

INTRODUCTION TO
MACHINE LEARNING
COMPSCI 4ML3

LECTURE 15

HASSAN ASHTIANI

LINEARLY SEPARABLE DATA

- A BINARY CLASSIFICATION DATA SET $Z = \{(x^i, y^i)\}_{i=1}^n$ IS LINEARLY SEPARABLE IF
 - THERE EXISTS W^* SUCH THAT
 - FOR EVERY $i \in [n]$ WE HAVE $\text{sgn}(\langle x^i, W^* \rangle) = y^i$
 - OR EQUIVALENTLY, FOR EVERY $i \in [n]$ WE HAVE $(W^{*T} x^i) y^i > 0$
- IN OTHER WORDS, THE CLASSIFICATION ERROR ON Z IS 0
- CAN WE FIND W^* EFFICIENTLY FOR LINEARLY SEPARABLE DATA?

LP FOR LINEAR CLASSIFICATION

- DEFINE $A = [x_j^i y^i]_{n \times d}$
- THEN FINDING THE OPTIMAL W IS EQUIVALENT TO

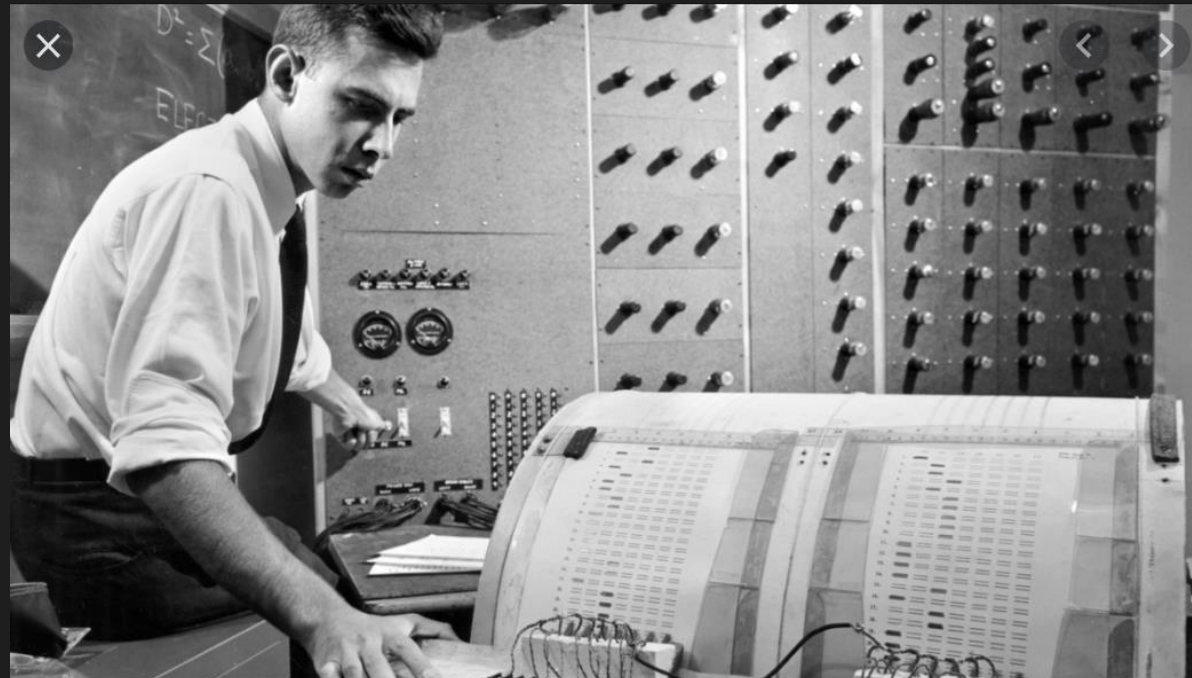
$$\max_{w \in \mathbb{R}^d} \langle \vec{0}, w \rangle$$

$$s. t. Aw \geq \vec{1}$$

WE CAN USE OFF-THE-SHELF LP SOLVERS!

APPROACH 2: PERCEPTRON

- PROPOSED IN 50'S BY ROSENBLATT
- PREDECESSOR OF NEURAL NETWORKS
 - MULTI-LAYER PERCEPTRON!



