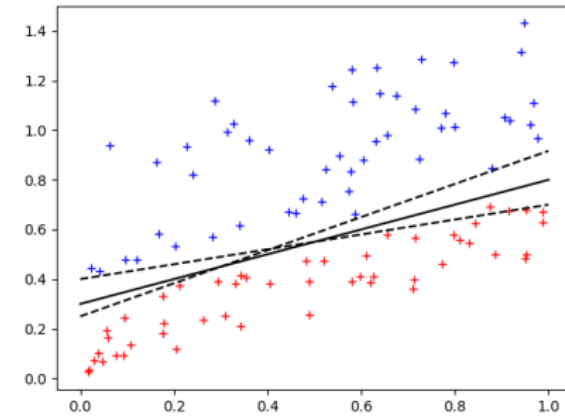
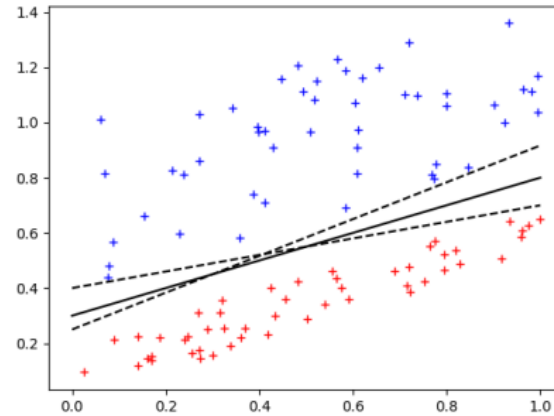
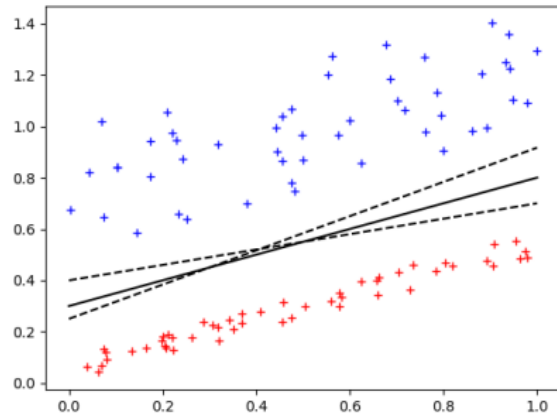


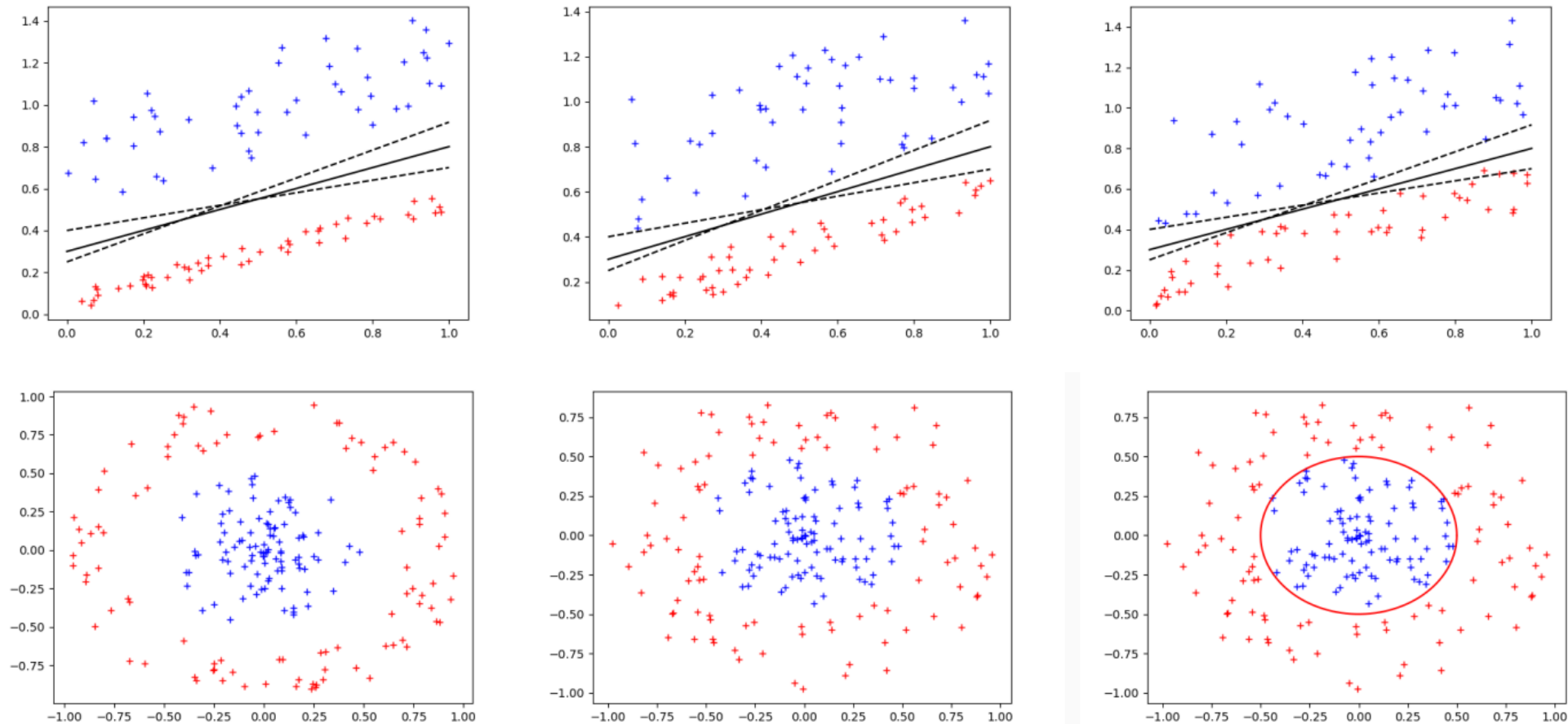
# Supervised Learning: SVM Method + Model evaluation methods

