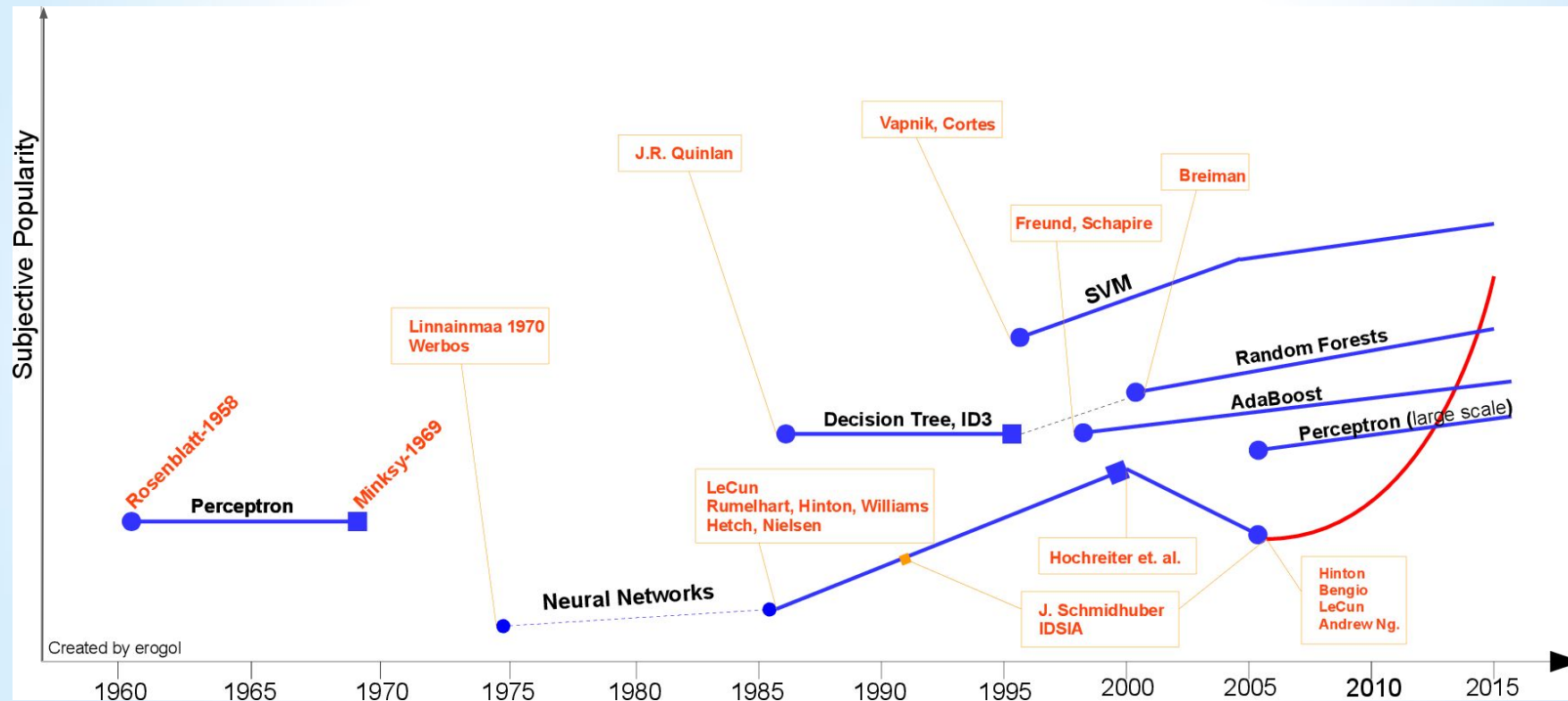


Support Vector Machines

Lior Sidi & Efrat Egozi



Models History



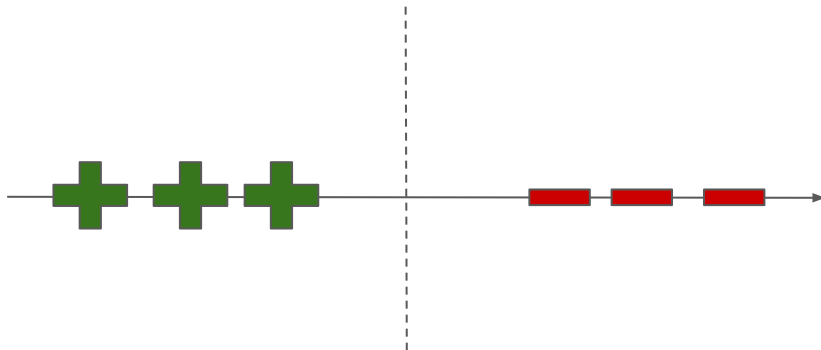
Motivation

Good for difficult problems with limited data (<10K data points)

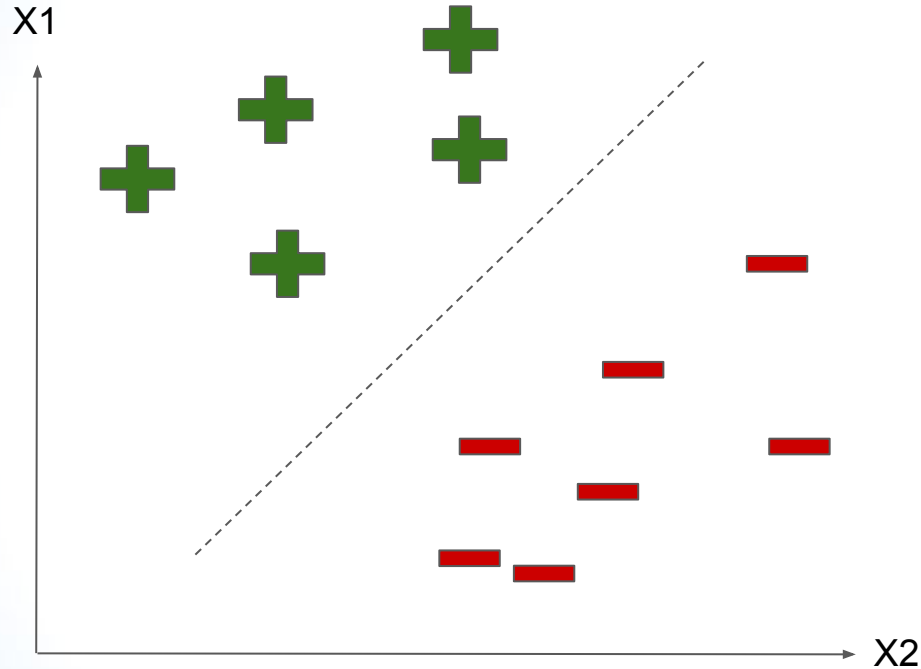
- Face Detection
- Text Classification
- Protein Fold and Remote Homology Detection
- Handwriting Recognition