ROSENBLATT'S PERCEPTRON

Batch Perceptron

input: A training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

initialize: $\mathbf{w}^{(1)} = (0, \dots, 0)$

for $t = 1, 2, \dots$

if $(\exists i \text{ s.t. } y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle \leq 0)$ then

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + y_i \mathbf{x}_i$$

else

output $\mathbf{w}^{(t)}$

- IN EACH UPDATE, W BECOMES “MORE CORRECT” ON x^i
- [HTTPS://PHIRESKY.GITHUB.IO/KOGSYS-DEMOS/NEURAL-NETWORK-DEMO/?PRESET=ROSENBLATT+PERCEPTRON](https://phiresky.github.io/kogsys-demos/neural-network-demo/?preset=rosenblatt+perceptron)

THE GREEDY UPDATE

- IN EACH UPDATE, W BECOMES “MORE CORRECT” ON x^i :
- WHAT ABOUT OTHER x^j 'S?

NOVIKOFF, 1962

UNCLASSIFIED

AD 298 258

*Reproduced
by the*

**ARMED SERVICES TECHNICAL INFORMATION AGENCY
ARLINGTON HALL STATION
ARLINGTON 12, VIRGINIA**

Technical Report

ON CONVERGENCE PROOFS FOR PERCEPTRONS

Prepared for:

OFFICE OF NAVAL RESEARCH
WASHINGTON, D.C.

CONTRACT Nonr 3438(00)