KASHTANOVA  
VICTORIYA  
INRIA, EPIONE

# SVM (Support Vector Machine) Method

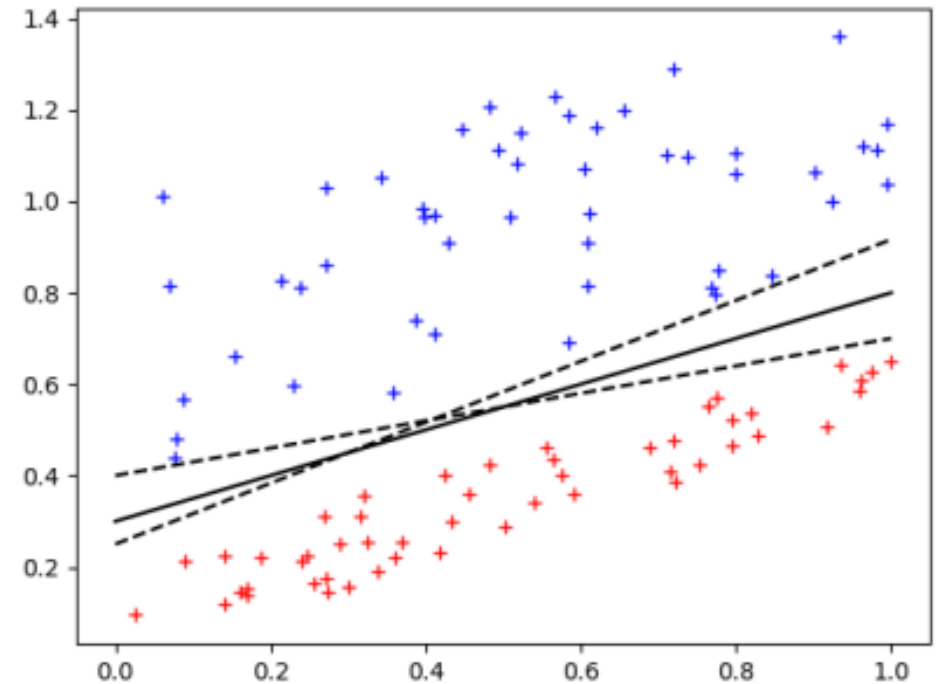


# SVM (Support Vector Machine) Method



# Linear SVM (Separable classes)

- Every data-sample :  $x \in \mathbb{R}^D$
- Decision Boundary :  $\mathcal{H} : w^T x + b = 0$
- Distance Measure of Hyper plane :  $d_{\mathcal{H}}(x_0) = \frac{|w^T x + b|}{\|w\|_2}$
- **Goal** :  $w^* = \operatorname{argmax}_w [\min_n d_{\mathcal{H}}(x_n)]$



# Linear SVM (Separable classes)

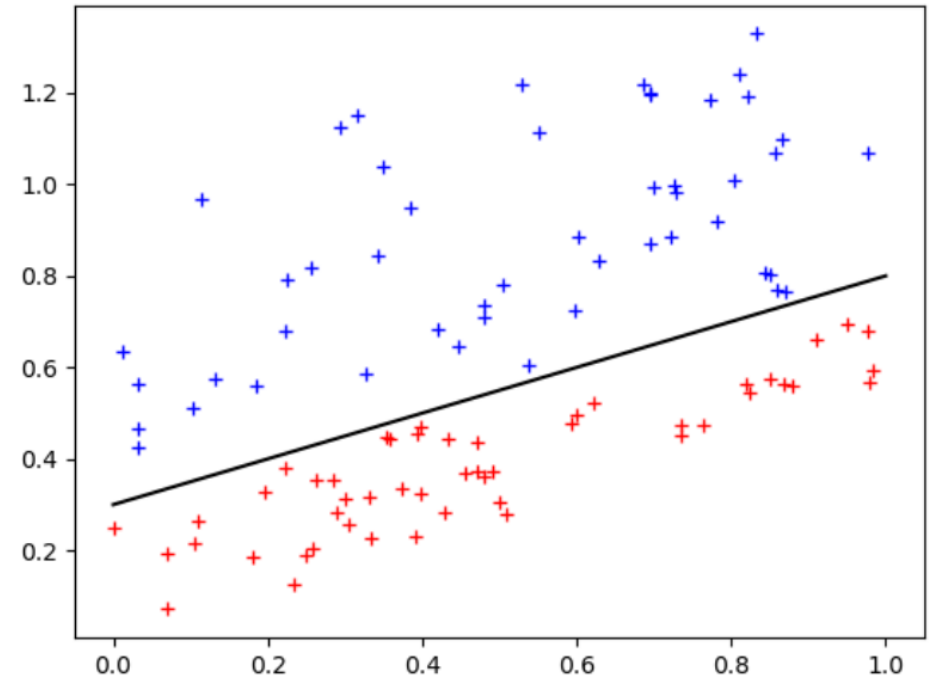
- Every data-sample :  $x \in \mathbb{R}^D$
- Decision Boundary :  $\mathcal{H} : w^T x + b = 0$
- Distance Measure of Hyper plane :  $d_{\mathcal{H}}(x_0) = \frac{|w^T x + b|}{\|w\|_2}$

- **Goal** :  $w^* = \operatorname{argmax}_w [\min_n d_{\mathcal{H}}(x_n)]$   
↓  $y_n[w^T x + b] = \begin{cases} \geq 0, & \text{class A} \\ < 0, & \text{class B} \end{cases}$