Linear Classification - 1 dim



Linear Classification - 2 dim

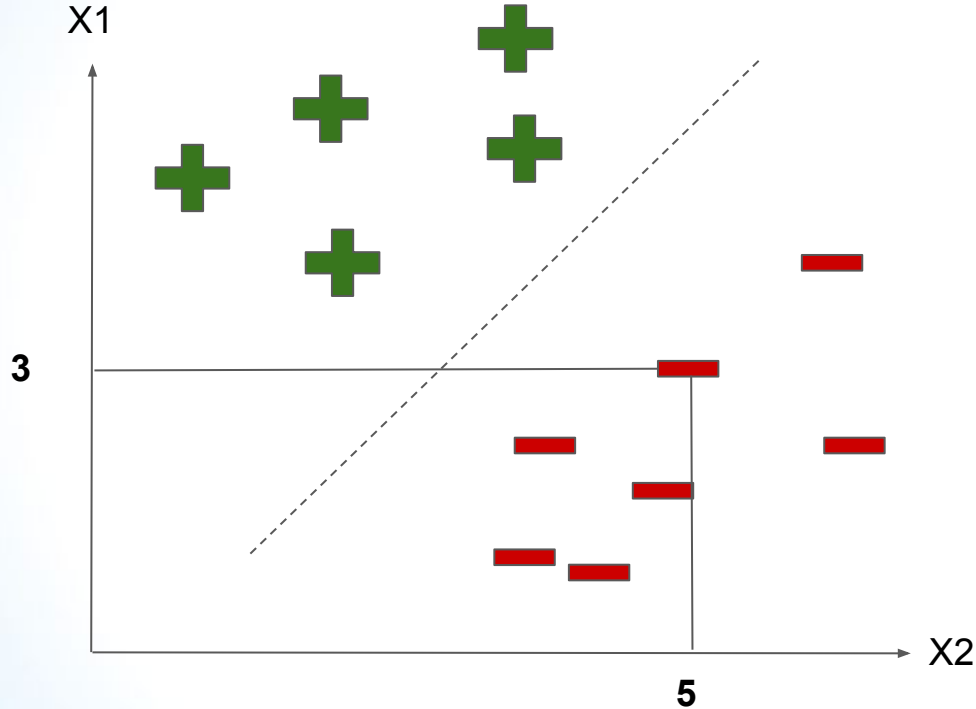


$$w_1 \cdot x_1 = w_2 \cdot x_2 + b$$

$$2 \cdot x_1 = 4 \cdot x_2 - 3$$

$$2 \cdot x_1 - 4 \cdot x_2 + 3 = 0$$

Linear Classification - 2 dim



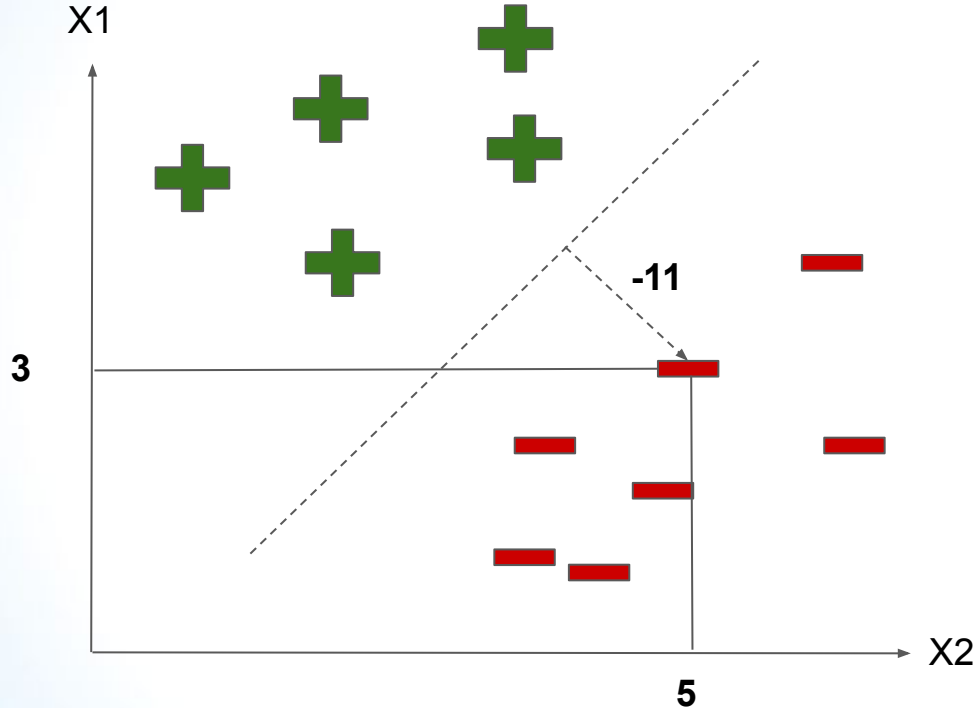
$$w_1 \cdot x_1 = w_2 \cdot x_2 + b$$

$$2 \cdot x_1 = 4 \cdot x_2 - 3$$

$$2 \cdot x_1 - 4 \cdot x_2 + 3 = 0$$

$$2 \cdot 3 - 4 \cdot 5 + 3 =$$

Linear Classification - 2 dim



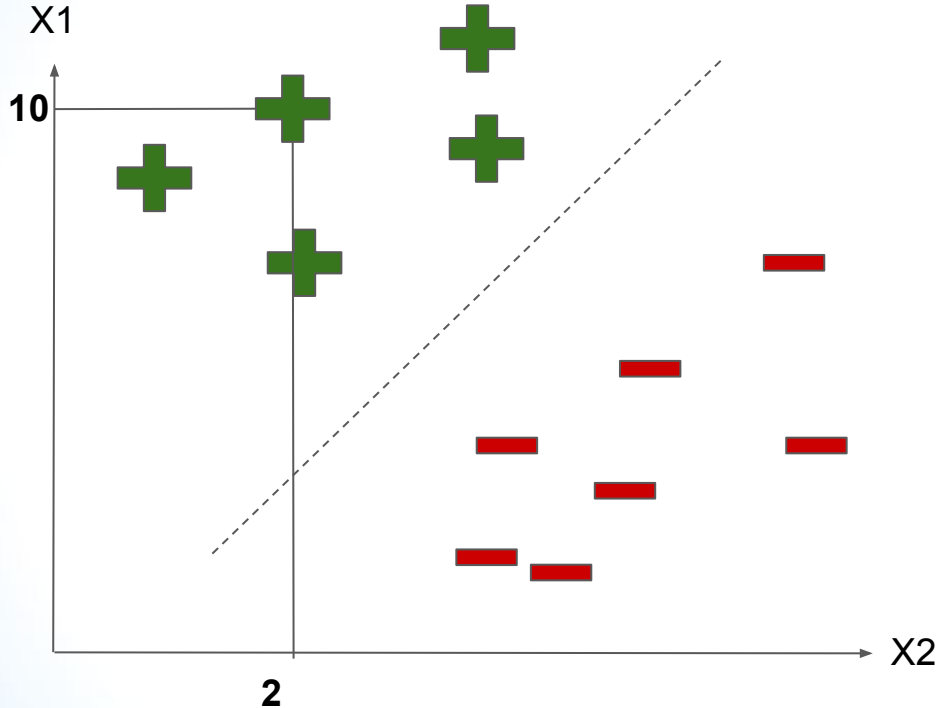
$$w_1 \cdot x_1 = w_2 \cdot x_2 + b$$

$$2 \cdot x_1 = 4 \cdot x_2 - 3$$

$$2 \cdot x_1 - 4 \cdot x_2 + 3 = 0$$

$$2 \cdot 3 - 4 \cdot 5 + 3 = 6 - 20 + 3 = -11$$

Linear Classification - 2 dim



$$w_1 \cdot x_1 = w_2 \cdot x_2 + b$$

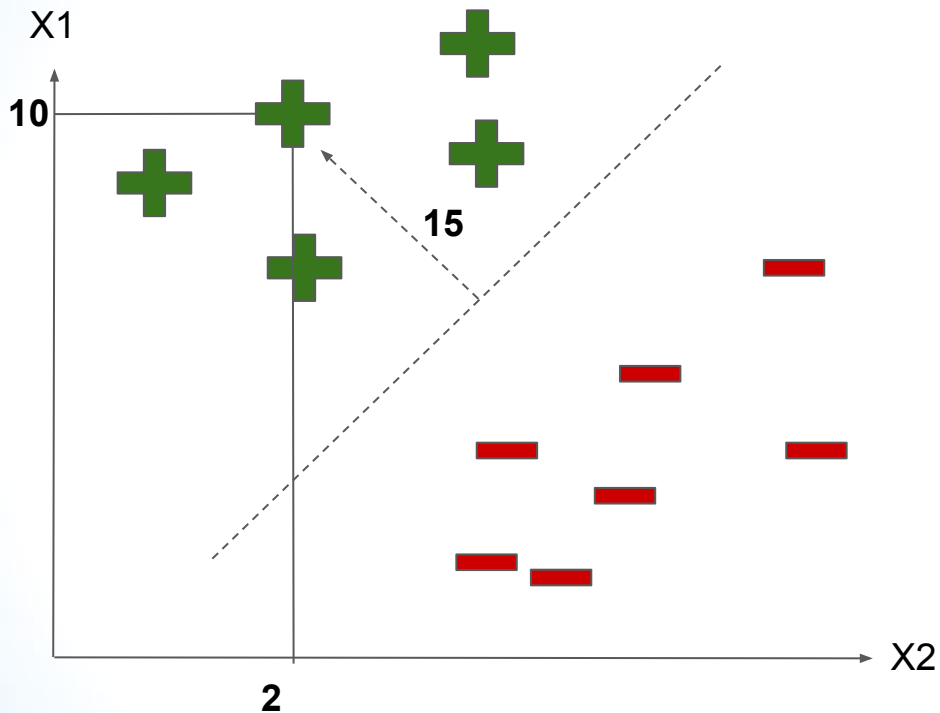
$$2 \cdot x_1 = 4 \cdot x_2 - 3$$

$$2 \cdot x_1 - 4 \cdot x_2 + 3 = 0$$

$$2 \cdot 3 - 4 \cdot 5 + 3 = 6 - 20 + 3 = -11$$

$$2 \cdot 10 - 4 \cdot 2 + 3 = 20 - 8 + 3 = 15$$

Linear Classification - 2 dim



$$w_1 \cdot x_1 = w_2 \cdot x_2 + b$$

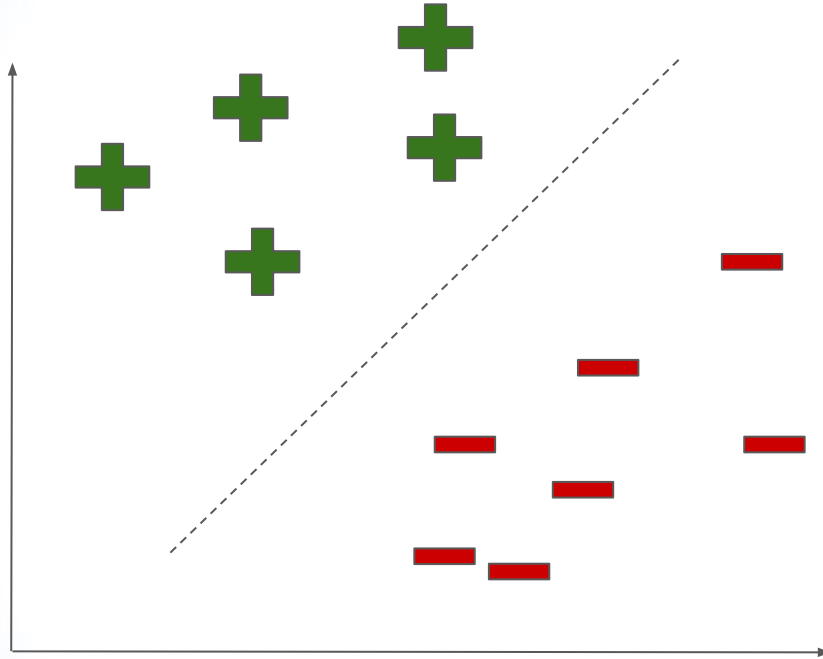
$$2 \cdot x_1 = 4 \cdot x_2 - 3$$

$$2 \cdot x_1 - 4 \cdot x_2 + 3 = 0$$

$$2 \cdot 3 - 4 \cdot 5 + 3 = 6 - 20 + 3 = -11$$

$$2 \cdot 10 - 4 \cdot 2 + 3 = 20 - 8 + 3 = 15$$

Linear Classification - 2 dim

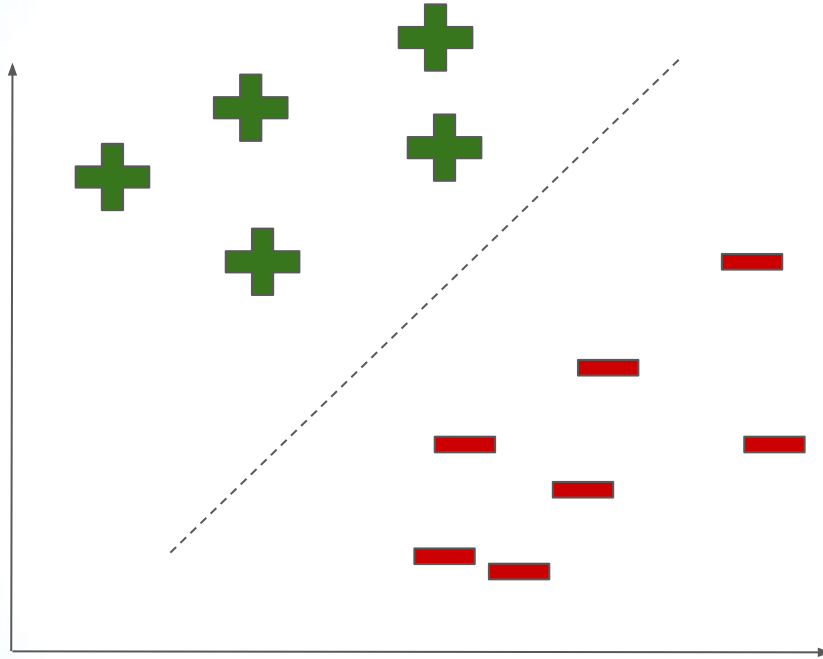


$$w_1 \cdot x_1 = w_2 \cdot x_2 + b$$

$$f(X, W) = w_1 \cdot x_1 + w_2 \cdot x_2 =$$

For simplicity We are going to eliminate the bias term
Which can be added as a vectors of ones

