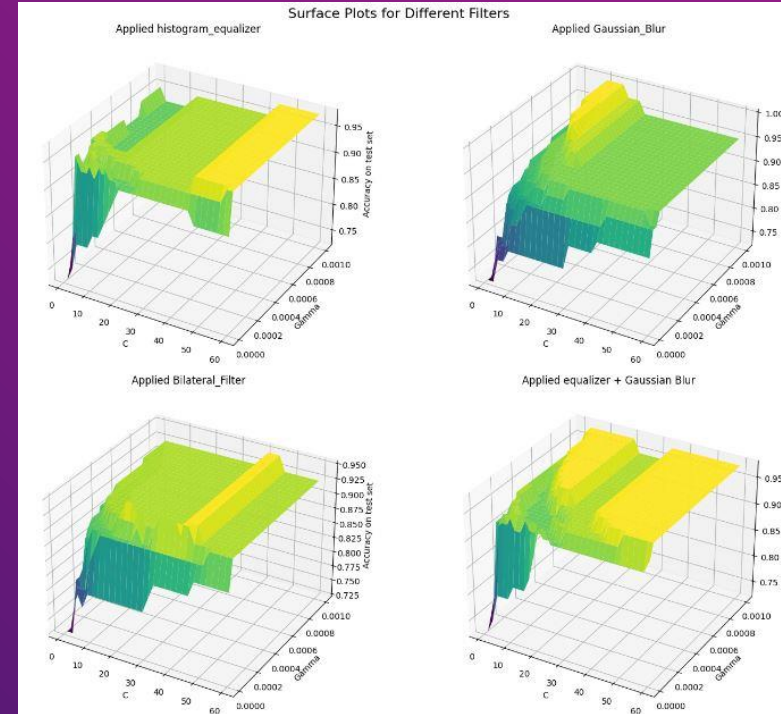**Histogram equalizer + Gaussian Blur**

# ESTIMATION OF PARAMETERS C AND GAMMA

C beetween 1 and 60

Gamma beetween 1e-05 e 0.01

The model's accuracy was evaluated on the test set for each type of filter *



*Considered a smaller dataset for computational time

# 04

# PERFORMANCE MEASURES

*Accuracy – Precision - Recall - F1- Score*

# RESULTS

Results of the model with:

- **Histogram Equalizer + Gaussian Blur**
- **C = 47**
- **Gamma = 0.01**
- Accuracy = 0.96
- Precision = 0.97
- Recall = 0.97
- F1 score = 0.97



Confusion Matrix

# 5 FOLD CROSS VALIDATION RESULTS

| | Training Sets | | | | Test Set | |
| Iteration 1 | | | | | | → $Error_1$ |
| Iteration 2 | | | | | | → $Error_2$ |
| Iteration 3 | | | | | | → $Error_3$ |
| Iteration 4 | | | | | | → $Error_4$ |
| Iteration 5 | | | | | | → $Error_5$ |

$$Error = \frac{1}{5}\sum_{i=1}^{5} Error_i$$

| | fit_time | score_time | test_accuracy | test_precision | test_recall | test_f1 |
|---|---|---|---|---|---|---|
| 1 fold | 264.046500 | 146.738816 | 0.943498 | 0.977893 | 0.942356 | 0.959796 |
| 2 fold | 219.482963 | 140.404387 | 0.963229 | 0.987130 | 0.961153 | 0.973968 |
| 3 fold | 222.070816 | 145.075061 | 0.965919 | 0.976190 | 0.976190 | 0.976190 |
| 4 fold | 223.573325 | 147.655448 | 0.966786 | 0.984713 | 0.968672 | 0.976627 |
| 5 fold | 168.155565 | 109.853444 | 0.858169 | 0.845572 | 0.981203 | 0.908353 |
| Mean results | 219.465834 | 137.945431 | 0.939520 | 0.954300 | 0.965915 | 0.958987 |

Thanks for your attention!

# OUR TEAM



## GIANMARCO BORRATA

## PASQUALE PIPICIELLO