

Review: ① Define machine learning,

② a) what are features?

b) what is a feature-vector?

③ what is a classifier?

④ what sort of tradeoffs do we have to make when choosing a classifier?

2-1

Lecture 2 - Linear Classifiers and the Perceptron Algorithm

Today:

Classification problem setup

• Classification problem setup

Training set: $S_n = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$

• Linear classifiers

- Definition

$x^{(i)}$ = feature vector for i^{th} data point

- Examples

$y^{(i)}$ = label for i^{th} data point

• Perceptron algorithm

- Algorithm statement

- Examples

- Perceptron with offset

- convergence and realizable data
Average perceptron

$$x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_d^{(i)} \end{bmatrix}$$

$\in \mathbb{R}^d$

$y^{(i)} \in \{-1, +1\}$

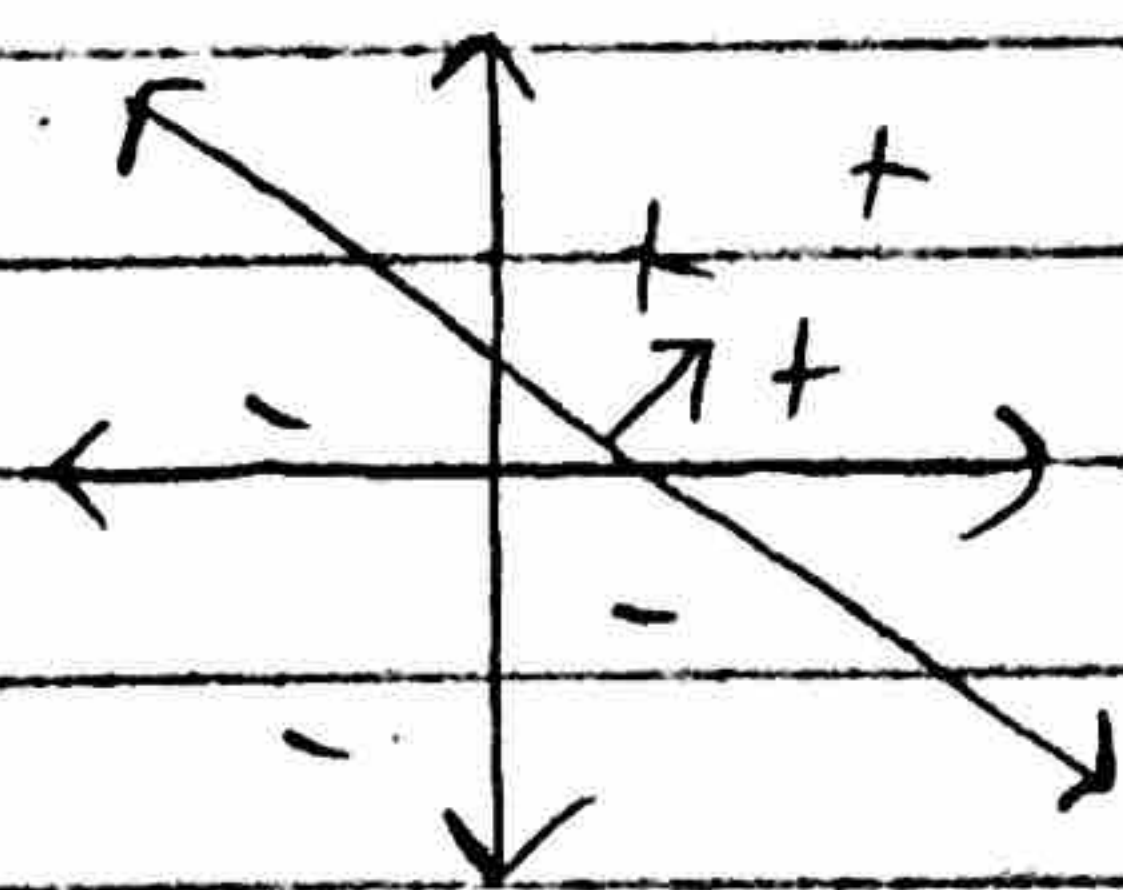
Classifier: $h: \mathbb{R}^d \rightarrow \{-1, +1\}$