Linear Classification - 2 dim

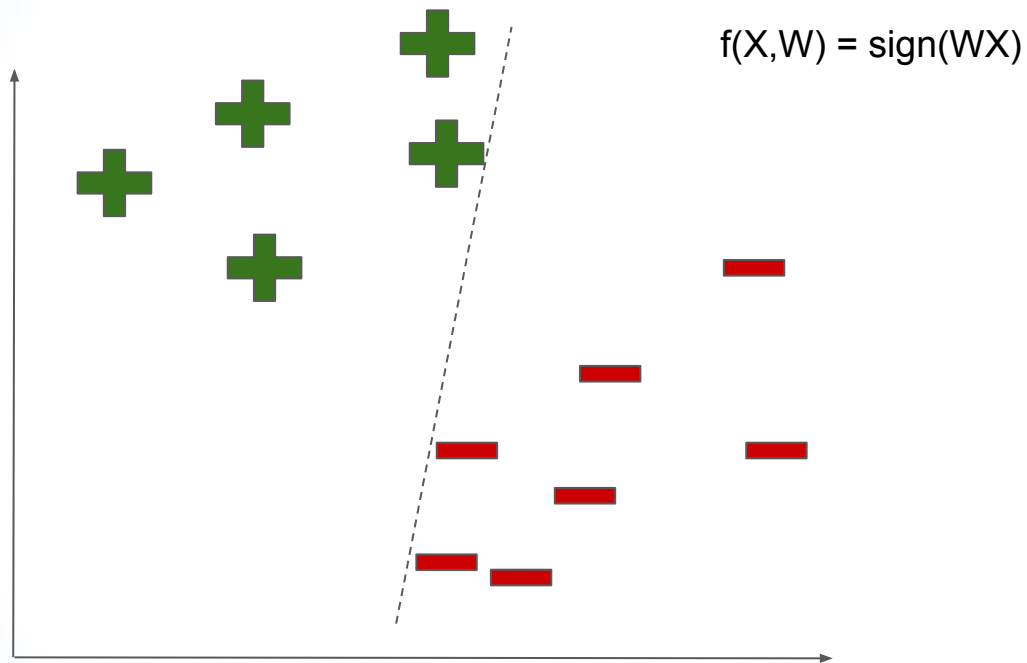


$$w_1 \cdot x_1 = w_2 \cdot x_2 + b$$

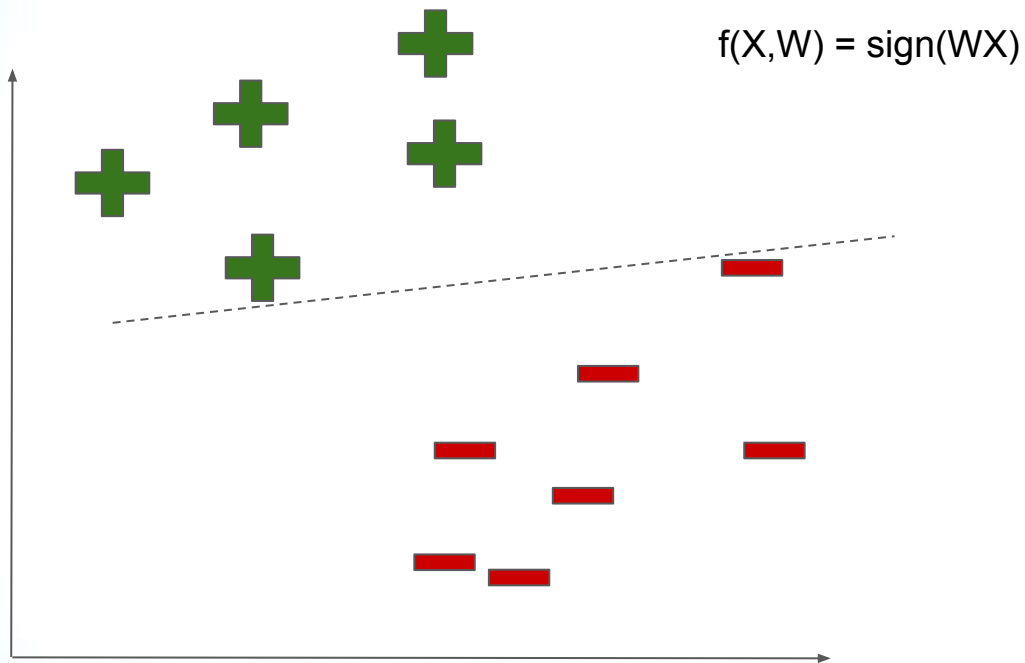
$$\begin{aligned} f(X, W) &= w_1 \cdot x_1 + w_2 \cdot x_2 = \\ &= \sum W X = W^t X = 0 \\ &\Rightarrow \text{sign}(W X) \end{aligned}$$

For simplicity We write
 $W^t X$ as $W X$

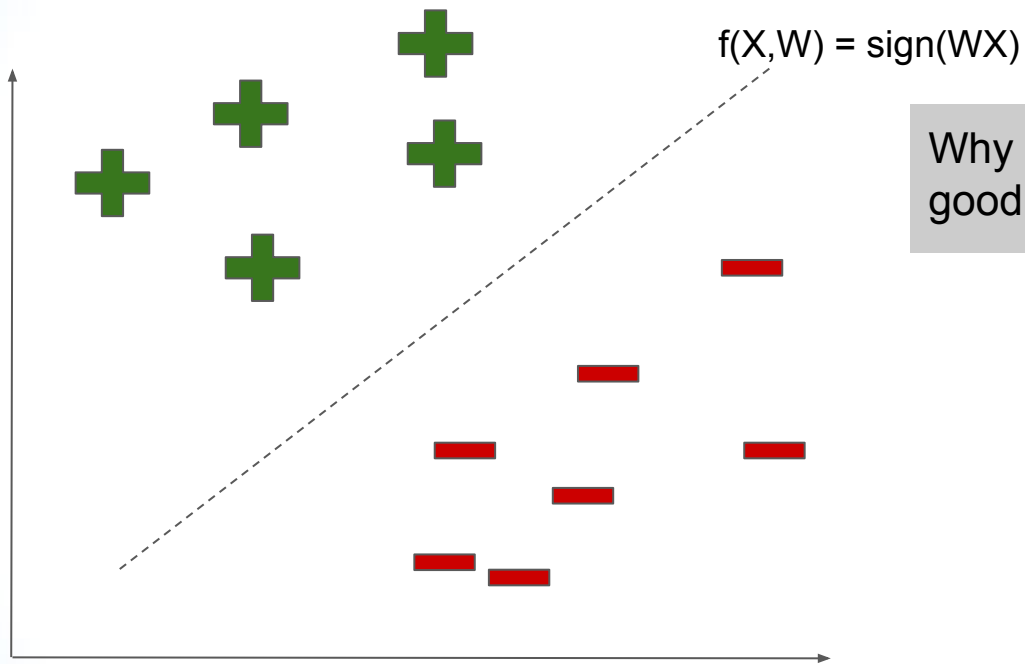
Linear Classification



Linear Classification

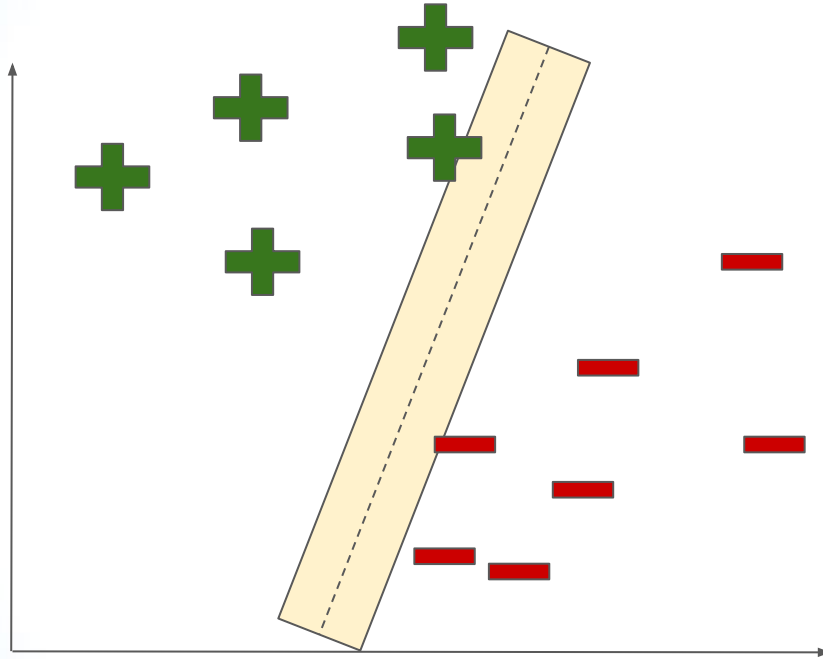


Linear Classification



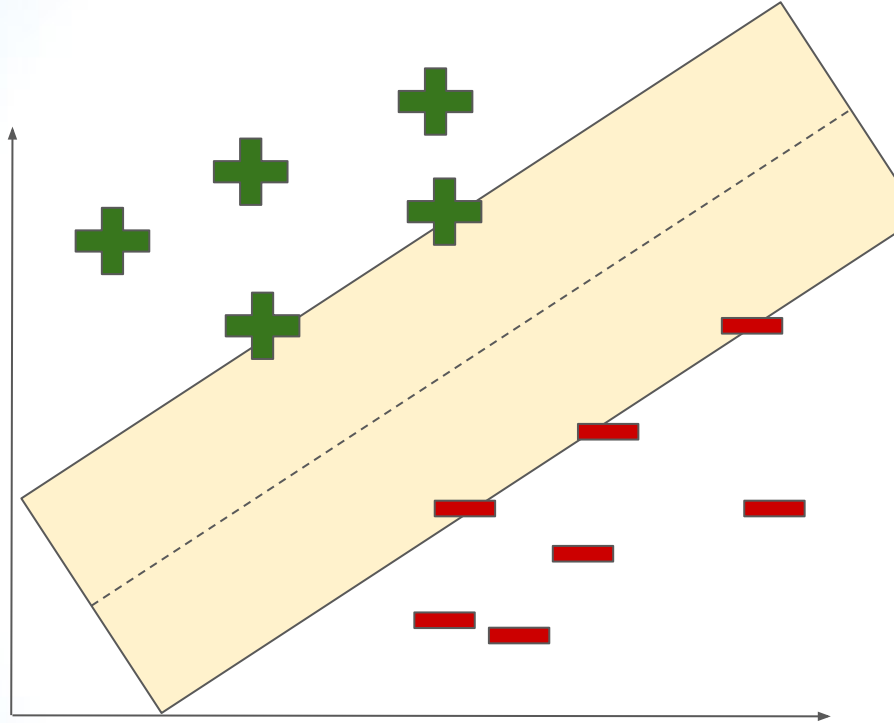
Why is this seems as a good separator?