$$w^* = \operatorname{argmax}_w \frac{1}{\|w\|_2} [\min_n y_n [w^T x + b]]$$

↓ Let  $\min_n y_n [w^T x + b] = 1$

$$w^* = \operatorname{argmax}_w \frac{1}{\|w\|_2}$$



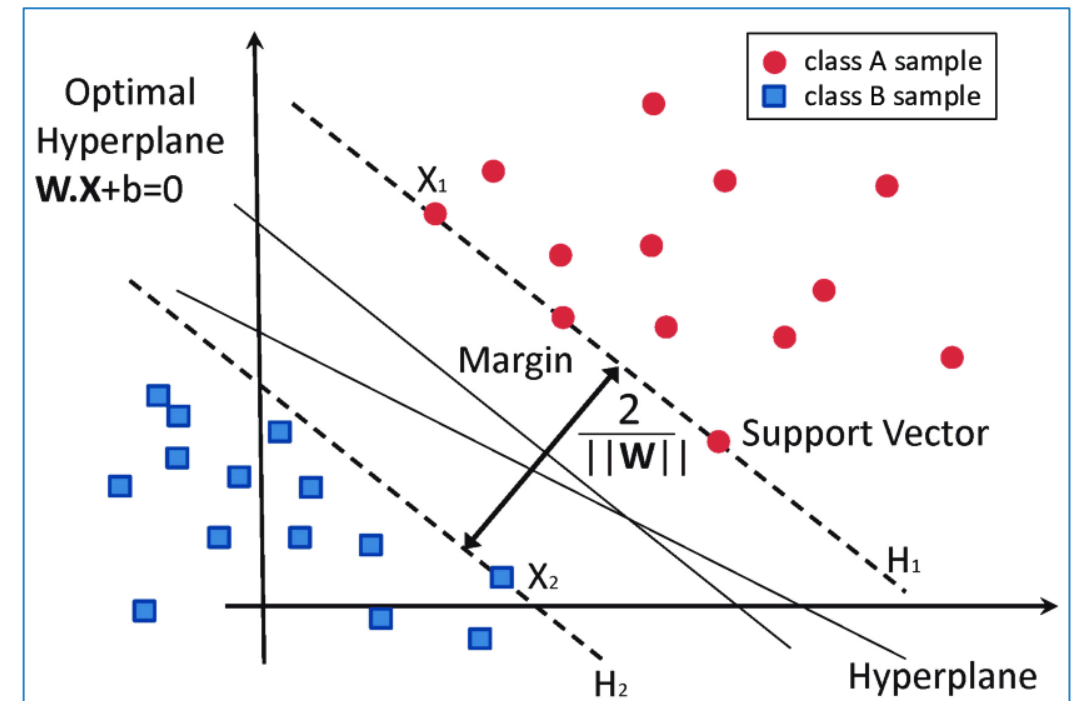
# Linear SVM (Separable classes)

- Every data-sample :  $x \in \mathbb{R}^D$
- Decision Boundary :  $\mathcal{H} : w^T x + b = 0$
- Distance Measure of Hyper plane :  $d_{\mathcal{H}}(x_0) = \frac{|w^T x + b|}{\|w\|_2}$

- **Goal** :  $w^* = \operatorname{argmax}_w [\min_n d_{\mathcal{H}}(x_n)]$

$$\Downarrow \quad y_n[w^T x + b] = \begin{cases} \geq 0, & \text{class A} \\ < 0, & \text{class B} \end{cases}$$
$$w^* = \operatorname{argmax}_w \frac{1}{\|w\|_2} [\min_n y_n[w^T x + b]]$$

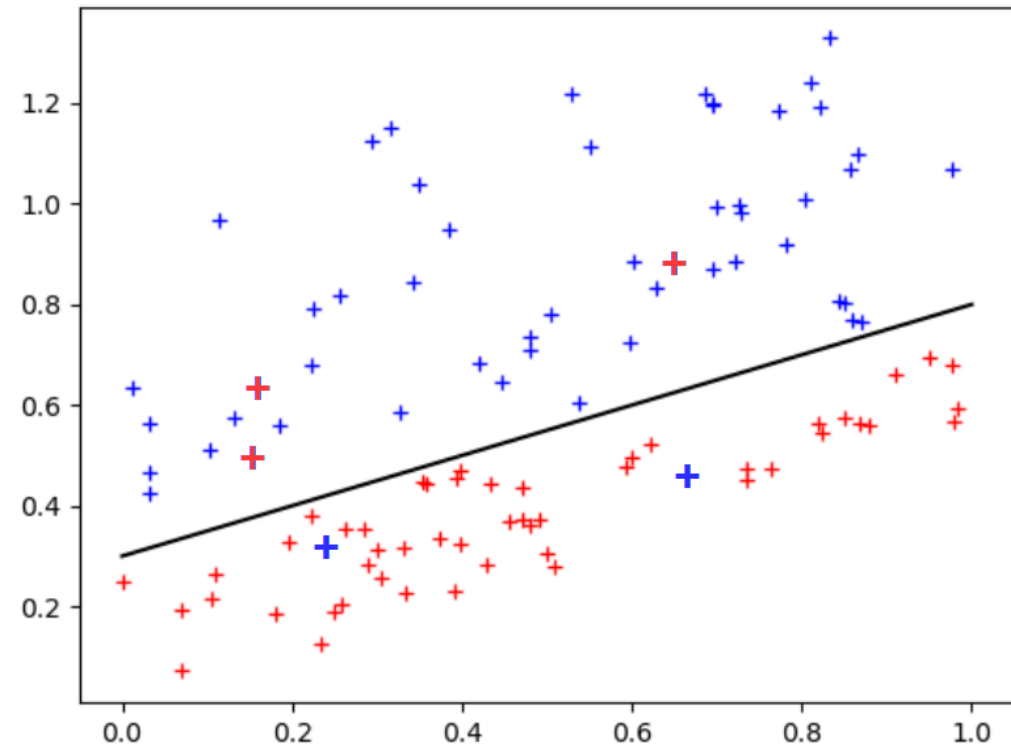
$$\Downarrow \quad y_n[w^T x + b] \geq 1, \forall n$$
$$\min_w \frac{1}{2} \|w\|_2^2$$