By: Albert B. J. Novikoff

CATALOGUE BY
AS AD NO.---

STANFORD RESEARCH INSTITUTE

MENLO PARK, CALIFORNIA

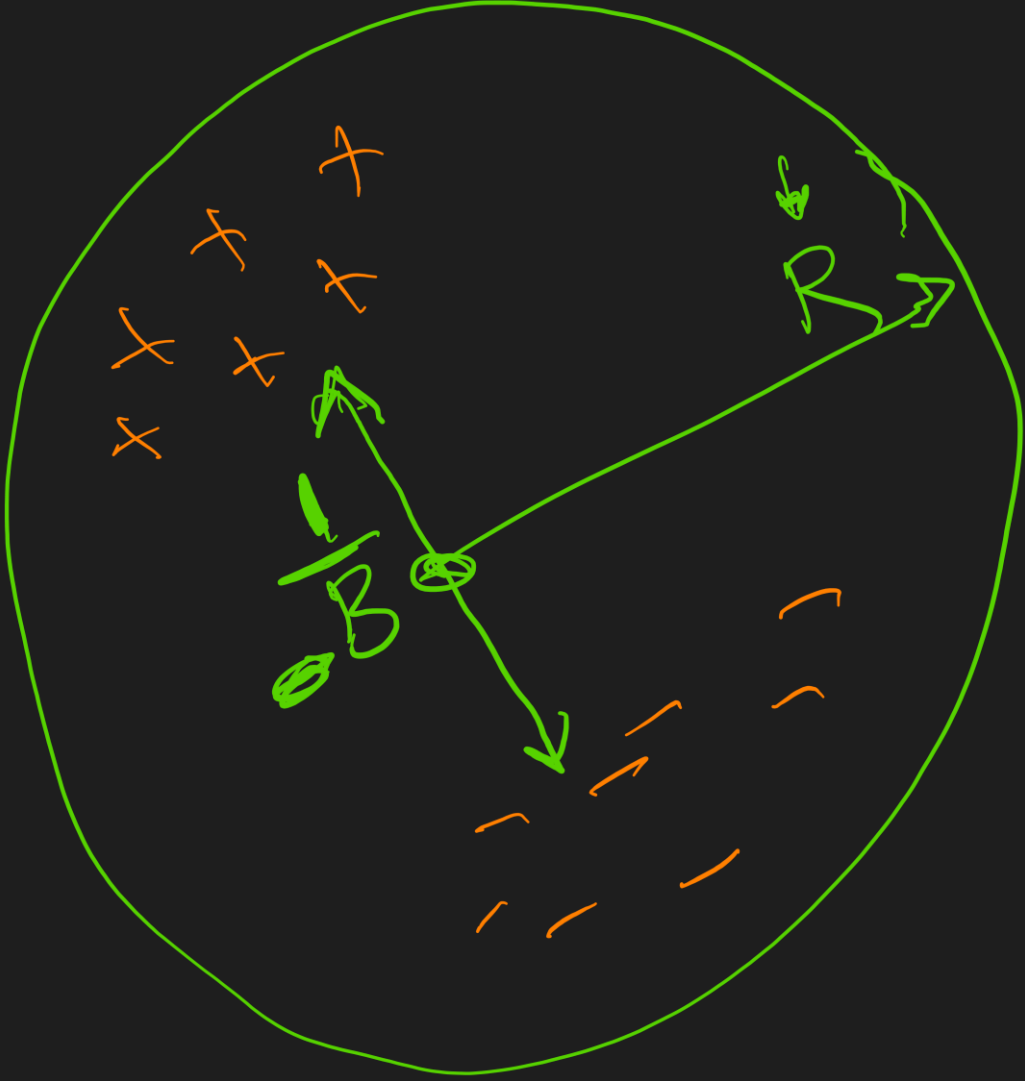


Novikoff, A. B. J. (1962). On convergence proofs on perceptrons. In *Proceedings of the Symposium on the Mathematical Theory of Automata*, Volume 12, pp. 615–622. Polytechnic Institute of Brooklyn.

CONVERGENCE OF PERCEPTRON

THEOREM 9.1 Assume that $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ is separable, let $B = \min\{\|\mathbf{w}\| : \forall i \in [m], y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1\}$, and let $R = \max_i \|\mathbf{x}_i\|$. Then, the Perceptron algorithm stops after at most $(RB)^2$ iterations, and when it stops it holds that $\forall i \in [m], y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle > 0$.

- #STEPS DOES NOT EXPLICITLY DEPEND ON d
- YOU CAN FIND MORE DETAILS ABOUT THIS LECTURE IN
 - UNDERSTANDING MACHINE LEARNING, CHAPTER 9
 - [HTTPS://WWW.CS.HUJI.AC.IL/~SHAIS/UNDERSTANDINGMACHINELEARNING/UNDERSTANDING-MACHINE-LEARNING-THEORY-ALGORITHMS.PDF](https://www.cs.huji.ac.il/~shais/understandingmachinelearning/understanding-machine-learning-theory-algorithms.pdf)

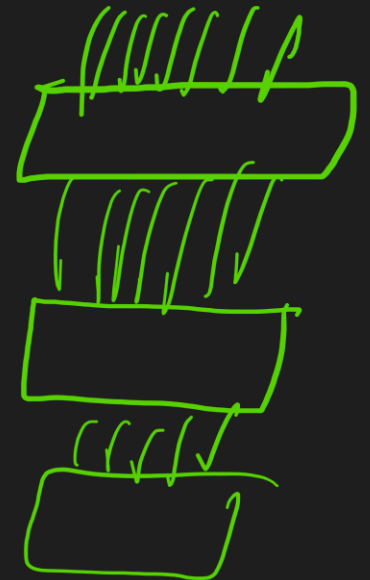


- IN 1969, MARVIN MINSKY AND SEYMOUR PAPERT ARGUED THAT IT IS IMPOSSIBLE TO LEARN XOR FUNCTION USING MULTILAYER PERCEPTRON...