Classifier Margin



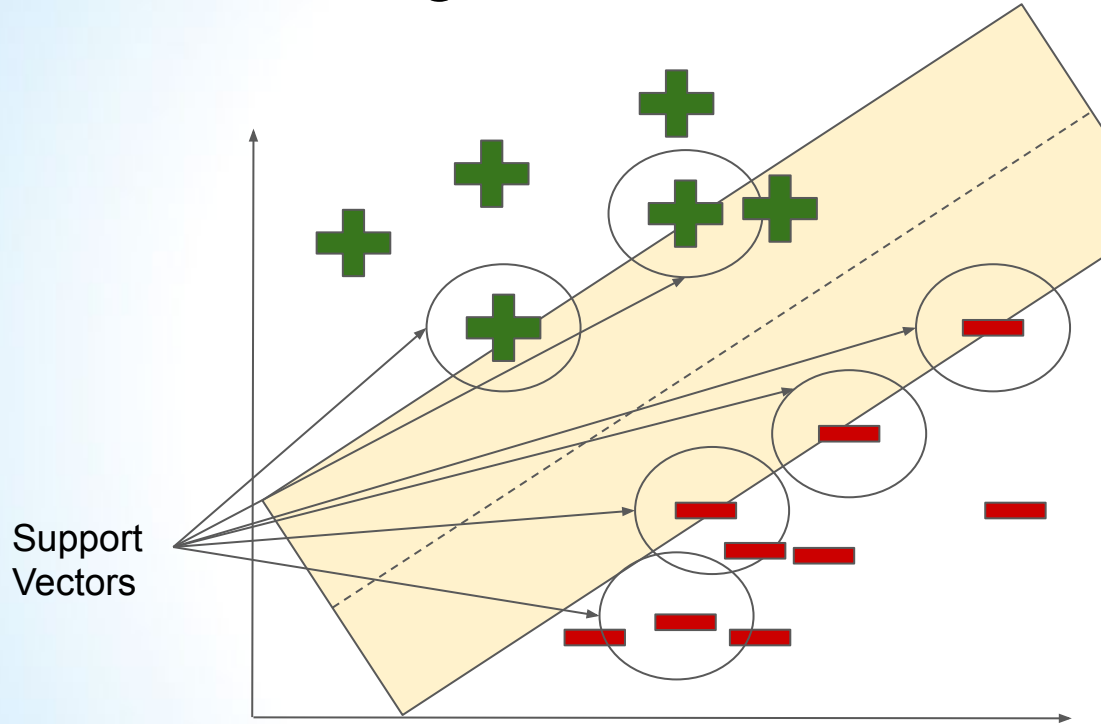
A margin in linear classifiers is the boundary width the touches the datapoint

Maximum margin



A maximum margin in linear classifiers is the Max boundary width the touches the datapoint

Maximum Margin

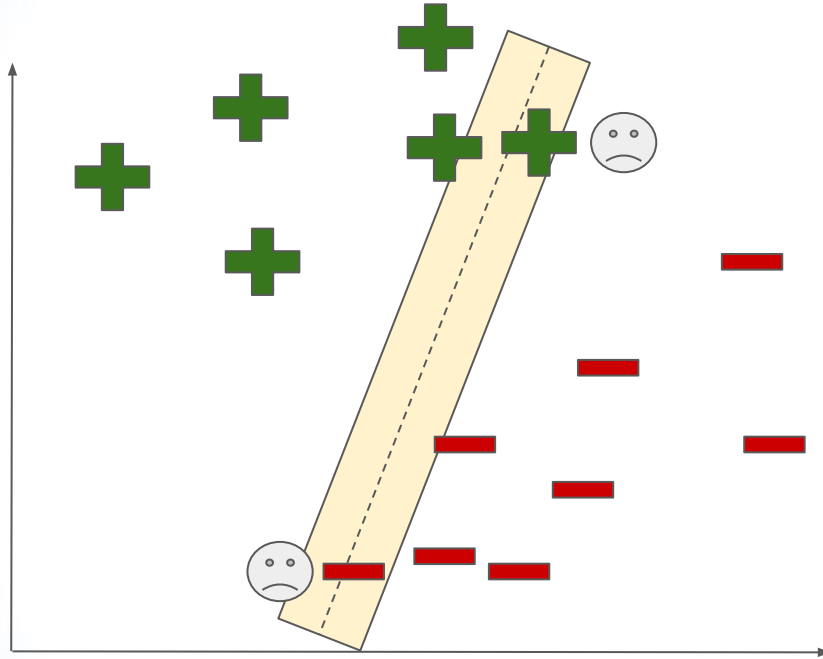


A maximum margin in linear classifiers is the Max boundary width the touches the datapoint

The points on the margins are called Support Vectors

***VC dimension* can show that the maximum margin is a good approach to linearly separable problems.**

Classifier margin



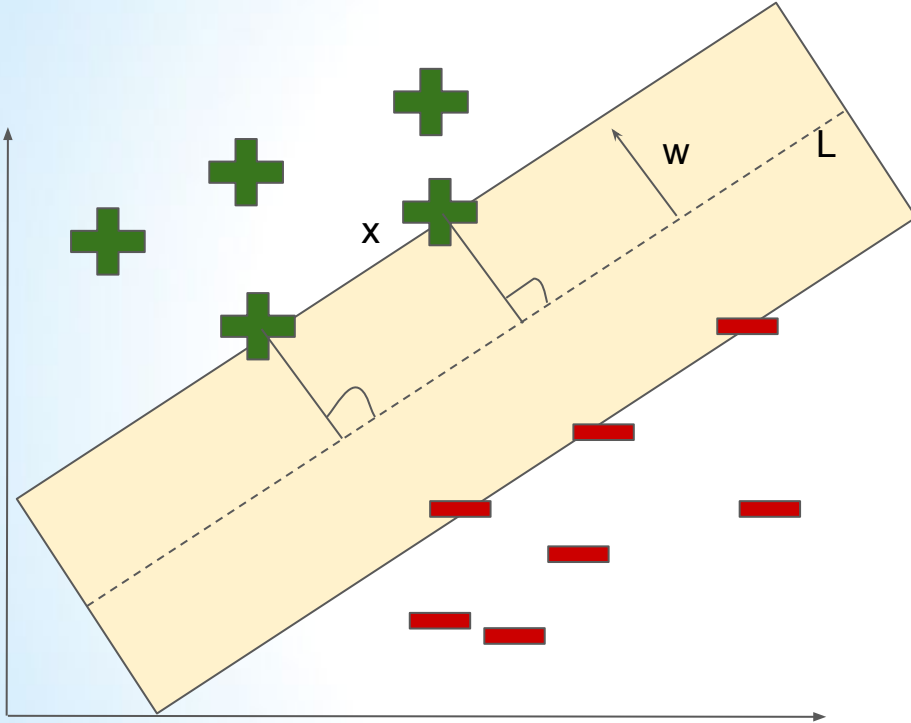
A maximum margin in linear classifiers is the Max boundary width the touches the datapoint

The points on the margins are called Support Vectors

VC dimension can show that the maximum margin is a good approach to linearly separable problems.

Allows a more flexibility around the decision boundary

Hard SVM

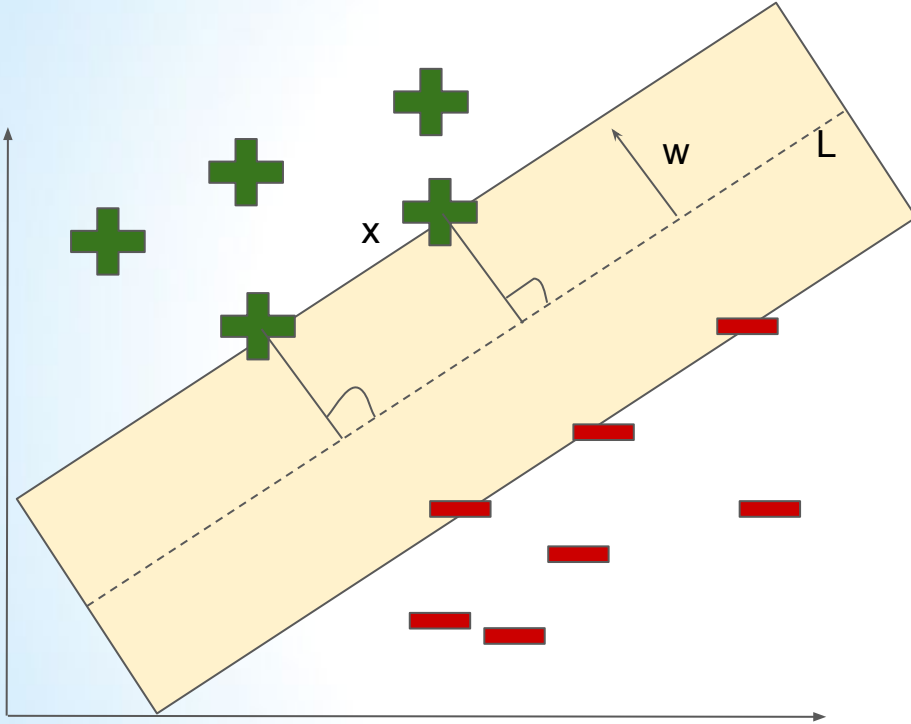


A distance of datapoint i can be defined by

$$\rho_i = y_i w x_i$$

y is the point label which
can be 1 or -1

Hard SVM



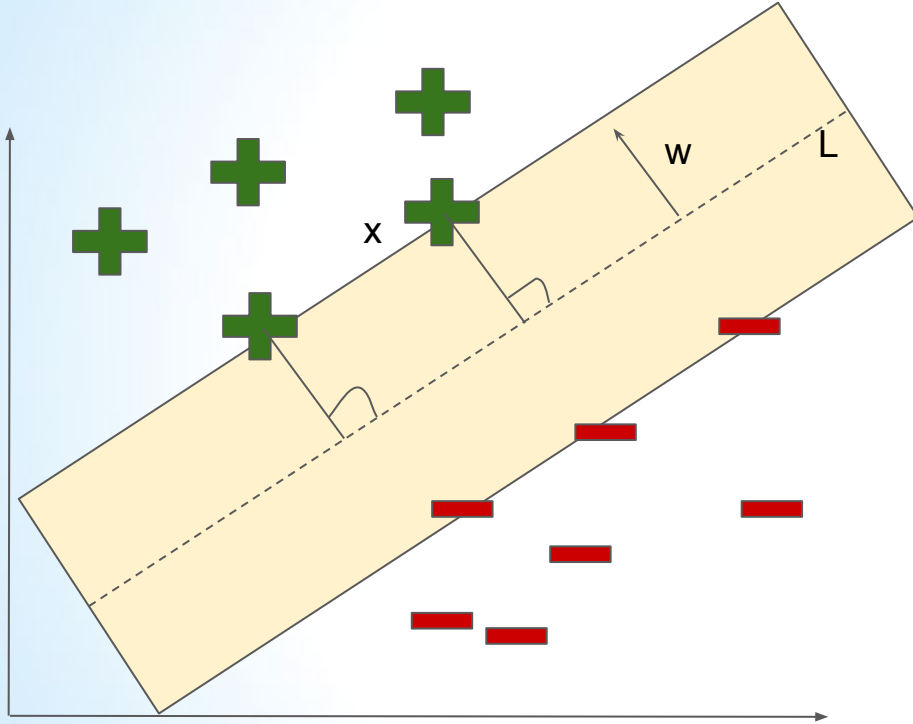
A distance of datapoint i can be defined by