# Linear SVM (General case)

**Final form SVM optimization problem :**

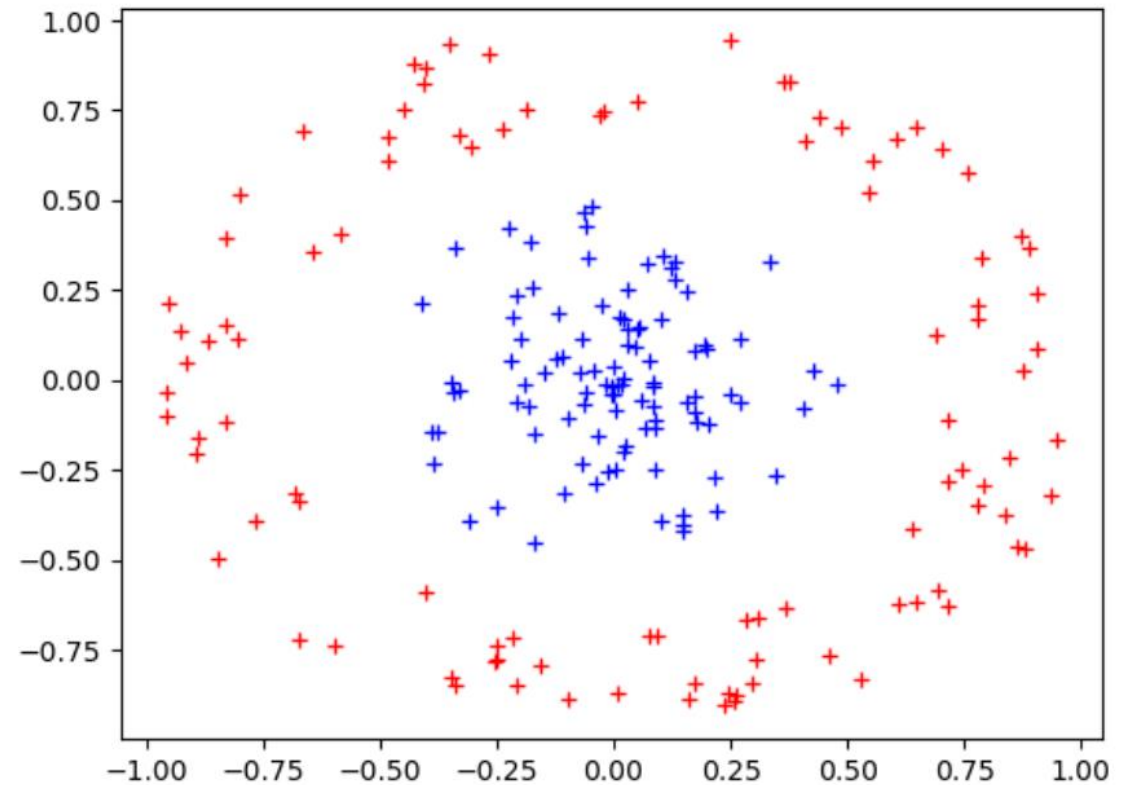
$$\begin{aligned} \min_w \quad & \frac{1}{2} \|w\|_2^2 + C \sum_n \xi_n \\ \text{s. t. } & y_n [w^T x + b] \geq 1 - \xi_n, \forall n \\ & \xi_n \geq 0, \forall n \end{aligned}$$



# Linear SVM (General case)

**Final form SVM optimization problem :**

$$\begin{aligned} \min_w \quad & \frac{1}{2} \|w\|_2^2 + C \sum_n \xi_n \\ \text{s. t. } & y_n [w^T x + b] \geq 1 - \xi_n, \forall n \\ & \xi_n \geq 0, \forall n \end{aligned}$$



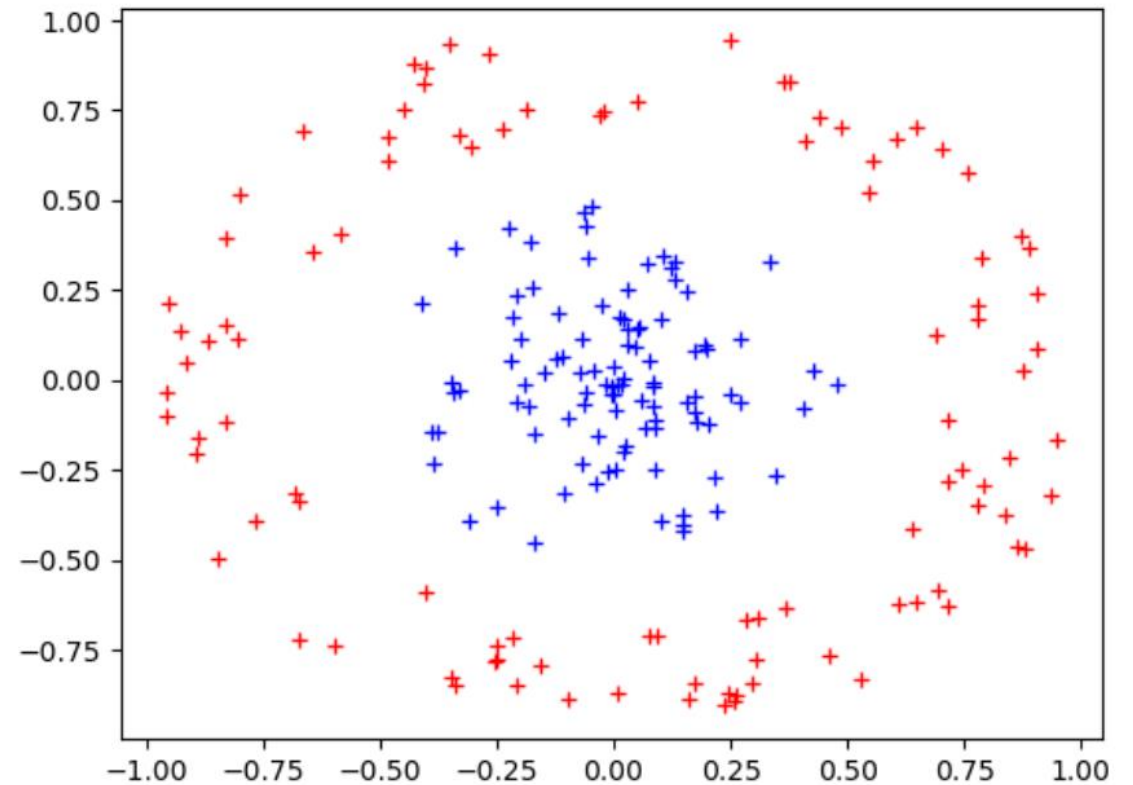


# Linear SVM (General case)

**Final form SVM optimization problem :**

$$\begin{aligned} \min_w \quad & \frac{1}{2} \|w\|_2^2 + C \sum_n \xi_n \\ \text{s. t. } & y_n [w^T x + b] \geq 1 - \xi_n, \forall n \\ & \xi_n \geq 0, \forall n \end{aligned}$$