- ONLY GOOD FOR LINEARLY SEPARABLE DATA

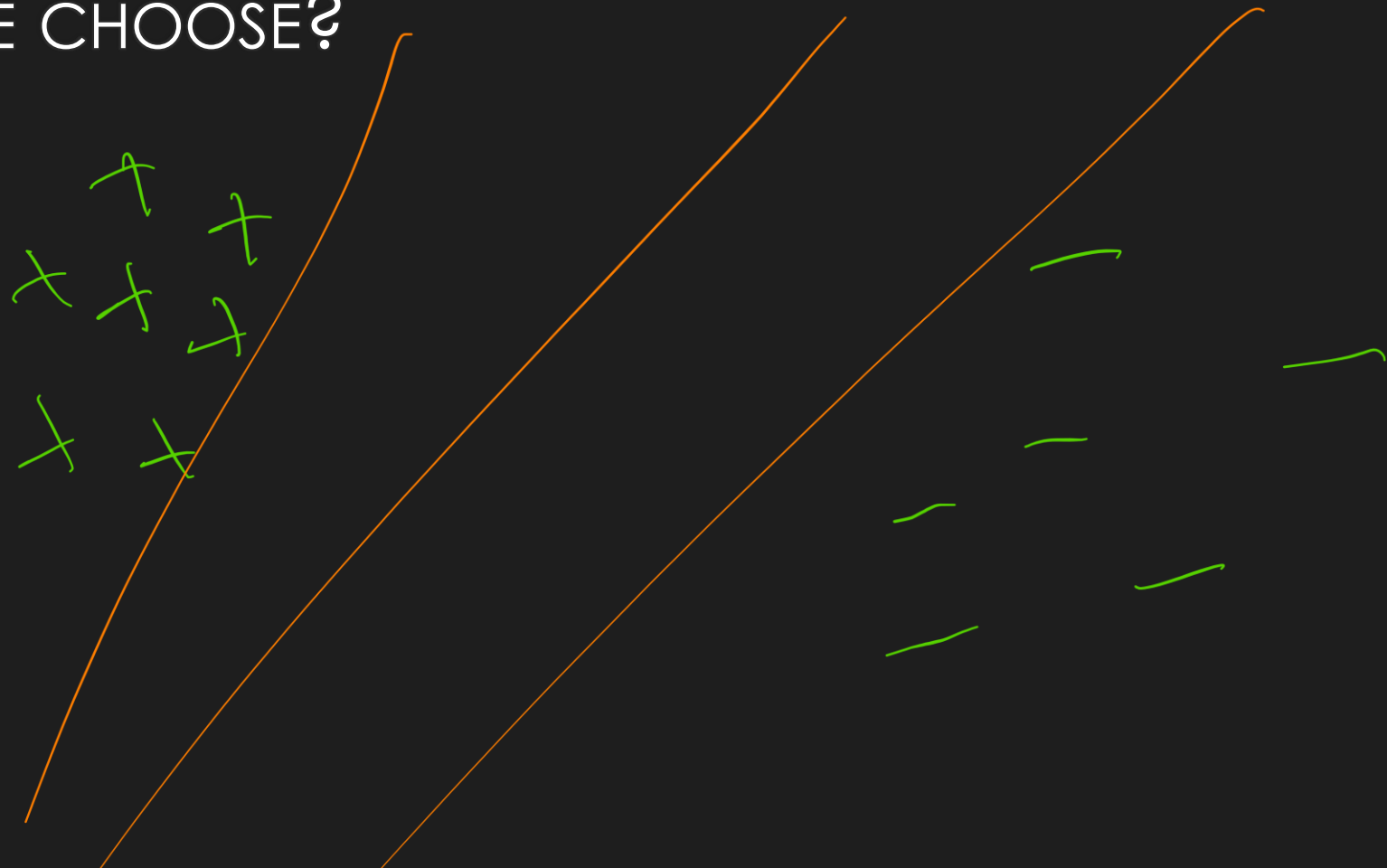
- STACKING PERCEPTRONS?