$$\rho_i = y_i w x_i$$

The support vectors has the minimum margin

$$\rho = \min y_i w x_i$$

Hard SVM



A distance of datapoint i can be defined by

$$\rho_i = y_i w x_i$$

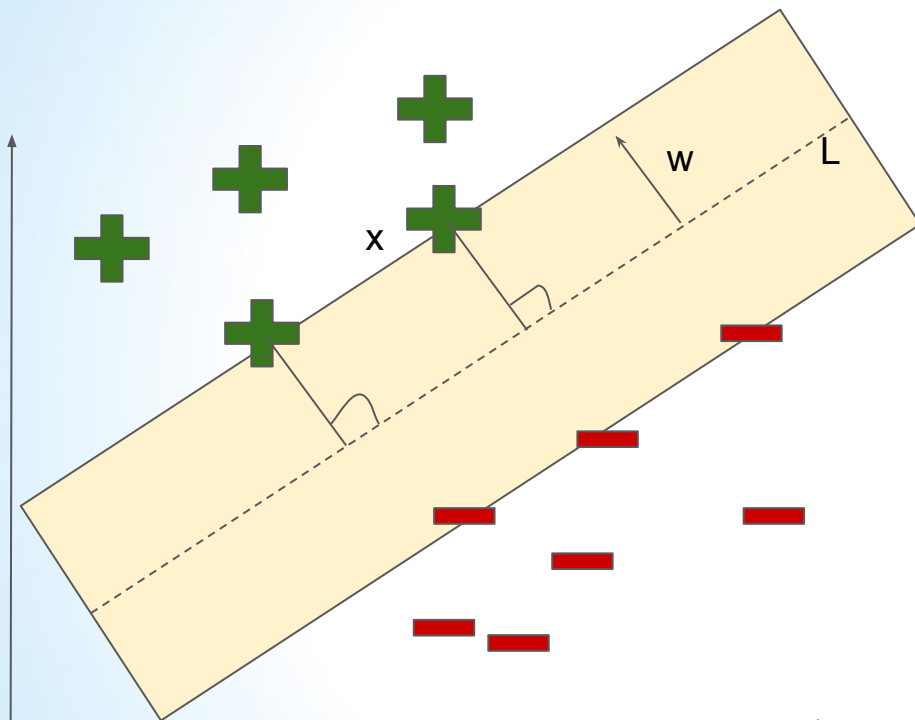
The support vectors has the minimum margin

$$\rho = \min y_i w x_i$$

We want to find the w that maximize the margin

$$\operatorname{argmax}_w \rho$$

Hard SVM



Add $\|w\|$ to constraint the w scaling
Define a margin that not depends on the scale of w

A distance of datapoint i can be defined by

$$\rho_i = \frac{y_i w x_i}{\|w\|}$$

The support vectors has the minimum margin

$$\rho = \min \frac{y_i w x_i}{\|w\|}$$

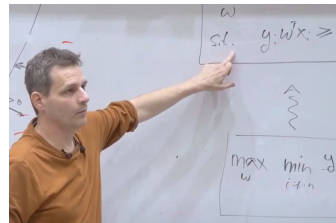
We want to find the w that maximize the margin

$$\operatorname{argmax}_w \rho$$

$$\operatorname{argmax}_w \min \frac{y_i w x_i}{\|w\|}$$

$$\operatorname{argmax}_w \min \frac{y_i w x_i}{\|w\|}$$

equivalent to

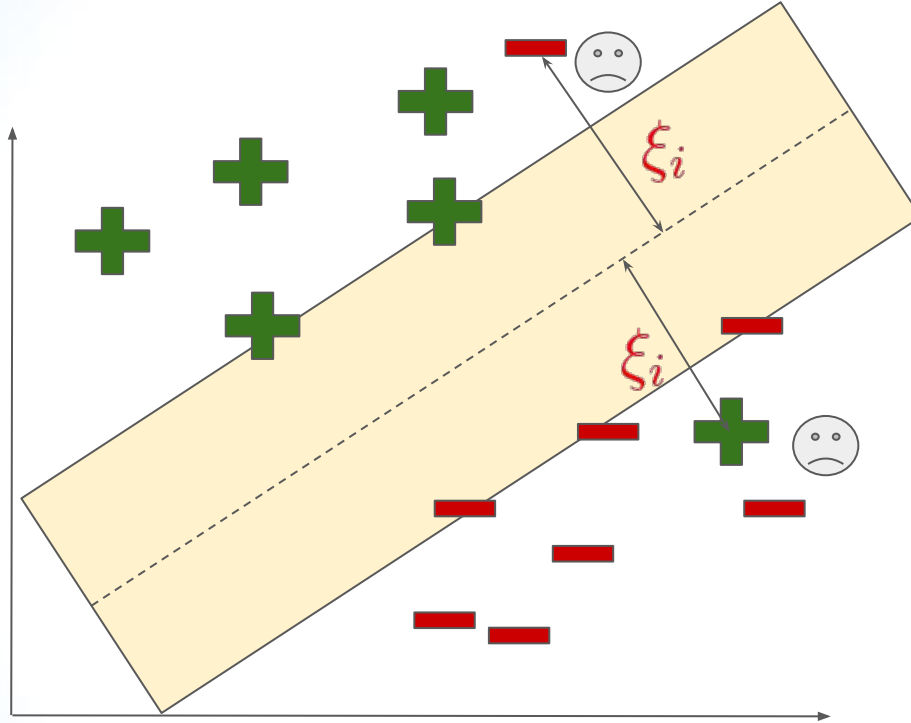


<https://youtu.be/LceLJvKMbBk?t=3613>

$$\min_w \|w\|^2 \quad s.t \quad \forall i, y_i w x_i \geq 1$$

*Assumes full separability (1.1 in HW)

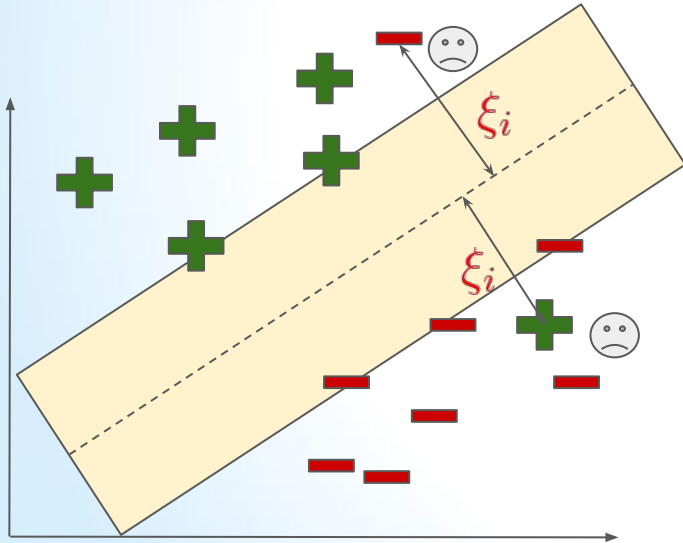
What if there is no linear separability?



Soft SVM - Supporting error data points

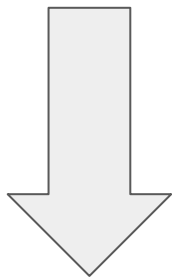
$$\operatorname{argmin}_{w, \xi} \left(\lambda \|w\|^2 + C \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$

$$s.t \ \forall y_i (wx_i) \geq 1 - \xi_i$$



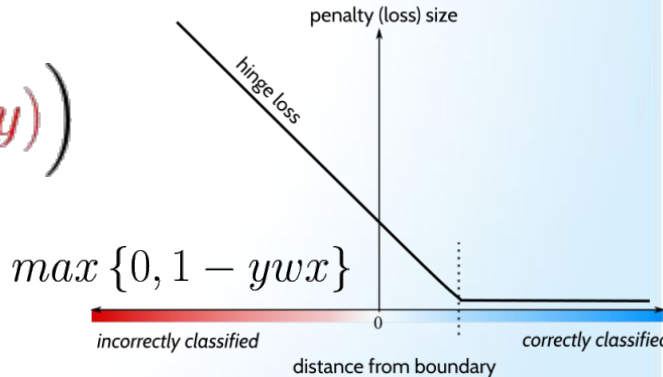
Soft SVM - Supporting error data points

$$\operatorname{argmin}_{w, \xi} \left(\lambda \|w\|^2 + C \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$



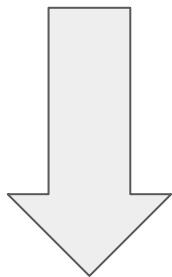
$$s.t \ \forall y_i (wx_i) \geq 1 - \xi_i$$

$$\operatorname{argmin}_{w, \xi} \left(\lambda \|w\|^2 + \text{hinge}(wx, y) \right)$$



Soft SVM - Supporting error data points

$$\operatorname{argmin}_{w, \xi} \left(\lambda \|w\|^2 + C \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$



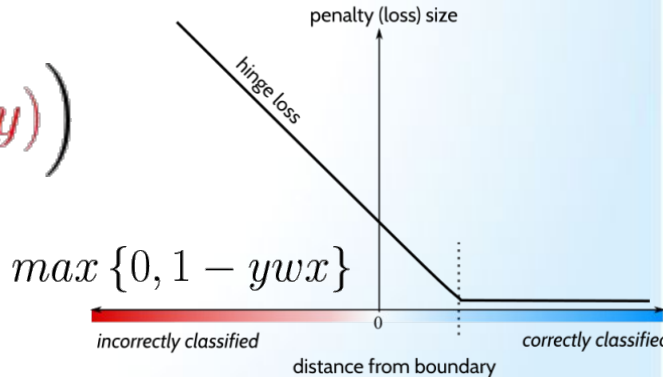
$$s.t \ \forall y_i (wx_i) \geq 1 - \xi_i$$

$$\operatorname{argmin}_{w, \xi} \left(\lambda \|w\|^2 + \text{hinge}(wx, y) \right)$$

regularization

objective

But came from problem definition



How we can solve it?

Gradient descent

Quadratic Programing

SGD for solving Soft-SVM

goal: Solve $\operatorname{argmin}_{\mathbf{w}} \left(\frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x}_i \rangle\} \right)$

parameter: T

initialize: $\boldsymbol{\theta}^{(1)} = \mathbf{0}$

for $t = 1, \dots, T$

Let $\mathbf{w}^{(t)} = \frac{1}{\lambda t} \boldsymbol{\theta}^{(t)}$

Choose i uniformly at random from $[m]$

If $(y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle < 1)$

Set $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + y_i \mathbf{x}_i$

Else

Set $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)}$

output: $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)}$

Homework

1. **Implementing SVM called Pegasos (2011) - 55 (+ 10 bonus)**
 1. Implement Class - 35
 2. Test - 10
 3. Analyze param - 5
 4. Analyze learning - 5
 5. Mini-batch bonus - 10*
2. **The effect of imbalance on SVM - 15**
3. **Practical SVM in scikit-learn & hypertune - 10**
4. **Using different Kernels - 20**

Another Way to Solve SVM

And other tricks

From Primal to Dual

Primal

$$\min f_0(x)$$

$$s.t \ f_i(x) \leq 0 \quad \forall i = 1..m$$

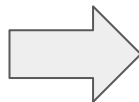
Original SVM definition

From Primal to Dual

Primal

$$\begin{aligned} \min f_0(x) \\ \text{s.t. } f_i(x) \leq 0 \quad \forall i = 1..m \end{aligned}$$

Original SVM definition



Dual

$$\begin{aligned} L(x, \alpha) &= f_0(x) + \sum_{i=1} \alpha_i f_i(x) \\ \alpha &\geq 0 \\ g(\alpha) &= \min_x L(x, \alpha) \end{aligned}$$