$$\varphi(x) : \mathbb{R}^D \rightarrow \mathbb{R}^M$$

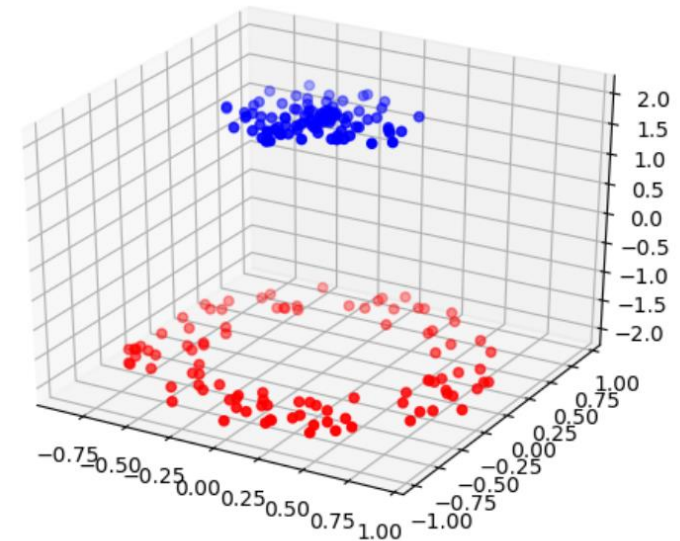
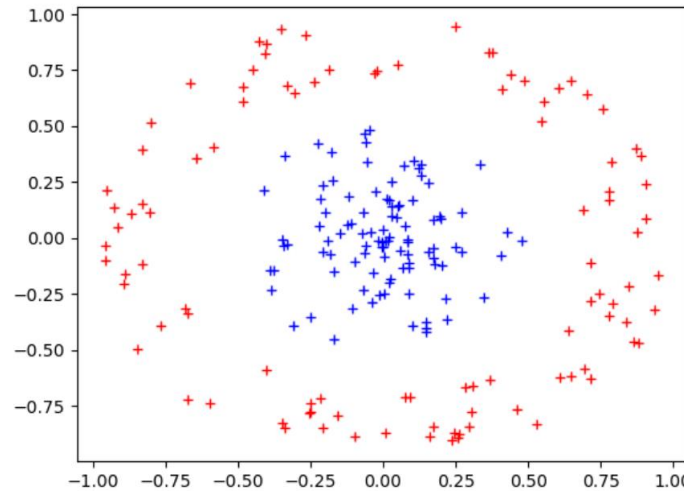


# Linear SVM (General case)

**Final form SVM optimization problem :**

$$\begin{aligned} \min_w \quad & \frac{1}{2} \|w\|_2^2 + C \sum_n \xi_n \\ \text{s. t. } & y_n [w^T x + b] \geq 1 - \xi_n, \forall n \\ & \xi_n \geq 0, \forall n \end{aligned}$$

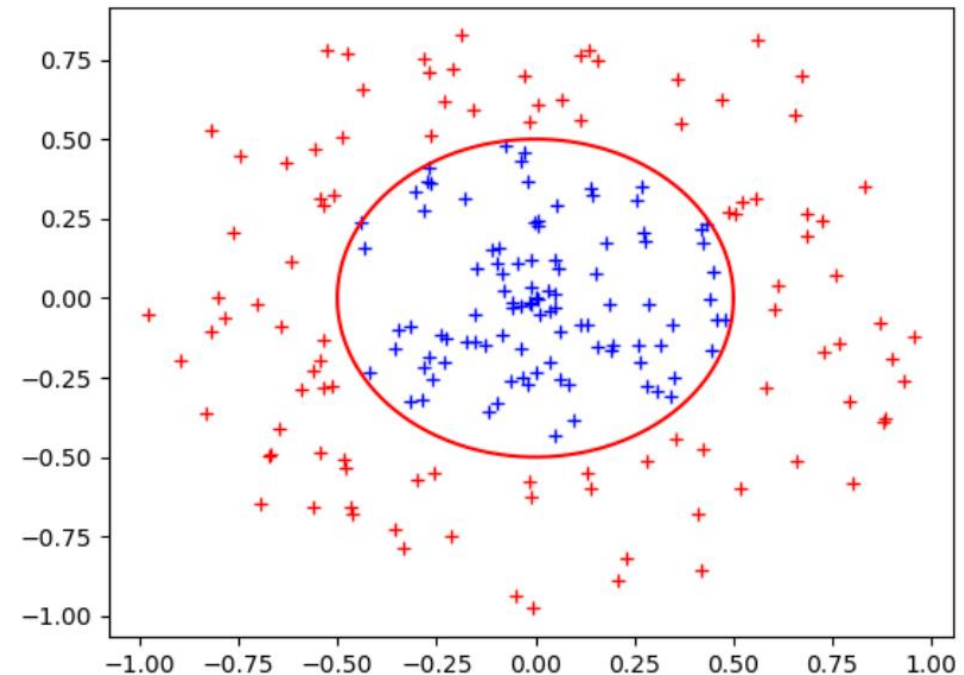
$$\varphi(x) : \mathbb{R}^D \rightarrow \mathbb{R}^M$$



# Linear SVM (General case)

**The Most Final form SVM optimization problem :**

$$\begin{aligned} \min_{w, \{\xi_n\}} & \frac{1}{2} \|w\|_2^2 + C \sum_n \xi_n \\ \text{s.t. } & y_n [w^T \varphi(x) + b] \geq 1 - \xi_n, \forall n \\ & \xi_n \geq 0, \forall n \end{aligned}$$

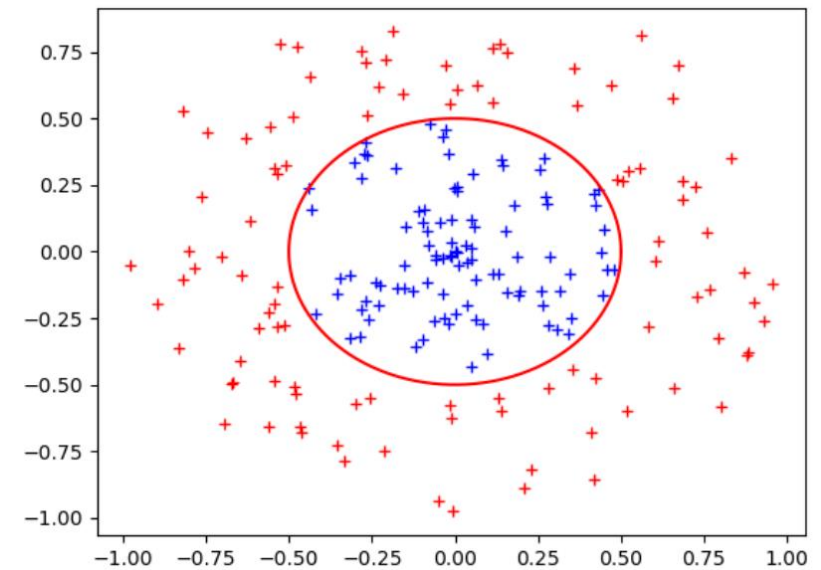


# Kernel SVM

**Problem :** Don't know the nature of  $\varphi(x)$  transformation

**Idea :** Rewrite Optimization problem without determining  $\varphi(x)$

It's possible via applying Lagrange multipliers optimization method, Dual Form of SVM, Mercer Theorem etc.