$$h(x^{(i)}) = \pm 1$$

$h(x^{(i)}) = y^{(i)} \rightarrow$ correct prediction

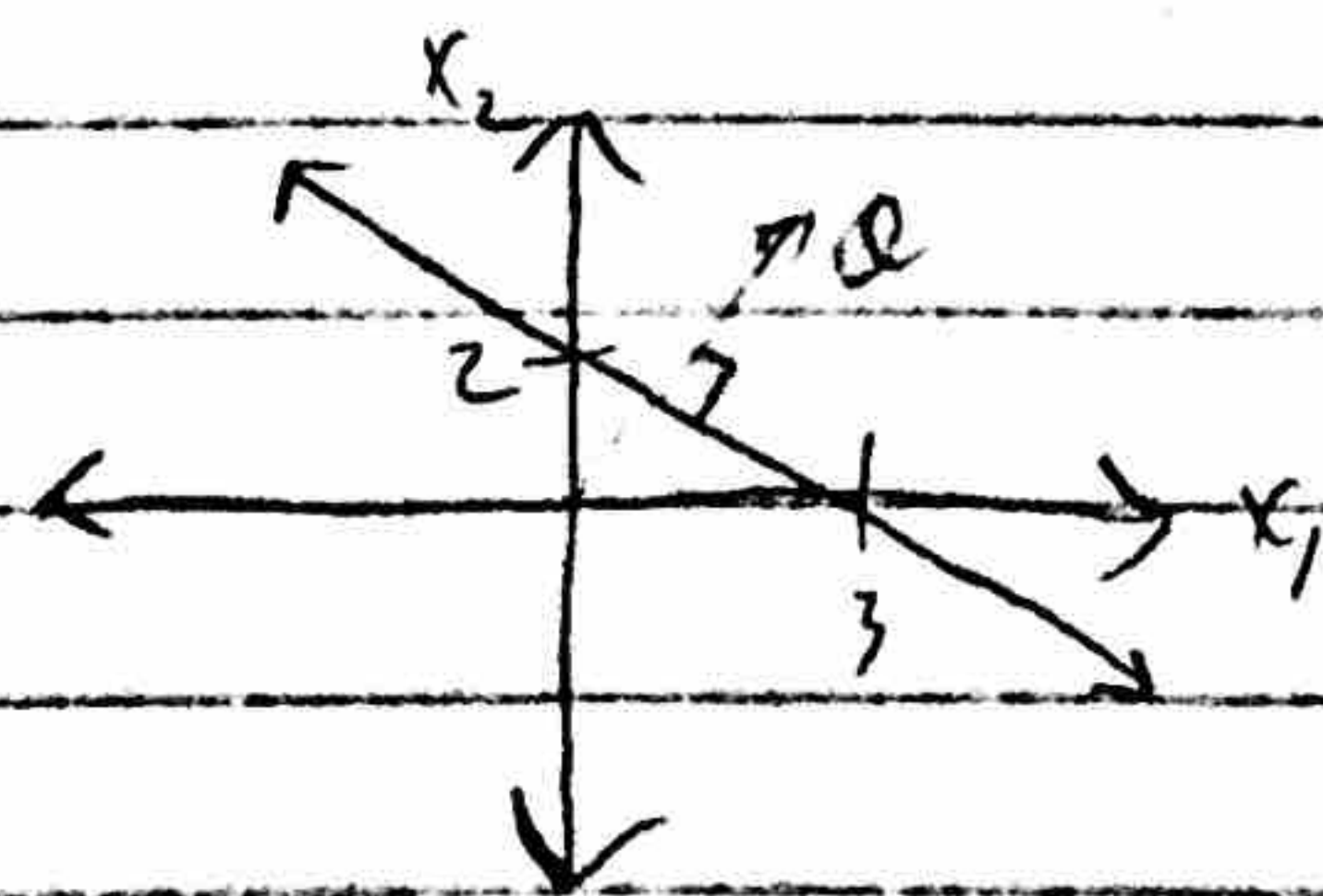
$h(x^{(i)}) \neq y^{(i)} \rightarrow$ incorrect prediction

Linear classifier (separator)



(Arrow points in direction of +)

Equation of a line



$$2x_1 + 3x_2 - 6 = 0$$

$$(x_2 = -\frac{2}{3}x_1 + 2)$$

$$\begin{pmatrix} 2 \\ 3 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + (-6) = 0$$

$$\theta \cdot x + \theta_0 = 0$$

Note: θ is perpendicular to (normal to) the line.

① What is the equation of the line

I draw?

② What is dot product?

③ Draw θ .

Example:

$d=2$



$$-3x_1 + 2x_2 + 2 = 0$$

$$x_2 = 3x_1 - 2$$

(?) Compute $\text{sign}(w \cdot x + b_0)$ and classify. Draw the point and confirm that it matches your expectations

1) $x = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \rightarrow +$ 2) $x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

3) $x = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \rightarrow -$ 4) $x = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \rightarrow -$

In general, a hyperplane in d -dimensions can be written as

$$w \cdot x + b_0 = 0$$

$$w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^d \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \in \mathbb{R}^d \quad b_0 \in \mathbb{R}$$

In previous example, we had $w = \begin{pmatrix} -3 \\ 2 \end{pmatrix} \in \mathbb{R}^2$ and $x = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \in \mathbb{R}^2$ and $b_0 = 2 \in \mathbb{R}$

If x is on the line (hyperplane): $w \cdot x + b_0 = 0$

If x is in the direction of w : $w \cdot x + b_0 > 0$

If x is in the opposite direction of w : $w \cdot x + b_0 < 0$

Definition:

$$\text{Def: } \text{sign}(w \cdot x + b_0) = \begin{cases} +1, & w \cdot x + b_0 > 0 \\ -1, & w \cdot x + b_0 \leq 0 \end{cases}$$

The sign function causes the line (hyperplane) to act as a decision boundary.