Dual definition that
solvable with
linear solvers

From Primal to Dual

Primal

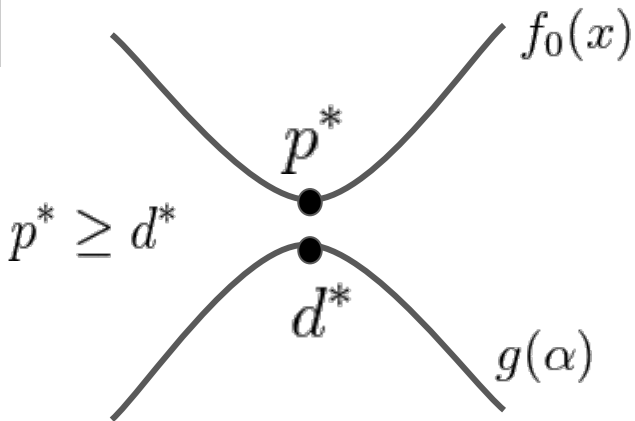
$$\begin{aligned} \min f_0(x) \\ \text{s.t. } f_i(x) \leq 0 \quad \forall i = 1..m \end{aligned}$$

Original SVM definition

Dual

$$\begin{aligned} L(x, \alpha) &= f_0(x) + \sum_{i=1} \alpha_i f_i(x) \\ g(\alpha) &= \min_x L(x, \alpha) \end{aligned}$$

Dual definition that
solvable with
linear solvers

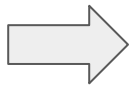


From Primal to Dual

Primal

$$\min f_0(x)$$

$$s.t \ f_i(x) \leq 0 \quad \forall i = 1..m$$



Dual

$$L(x, \alpha) = f_0(x) + \sum_{i=1}^m \alpha f_i(x)$$

$$g(x) = \min_x L(x, \alpha)$$

**SVM
definition**

$$\min \|w\|^2$$

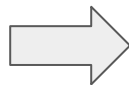
$$s.t \quad y_i w x_i \geq 1 \quad \forall i = 1..m$$

From Primal to Dual

Primal

$$\min f_0(x)$$

$$s.t. f_i(x) \leq 0 \quad \forall i = 1..m$$



Dual

$$L(x, \alpha) = f_0(x) + \sum_{i=1}^m \alpha f_i(x)$$

$$g(x) = \min_x L(x, \alpha)$$

**SVM
definition**

$$\min \|w\|^2$$

$$s.t. y_i w x_i \geq 1 \quad \forall i = 1..m$$

Modify

$$\min \frac{1}{2} \|w\|^2$$

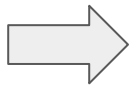
$$s.t. 1 - y_i w x_i \leq 0 \quad \forall i = 1..m$$

From Primal to Dual

Primal

$$\min f_0(x)$$

$$s.t. f_i(x) \leq 0 \quad \forall i = 1..m$$



Dual

$$L(x, \alpha) = f_0(x) + \sum_{i=1}^m \alpha f_i(x)$$

$$g(x) = \min_x L(x, \alpha)$$

**SVM
definition**

$$\min \|w\|^2$$

$$s.t. y_i w x_i \geq 1 \quad \forall i = 1..m$$

Modify

$$\min \frac{1}{2} \|w\|^2$$

$$s.t. 1 - y_i w x_i \leq 0 \quad \forall i = 1..m$$



$$L(x, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i w x_i - 1)$$

$$g(\alpha) = \min_x \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i w x_i - 1)$$

From Dual to Primal

$$L(x, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i w x_i - 1)$$

$$g(\alpha) = \min_x \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i w x_i - 1)$$

From Dual to Primal

$$L(x, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i w x_i - 1)$$

$$g(\alpha) = \min_x \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i w x_i - 1)$$

$$\frac{\partial L}{\partial w} = 0$$

From Dual to Primal

$$L(x, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i w x_i - 1)$$

$$g(\alpha) = \min_x \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i w x_i - 1)$$

$$\frac{\partial L}{\partial w} = 0$$

$$w - \sum_{i=1}^m \alpha_i y_i x_i = 0$$

$$w^* = \sum_{i=1}^m \alpha_i y_i x_i$$

$$\alpha_i \geq 0$$

Lagrange coefficients for each sample in the training set

**Optimization
Definition**

$$\max_{\alpha_k} \sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l y_k y_l x_k x_l$$

All possible
pairs in the
training set

Constraints

$$s.t \quad 0 \leq \alpha_k \leq C \quad \sum_{k=1}^R \alpha_k y_k = 0$$

**Optimization
Definition**

$$\max_{\alpha_k} \sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l y_k y_l x_k x_l$$

Constraints

$$s.t \quad 0 \leq \alpha_k \leq C \quad \sum_{k=1}^R \alpha_k y_k = 0$$

**Back to
Primal**

$$w = \sum_{k=1}^R \alpha_k y_k x_k$$

**Optimization
Definition**

$$\max_{\alpha_k} \sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l y_k y_l x_k x_l$$

Constraints

$$s.t \quad 0 \leq \alpha_k \leq C \quad \sum_{k=1}^R \alpha_k y_k = 0$$

**Back to
Primal**

$$w = \sum_{k=1}^R \alpha_k y_k x_k$$

Predict

$$f(x, w) = \text{sign}(w, x)$$

**Optimization
Definition**

$$\max_{\alpha_k} \sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l y_k y_l x_k x_l$$

Constraints

$$s.t \quad 0 \leq \alpha_k \leq C \quad \sum_{k=1}^R \alpha_k y_k = 0$$

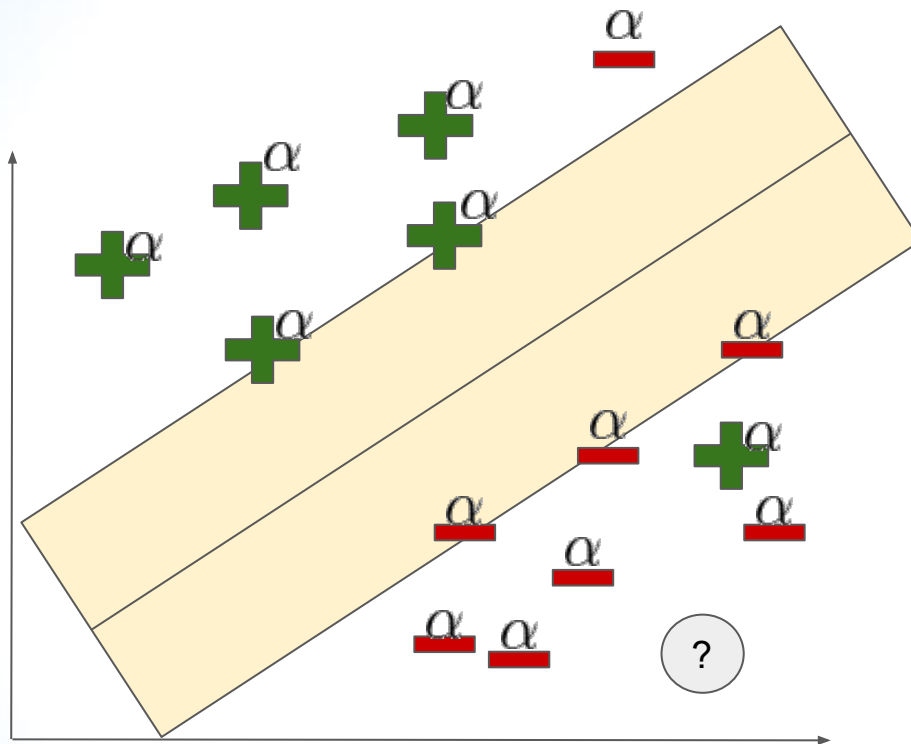
**Back to
Primal**

$$w = \sum_{k=1}^R \alpha_k y_k x_k$$

Predict

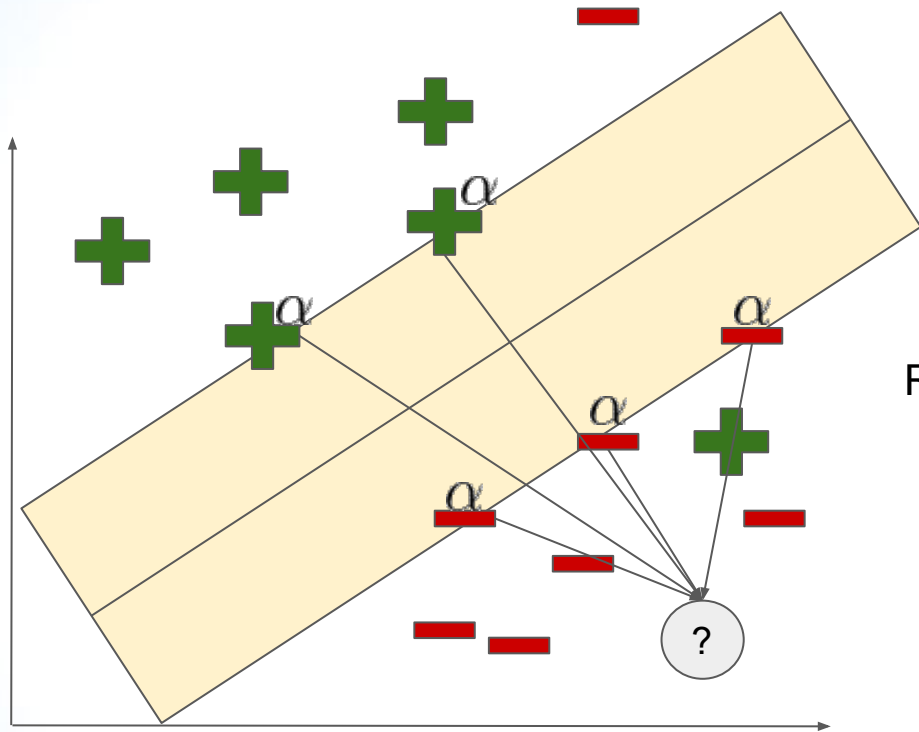
$$f(x, w) = \text{sign}(w, x) \Rightarrow f(x, w) = \text{sign}\left(\sum_{k=1}^R \alpha_k y_k x_k x\right)$$

Applying



$$f(x, w) = \text{sign}\left(\sum_{k=1}^R \alpha_k y_k x_k x\right)$$

Applying

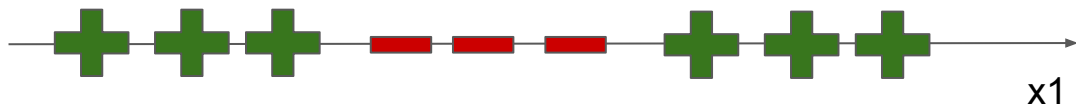


$$f(x, w) = \text{sign}\left(\sum_{k=1}^R \alpha_k y_k x_k x\right)$$

Reminds you something?