# Kernel SVM

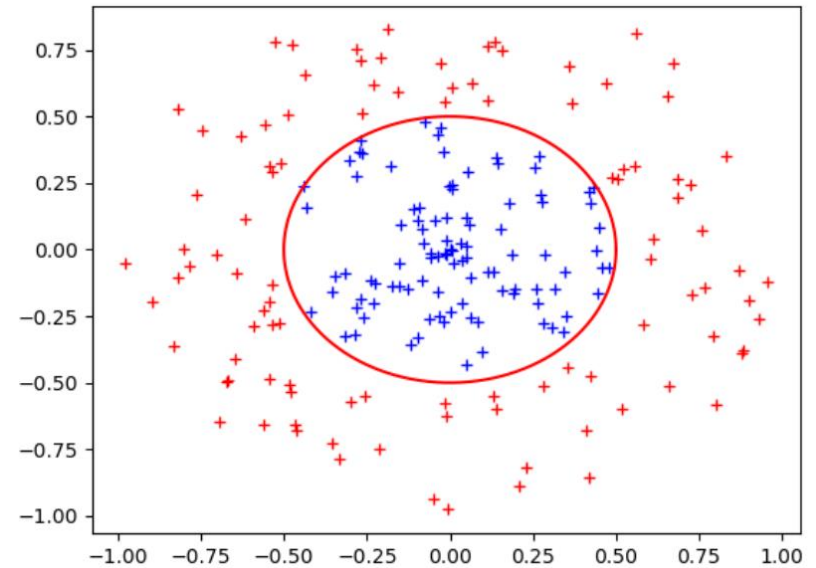
**Problem :** Don't know the nature of  $\varphi(x)$  transformation

**Idea :** Rewrite Optimization problem without determining  $\varphi(x)$

$$\begin{aligned} \min_{w, \{\xi_n\}} \quad & \frac{1}{2} \|w\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [w^T \varphi(x) + b] \geq 1 - \xi_n, \xi_n \geq 0, \forall n \end{aligned}$$

**Dual Form :**  $\varphi * \varphi^T = K(x, x')$

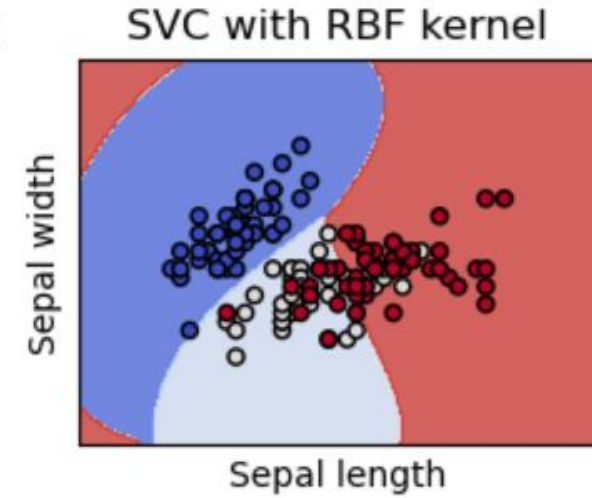
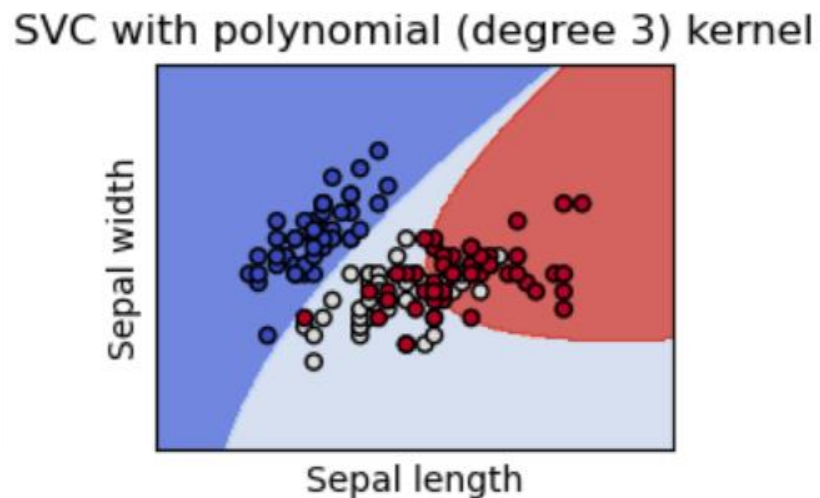
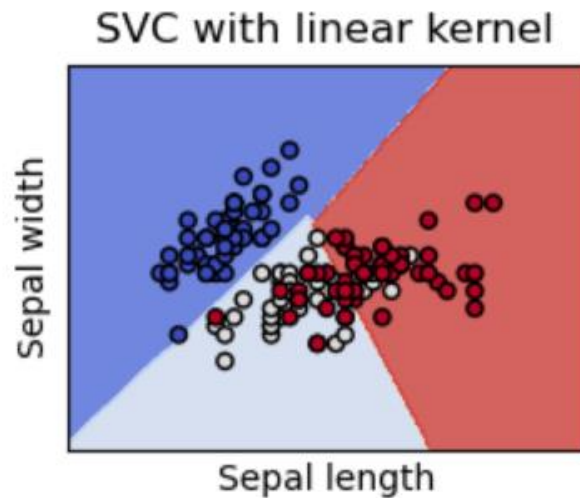
$$\begin{cases} \max L(\lambda) = \sum_{k=1}^m \lambda_k - \sum_{k=1}^m \sum_{l=1}^m \lambda_k \lambda_l y_k y_l K(x_k, x_l) \\ \text{s.t.} \quad 0 \leq \lambda_k \leq C \end{cases}$$