A linear classifier classifies a data point based on which side of the decision boundary it lies on. (Remember: w points towards positive)

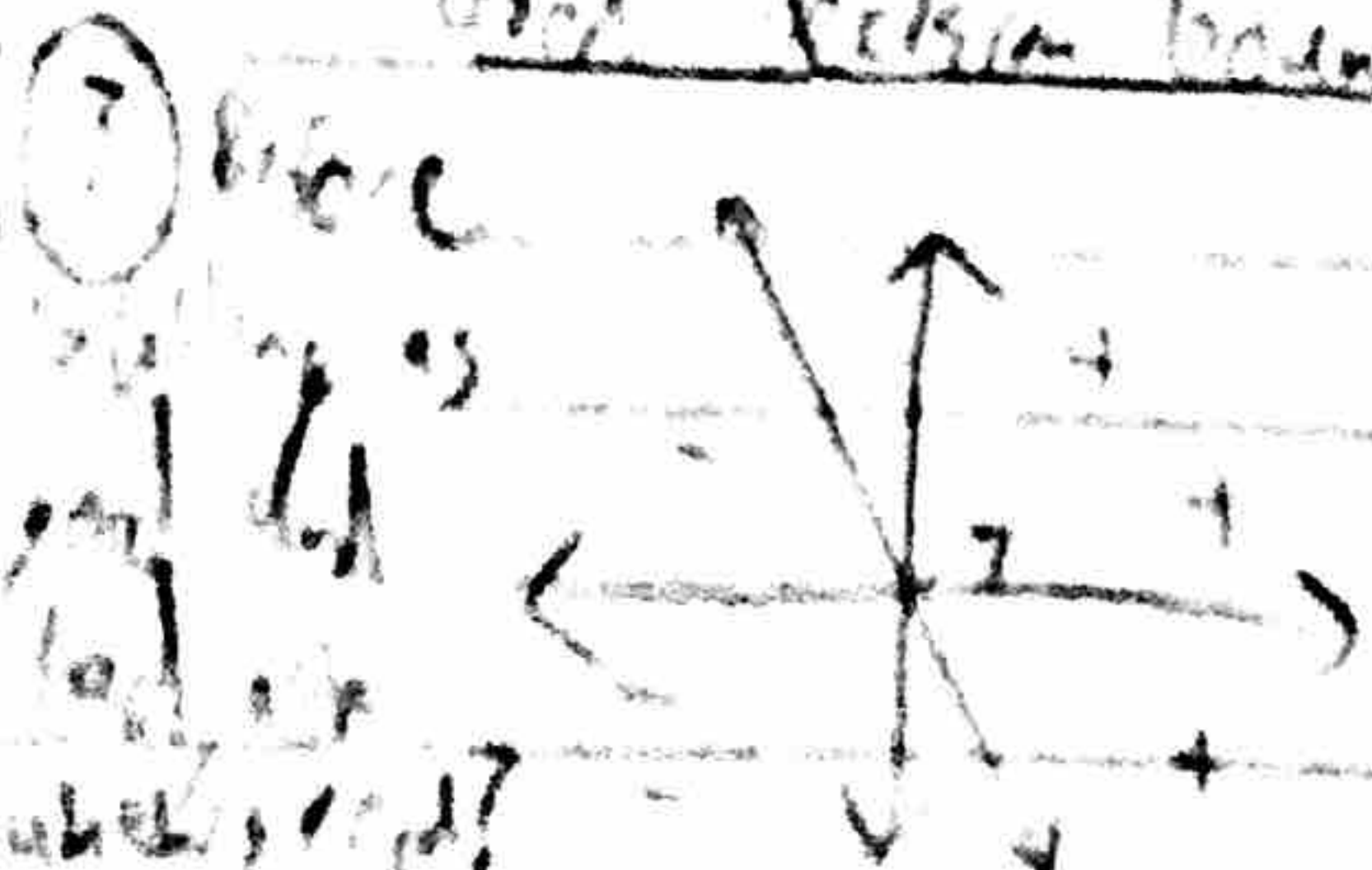
For a decision boundary parameterized by w and b_0 , a linear classifier separator is the function

$$h(x; w, b_0) = \text{sign}(w \cdot x + b_0) = \begin{cases} +1, & w \cdot x + b_0 > 0 \\ -1, & w \cdot x + b_0 \leq 0 \end{cases}$$