- 70's: AI (CONNECTIONISM) WINTER



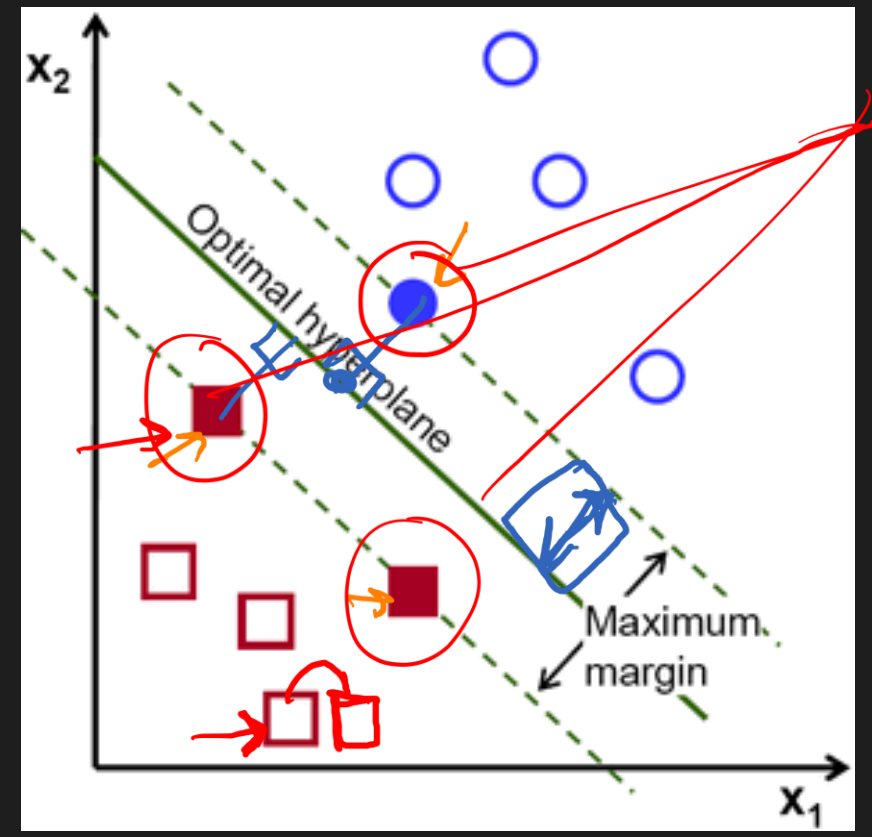
SUPPORT VECTOR MACHINES

- AMONG PERFECT LINEAR SEPARATORS, WHICH ONE SHOULD WE CHOOSE?



SUPPORT VECTOR MACHINES

- PICK THE LINEAR SEPARATOR THAT MAXIMIZES THE "MARGIN"
- MORE ROBUST TO "PERTURBATION"
- LESS PRONE TO OVERFITTING
 - WORKS WELL FOR HIGH-DIMENSIONAL DATA (?)
 - MORE ON THAT LATER!



DISTANCE OF A POINT TO A HYPERPLANE

THE EUCLIDEAN DISTANCE BETWEEN A POINT x AND THE HYPERPLANE PARAMETRIZED BY W IS (WHY?)

$$\rightarrow \frac{|W^T x + b|}{\|W\|_2}$$

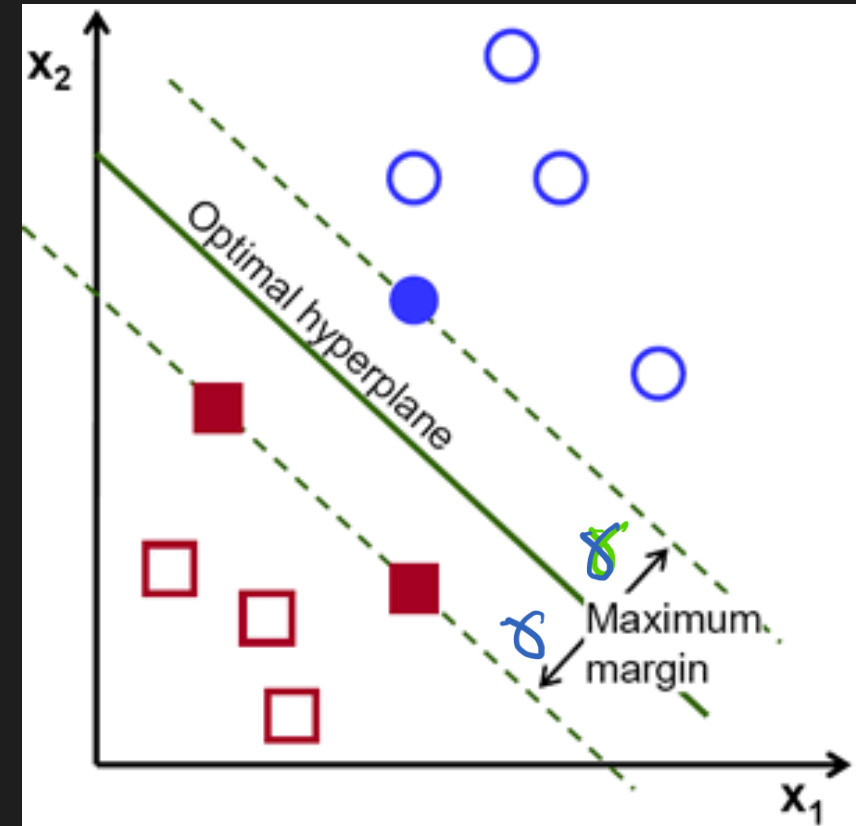


- THE DECISION BOUNDARY OF A LINEAR CLASSIFIER IS DETERMINED BY THE DIRECTION OF W (NOT $\|W\|_2$)
- ASSUME $\|W\|_2=1$, THEN THE DISTANCE IS