# Kernel SVM

Kernel exemples :

- linear:  $k(x, x') = x^T x'$
- polynomial:  $k(x, x') = (x^T x' + 1)^d$  of different power  $d = 2, 3, \dots$
- gauss-RBG:  $k(x, x') = \exp(-\frac{1}{2\sigma} |x - x'|^2)$



# SVM : Advantages & Drawbacks

---

- Good interpretability of the model
- Effective in high dimensional spaces
- Memory efficient



- **High sensitivity to the noise** in input data
- Slow training on large dataset



# Model evaluation methods



# Performance measures for Classification

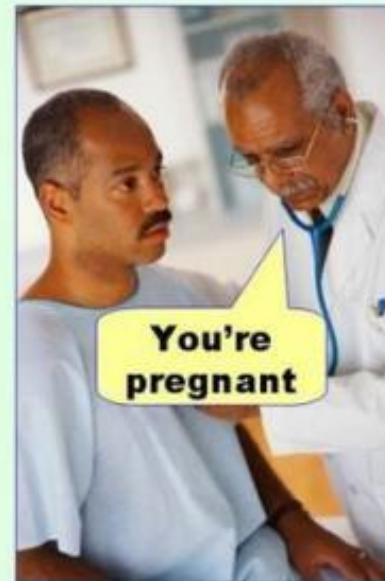
---

- Simple Accuracy
- Precision
- Recall
- F-beta measure
- ROC (and AUC)

# Confusion matrix

		Predicted class	
		$P$	$N$
Actual Class	$P$	True Positives (TP)	False Negatives (FN)
	$N$	False Positives (FP)	True Negatives (TN)

**Type I error**  
(false positive)



**Type II error**  
(false negative)



# Metrics

---

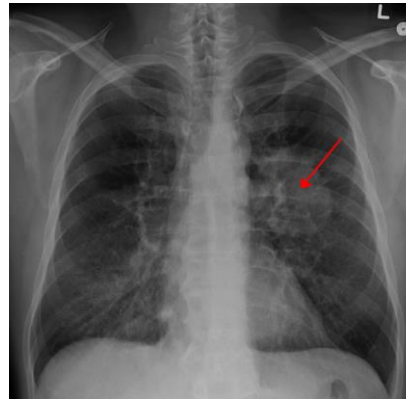
		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

- **Accuracy** =  $\frac{(TP + TN)}{(TP + TN + FP + FN)} = \frac{(TP + TN)}{N}$   
(fraction of correct predictions)

# Metrics

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

- **Accuracy** =  $\frac{(TP + TN)}{(TP + TN + FP + FN)} = \frac{(TP + TN)}{N}$   
(fraction of correct predictions)



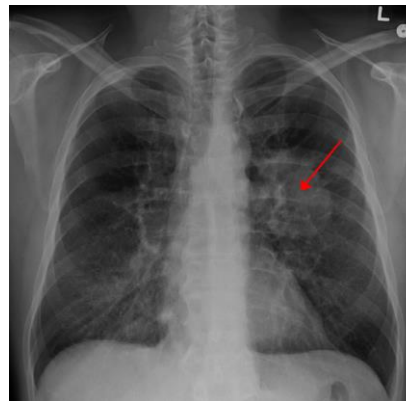
**9990** Not tumor

**10** Tumor

# Metrics

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

- $$\text{Accuracy} = \frac{(TP + TN)}{(TP + TN + FP + FN)} = \frac{(TP + TN)}{N}$$
  
(fraction of correct predictions)



9990 Not tumor  
10 Tumor

All is NOT tumor :  
Accuracy = 99,9%

# Metrics

---

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

- **Precision** =  $\frac{TP}{(TP+FP)}$   
(fraction of correctly predicted positive values to all values predicted positive)
- **Recall** =  $\frac{TP}{(TP + FN)}$   
(completeness, fraction of correctly predicted positive values to all positive values)

# Metrics

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

- **Precision** =  $\frac{TP}{(TP+FP)}$   
(fraction of correctly predicted positive values to all values predicted positive)  
  
 $\Rightarrow \text{Precision} = \frac{0}{0}$
- **Recall** =  $\frac{TP}{(TP + FN)}$   
(completeness, fraction of correctly predicted positive values to all positive values)  
  
 $\Rightarrow \text{Recall} = \frac{0}{0+10} = 0$

# Example 1&2

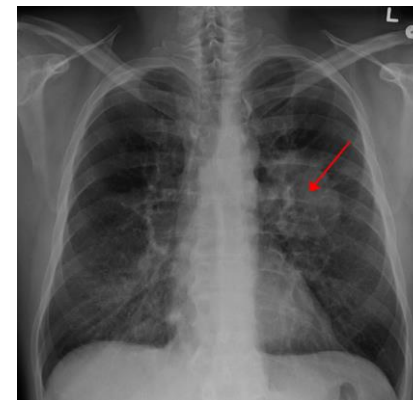
		Predicted class	
		$P$	$N$
Actual Class	$P$	True Positives (TP)	False Negatives (FN)
	$N$	False Positives (FP)	True Negatives (TN)

All is NOT Tumor :

Accuracy = 99,9%

Precision =  $\frac{0}{0}$

Recall = 0



9990 Not tumor  
10 Tumor



# Example 1&2

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

**All is NOT Tumor :**

Accuracy = 99,9%

Precision =  $\frac{0}{0}$

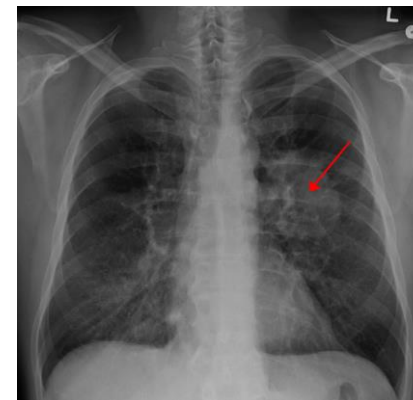
Recall = 0

**All is Tumor :**

Accuracy = 0,1%

Precision = 0,001

Recall = 1



9990 Not tumor  
10 Tumor

# Example 3

---

## System 1 :

- Precision = 70%
- Recall = 60%



## System 2 :

- Precision = 80%
- Recall = 50%

# Metrics

---

- $F_{\beta} = \frac{1}{\left(\beta * \frac{1}{\text{Precision}} + (1-\beta) * \frac{1}{\text{Recall}}\right)}$  (greater  $\beta$ , greater importance of Precision)
- $F_1 = \frac{2TP}{(2TP + FP + FN)}$  (harmonic mean of precision and recall,  $\beta = 0.5$ )