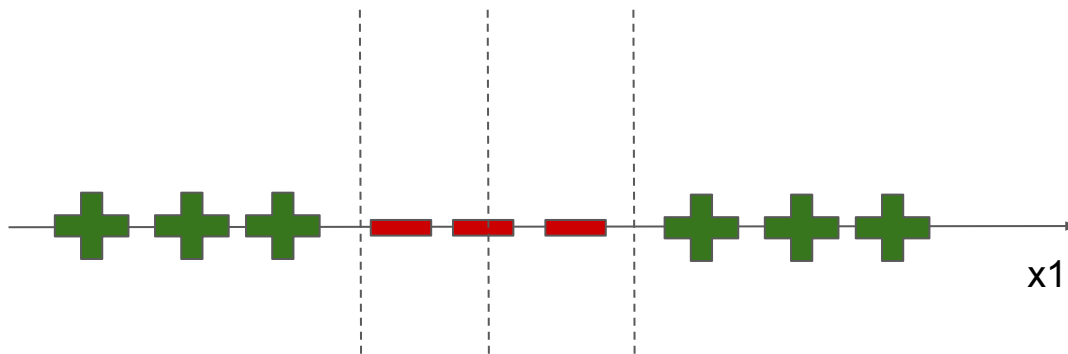
Dealing with non-linearity

How to separate with linear classifier?



Dealing with non-linearity

How to separate with linear classifier?

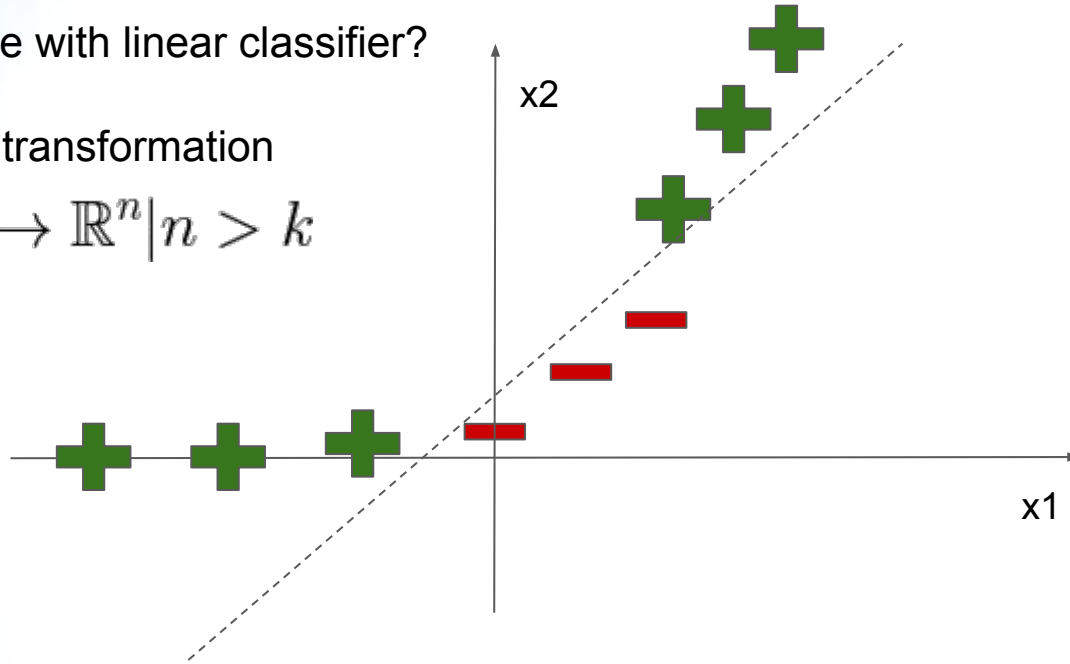


Dealing with non-linearity

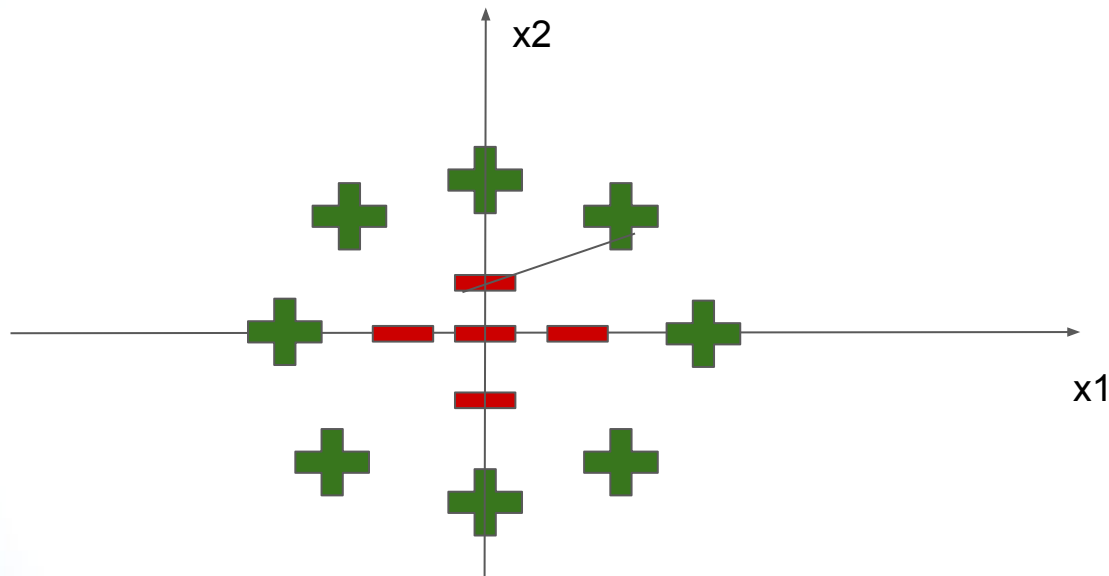
How to separate with linear classifier?

Use non-linear transformation

$$\phi(X) : \mathbb{R}^k \rightarrow \mathbb{R}^n | n > k$$

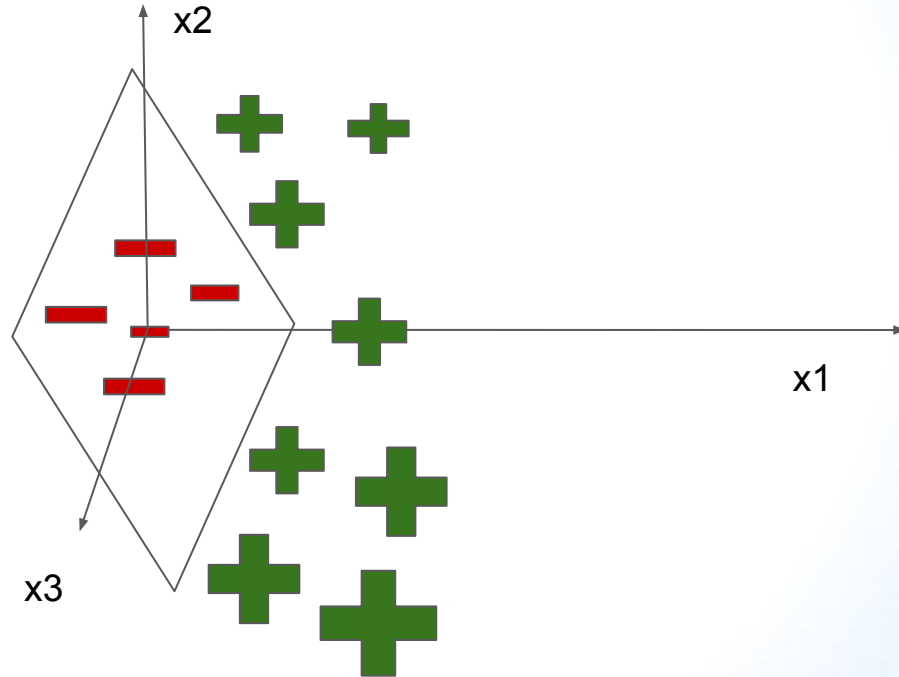


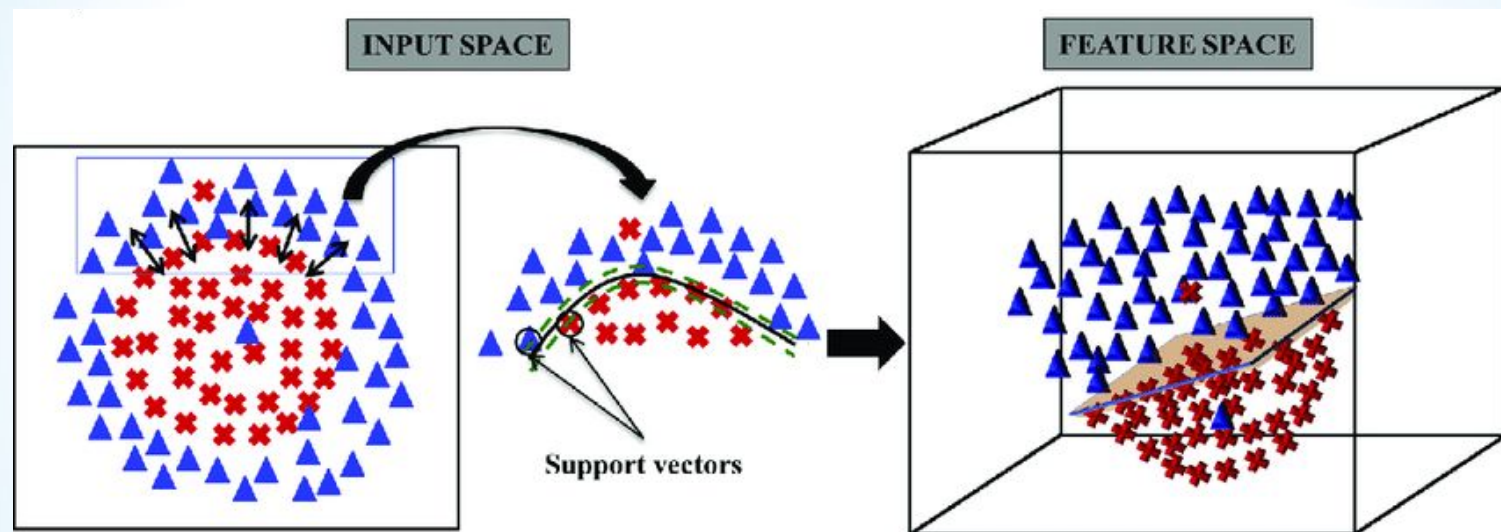
Dealing with non-linearity



Dealing with non-linearity

$$\phi(X) : \mathbb{R}^k \rightarrow \mathbb{R}^n | n > k$$





Dealing with non-linearity

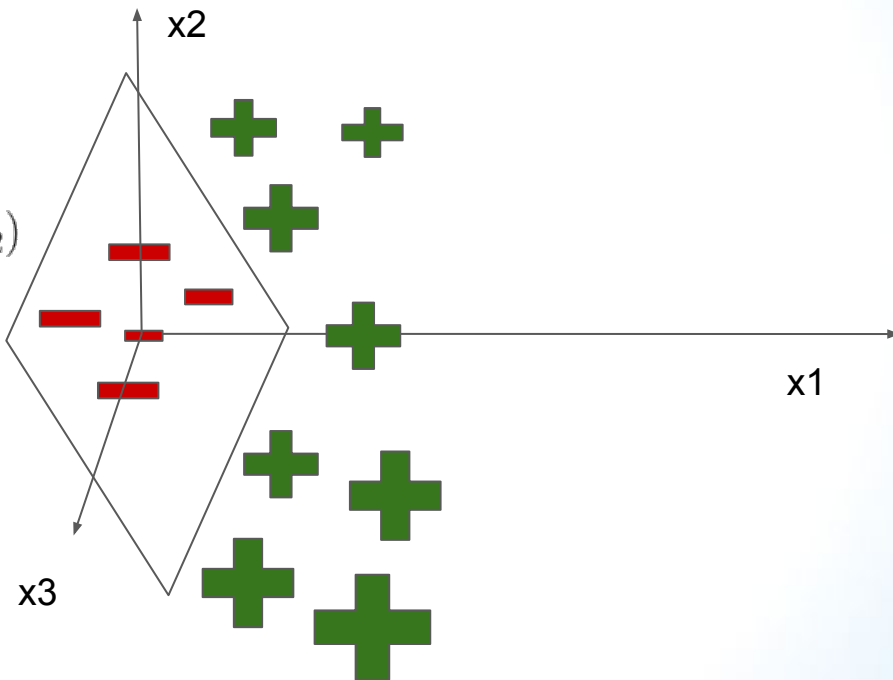
How to separate with linear classifier?

$$\phi(X) : \mathbb{R}^k \rightarrow \mathbb{R}^n | n > k$$

$$(x_1, x_2) \mapsto \theta(X) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$



A capital X = a datapoint



The Kernel Trick

To support this transformation We would need to compute all these features for each sample

The Kernel Trick

To support this transformation We would need to compute all these features for each sample

The Solution: the Kernel trick!