$$|W^T x + b| \leftarrow$$

MAXIMUM MARGIN HYPERPLANE

- LET THE HYPERPLANE BE PARAMETRIZED BY W and b
- ASSUME $\|W\|_2 = 1$
- W HAS A γ MARGIN IF
 - $W^T x + b \geq \gamma$ FOR EVERY BLUE x , AND
 - $W^T x + b \leq -\gamma$ FOR EVERY RED x



THE MARGIN

For simplicity, assume $b = 0$
 $y \in \{\pm 1\}$

- $Z = \{(x^i, y^i)\}_{i=1}^n, y \in \{-1, +1\}, \|W\|_2 = 1$

$$\text{Margin}(Z, W) = \min_{(x, y) \in Z} w^T x y =$$

$$= \text{Max } \delta$$

s.t. $\forall (x, y) \in Z, w^T x y \geq \delta$

$\text{Margin}(Z, W) < 0 \rightarrow Z$ is not linearly separable.

MAXIMIZING THE MARGIN

$$\begin{aligned} \text{Max Margin}(Z, w) &= \text{Max } \delta \\ w, \|w\|=1 & \text{ s.t. } \exists w, \left\{ \begin{array}{l} \forall (x, y) \in Z \\ \|w\|=1, \left\{ \begin{array}{l} w^T x y \geq \delta \end{array} \right\} \end{array} \right. \end{aligned}$$

$$\begin{aligned} &= \text{Max } \delta \\ \text{s.t. } \exists w, \|w\|=1, \forall (x, y) \in Z, \frac{1}{\delta} w^T x y &\geq 1 \end{aligned}$$

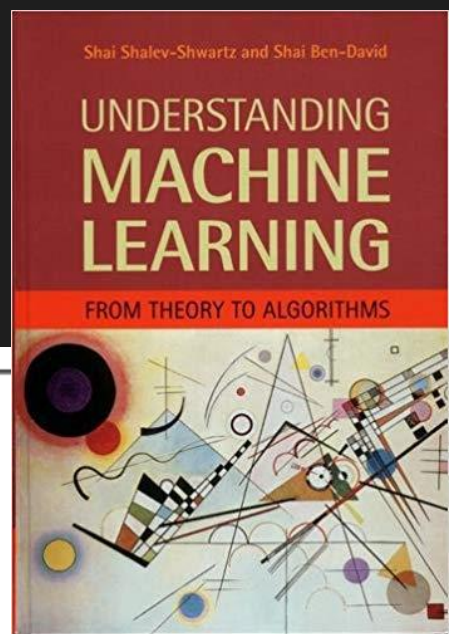
$$\begin{aligned} &= \text{Max } \frac{1}{\|w\|} \\ \text{s.t. } \forall (x, y) \in Z, w^T x y &\geq 1 \end{aligned}$$

so instead $\text{Min}_w \|w\|_2^2$

$$\text{s.t. } \forall (x, y) \in Z, w^T x y \geq 1$$

$$w^{new} = \frac{w^{old}}{\|w\|}$$

THE VERSION WITH "BIAS"



Hard-SVM

input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

solve:

$$(\mathbf{w}_0, b_0) = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|^2 \text{ s.t. } \forall i, y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \quad (15.2)$$

$$\text{output: } \hat{\mathbf{w}} = \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|}, \quad \hat{b} = \frac{b_0}{\|\mathbf{w}_0\|}$$

- WE COULD HAVE ALSO ADDED A DUMMY "1" FEATURE TO ALL POINTS SO AS TO ACCOUNT FOR THE BIAS/INTERCEPT

SENSITIVITY TO OUTLIERS