# Example 3

---

## System 1 :

- Precision = 70%
- Recall = 60%



## System 2 :

- Precision = 80%
- Recall = 50%

$$\beta = 0.5$$

$$F_{\beta} =$$

$$F_{\beta} =$$

# Example 3

---

## System 1 :

- Precision = 70%
- Recall = 60%



## System 2 :

- Precision = 80%
- Recall = 50%

$$\beta = 0.5$$

$$F_{\beta} = 0.6461$$



$$F_{\beta} = 0.6153$$

$$\beta = 0.95 \text{ (Precision is more important)}$$

$$F_{\beta} =$$

$$F_{\beta} =$$

# Example 3

---

## System 1 :

- Precision = 70%
- Recall = 60%



## System 2 :

- Precision = 80%
- Recall = 50%

**$\beta = 0.5$**

$$F_{\beta} = 0.6461$$



$$F_{\beta} = 0.6153$$

**$\beta = 0.95$**  (Precision is more important)

$$F_{\beta} = 0.6942$$



$$F_{\beta} = 0.7766$$

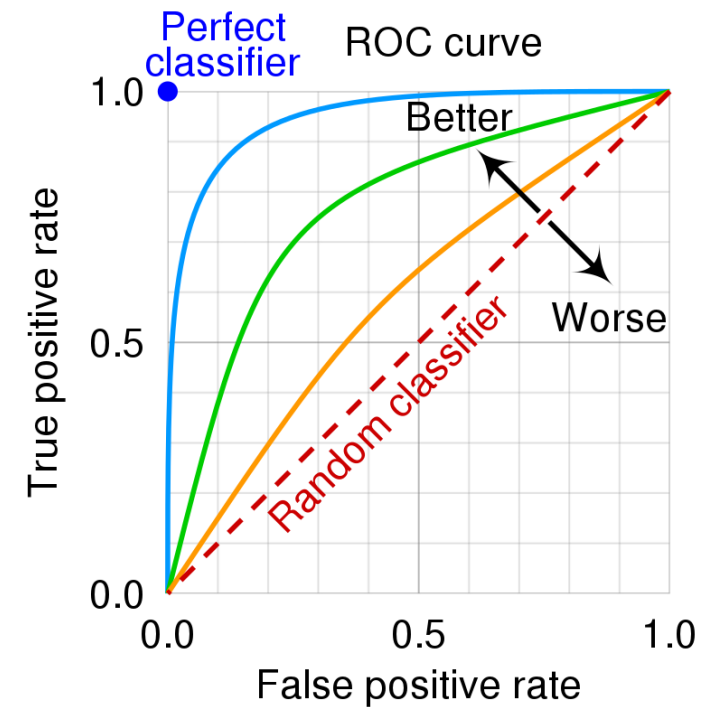
# ROC (Receiver operating characteristic)

- **Sensitivity** =  $\frac{TP}{P} = \frac{TP}{(TP + FN)} = \text{TPR}$  (True positive rate, TPR, Recall)
- **Specificity** =  $\frac{TN}{N} = \frac{TN}{(FP + TN)} = 1 - \text{FPR}$  (True negative rate, TNR)

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

# ROC (AUC)

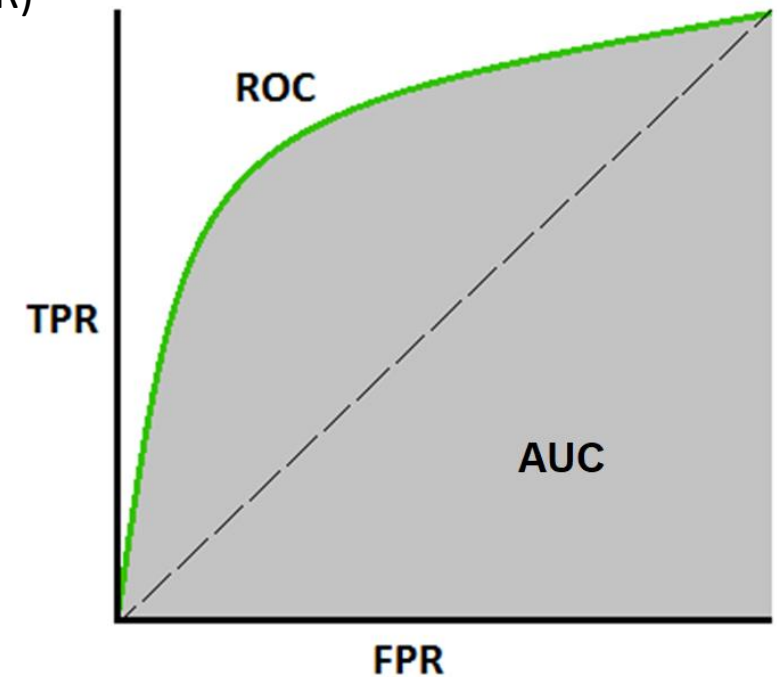
- **Sensitivity** =  $\frac{TP}{P} = \frac{TP}{(TP + FN)} = \mathbf{TPR}$  (True positive rate, TPR, Recall)
- **Specificity** =  $\frac{TN}{N} = \frac{TN}{(FP + TN)} = \mathbf{1 - FPR}$  (True negative rate, TNR)





# ROC (AUC)

- **Sensitivity** =  $\frac{TP}{P} = \frac{TP}{(TP + FN)} = \mathbf{TPR}$  (True positive rate, TPR, Recall)
- **Specificity** =  $\frac{TN}{N} = \frac{TN}{(FP + TN)} = \mathbf{1 - FPR}$  (True negative rate, TNR)



# Sources

---

- MIT course “Introduction to Computational Thinking and Data Science” (Prof. Eric Grimson, Prof. John Guttag)
- Open Machine Learning Course (by Yury Kashnitsky, [mlcourse.ai](https://mlcourse.ai))
- YouTube lectures “Algorithms and Concepts” (by CodeEmporium)