Use a dot product in feature space can be computed as kernel function

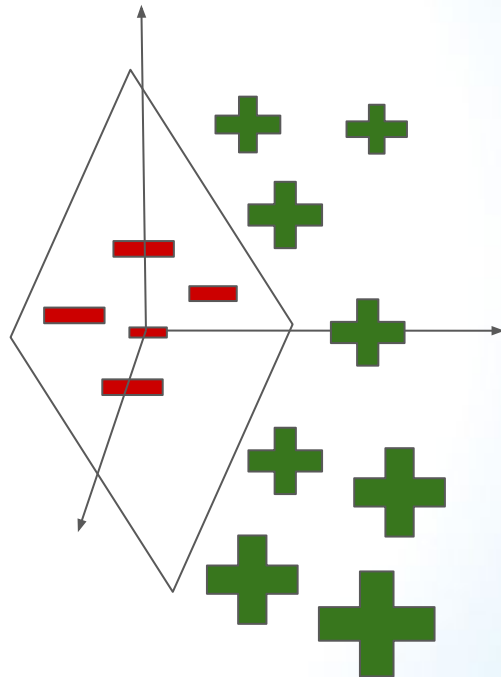
$$K(x_i, x_j) = \phi(x_i)\phi(x_j)$$

The Kernel Trick

$$(x_1, x_2) \mapsto \theta(X) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

$$K(X_i, X_j) = \phi(X_i)\phi(X_j) = (X_iX_j)^2$$

$$= (X_{i_1}^2, X_{i_2}^2, \sqrt{2}X_{i_1}X_{i_2})(X_{j_1}^2, X_{j_2}^2, \sqrt{2}X_{j_1}X_{j_2})^T = \\ (X_{i_1}X_{j_1} + X_{i_2}X_{j_2})^2 = (X_iX_j)^2$$



The Kernel Trick

We would need to compute all these features!

The Solution: the Kernel trick!

Use a dot product in feature space can be computed as kernel function

$$K(x_i, x_j) = \phi(x_i)\phi(x_j)$$

Where do we have dot product in SVM?

**Optimization
Definition**

$$\max_{\alpha_k} \sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l y_k y_l x_k x_l$$

Constraints

$$s.t \quad 0 \leq \alpha_k \leq C \quad \sum_{k=1}^R \alpha_k y_k = 0$$

**Back to
Primal**

$$w = \sum_{k=1}^R \alpha_k y_k x_k$$

Predict

$$f(x, w) = \text{sign}(w, x) \qquad f(x, w) = \text{sign}\left(\sum_{k=1}^R \alpha_k y_k x_k x\right)$$





**Optimization
Definition**

$$\max_{\alpha_k} \sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l y_k y_l x_k x_l$$

Constraints

$$s.t \quad 0 \leq \alpha_k \leq C \quad \sum_{k=1}^R \alpha_k y_k = 0$$

**Back to
Primal**

$$w = \sum_{k=1}^R \alpha_k y_k x_k$$

Predict

$$f(x, w) = \text{sign}(w, x) \quad f(x, w) = \text{sign}\left(\sum_{k=1}^R \alpha_k y_k x_k x\right)$$

**Optimization
Definition**

$$\max_{\alpha_k} \sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l y_k y_l K(x_k x_l)$$

Constraints

$$s.t \quad 0 \leq \alpha_k \leq C \quad \sum_{k=1}^R \alpha_k y_k = 0$$

**Back to
Primal**

$$w = \sum_{k=1}^R \alpha_k y_k x_k$$

Predict

$$f(x, w) = \text{sign}(w, x) \qquad f(x, w) = \text{sign}\left(\sum_{k=1}^R \alpha_k y_k x_k x\right)$$

**Optimization
Definition**

$$\max_{\alpha_k} \sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l y_k y_l K(x_k x_l)$$

Constraints

$$s.t \quad 0 \leq \alpha_k \leq C \quad \sum_{k=1}^R \alpha_k y_k = 0$$

We still need to compute all the kernels pairs, but we don't need to maintain the feature space

**Back to
Primal**

$$w = \sum_{k=1}^R \alpha_k y_k x_k$$

Predict

$$f(x, w) = \text{sign}(w, x) \qquad f(x, w) = \text{sign}\left(\sum_{k=1}^R \alpha_k y_k x_k x\right)$$

Good Kernel Functions for SVM

On top the polynomial kernel function there are more suitable ones:

Radial Basis Function (RBF):

$$K(X_i, X_j) = \exp \left(-\frac{(X_i - X_j)^2}{2\sigma^2} \right)$$

Tanh Function (nn):

$$K(X_i, X_j) = \tanh (\kappa X_i X_j - \delta)$$

Good Kernel Functions for SVM

On top the polynomial kernel function there are more suitable ones:

Radial Basis Function (RBF):

$$K(X_i, X_j) = \exp \left(-\frac{(X_i - X_j)^2}{2\sigma^2} \right)$$

Tanh Function (nn):

$$K(X_i, X_j) = \tanh (\kappa X_i X_j - \delta)$$

hypertune



SVM Pros & Cons

Good

- Learn many non linear pattern
- Pick "conservative" hypotheses, that are less likely to overfit the data,
- Good for Small datasets with though target pattern
- “Only” 2 params - C, Kernel Function

SVM Pros & Cons

Good

- Learn many non linear pattern
- Pick “conservative” hypotheses, that are less likely to overfit the data,
- Good for Small datasets with though target pattern
- “Only” 2 params - C, Kernel Function

Bad

- $O(n^2-3)$ runtime depends on C and the kernel (n number of datapoints)
- $O(n^2)$ memory to compute all the pairwise kernels - 5-10K datapoints
- different values for the Kernel prams

Q&A

VC Dimension - <https://www.youtube.com/watch?v=puDzy2XmR5c>

<https://winvector.github.io/margin/margin.pdf>

<https://youtu.be/LceLJvKMbBk?t=7311>

<https://youtu.be/fB47g3QM0sk?t=839>

Shalev-Shwartz, S., Singer, Y., Srebro, N., & Cotter, A. (2011). Pegasos: Primal estimated sub-gradient solver for svm. Mathematical programming, 127(1), 3-30. [[pdf](http://www.ee.oulu.fi/research/imag/courses/Vedaldi/ShalevSiSr07.pdf)]

Schölkopf, B., Williamson, R. C., Smola, A. J., Shawe-Taylor, J., & Platt, J. C. (2000). Support vector method for novelty detection. In Advances in neural information processing systems (pp. 582-588). [[pdf](http://papers.nips.cc/paper/1723-support-vector-method-for-novelty-detection.pdf)]

Livni, R., Crammer, K. & Globerson, A.. (2012). A Simple Geometric Interpretation of SVM using Stochastic Adversaries. Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics, in PMLR 22:722-730. [[pdf](http://proceedings.mlr.press/v22/livni12/livni12.pdf)]

From Primal to Dual

$$\begin{aligned} &\text{minimize} && (1/2)x^T P_0 x + q_0^T x + r_0 \\ &\text{subject to} && (1/2)x^T P_i x + q_i^T x + r_i \leq 0, \quad i = 1, 2, \dots, m \\ &&& Ax = b \end{aligned}$$

From Primal to Dual

$$\begin{aligned} &\text{minimize} && (1/2)x^T P_0 x + q_0^T x + r_0 \\ &\text{subject to} && (1/2)x^T P_i x + q_i^T x + r_i \leq 0, \quad i = 1, 2, \dots, m \\ &&& Ax = b \end{aligned}$$

Works in convex problems
Transform the problem from
minimize to maximize
Use lagrange coefficients

$$\begin{aligned} &\text{maximize} && -(1/2)q(\lambda)^T P(\lambda)^{-1} q(\lambda) + r(\lambda) \\ &\text{subject to} && \lambda \succeq 0 \end{aligned}$$

where:

$$P(\lambda) = P_0 + \sum_{i=1}^m \lambda_i P_i, \quad q(\lambda) = q_0 + \sum_{i=1}^m \lambda_i q_i, \quad r(\lambda) = r_0 + \sum_{i=1}^m \lambda_i r_i$$

From Primal to Dual

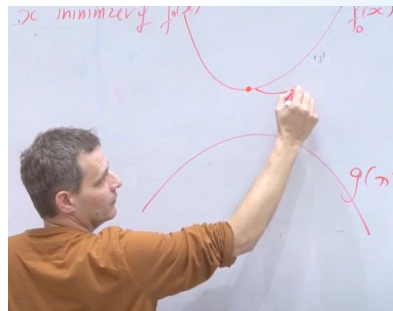
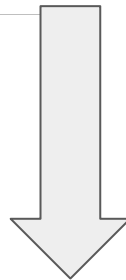
$$\begin{aligned} &\text{minimize} && (1/2)x^T P_0 x + q_0^T x + r_0 \\ &\text{subject to} && (1/2)x^T P_i x + q_i^T x + r_i \leq 0, \quad i = 1, 2, \dots, m \\ &&& Ax = b \end{aligned}$$

Works in convex problems
Transform the problem from
minimize to maximize
Use lagrange coefficients

$$\begin{aligned} &\text{maximize} && -(1/2)q(\lambda)^T P(\lambda)^{-1} q(\lambda) + r(\lambda) \\ &\text{subject to} && \lambda \succeq 0 \end{aligned}$$

where:

$$P(\lambda) = P_0 + \sum_{i=1}^m \lambda_i P_i, \quad q(\lambda) = q_0 + \sum_{i=1}^m \lambda_i q_i, \quad r(\lambda) = r_0 + \sum_{i=1}^m \lambda_i r_i$$



<https://youtu.be/LceLJvKMbBk?t=7311>
<https://youtu.be/fB47g3QM0sk?t=839>