Goal: Find the decision boundary (i.e. w and b_0) which best separates (classify) the data

Good decision boundary

Bad decision boundary



Simplification: Linear separator through the origin

(Offsets allow for more powerful classifiers, but linear classifiers without offsets are easier to reason about so we will often leave out the offset for simplicity)

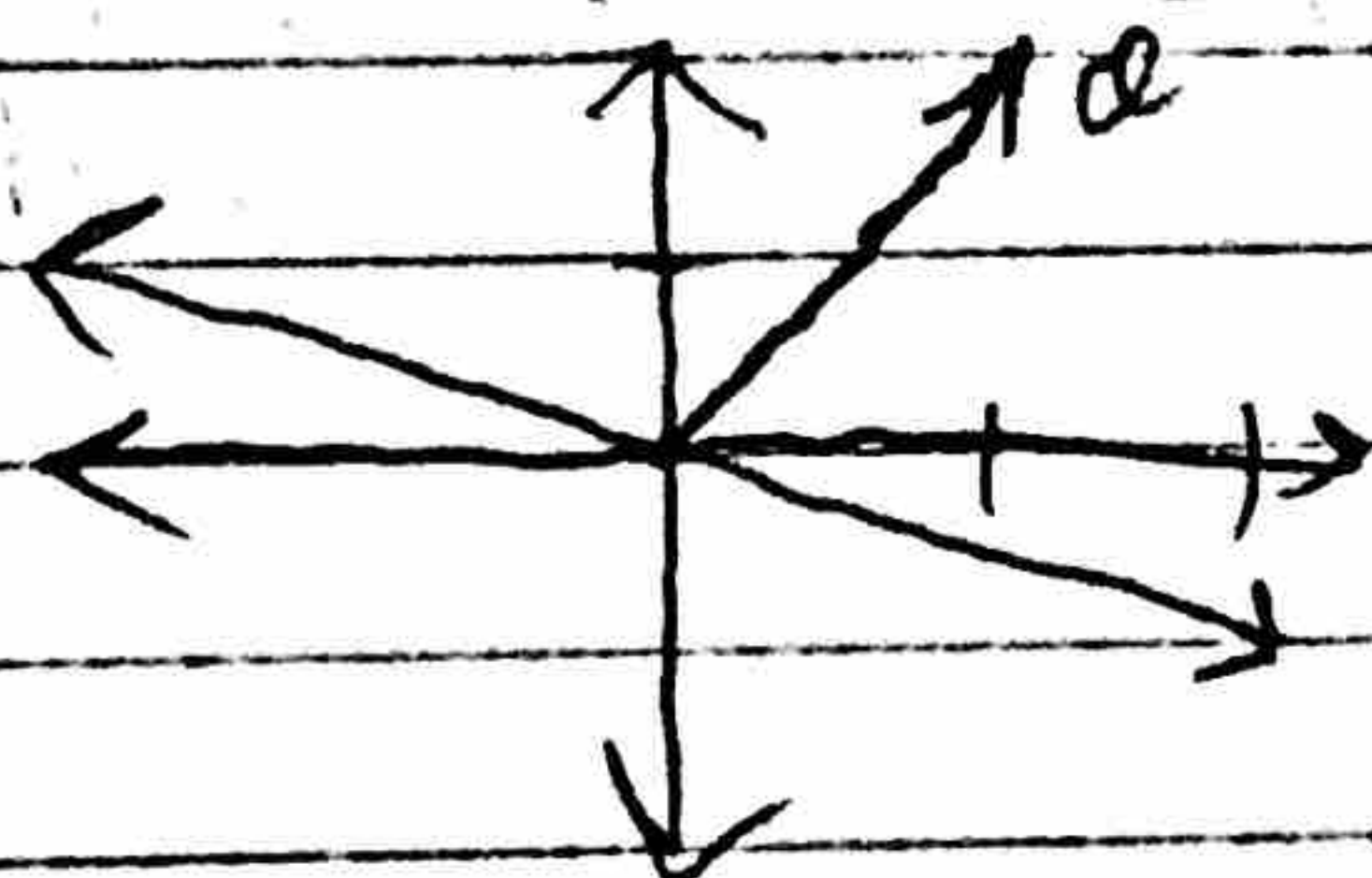
Linear classifier without offset: $h(x; \theta) = \text{sign}(\theta \cdot x) = \begin{cases} +1, & \theta \cdot x > 0 \\ -1, & \theta \cdot x \leq 0 \end{cases}$

Note: Same as $\theta_0 = 0$.

Example: $\theta = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$

$$x_1 + 2x_2 = 0$$

$$x_2 = -\frac{1}{2}x_1$$



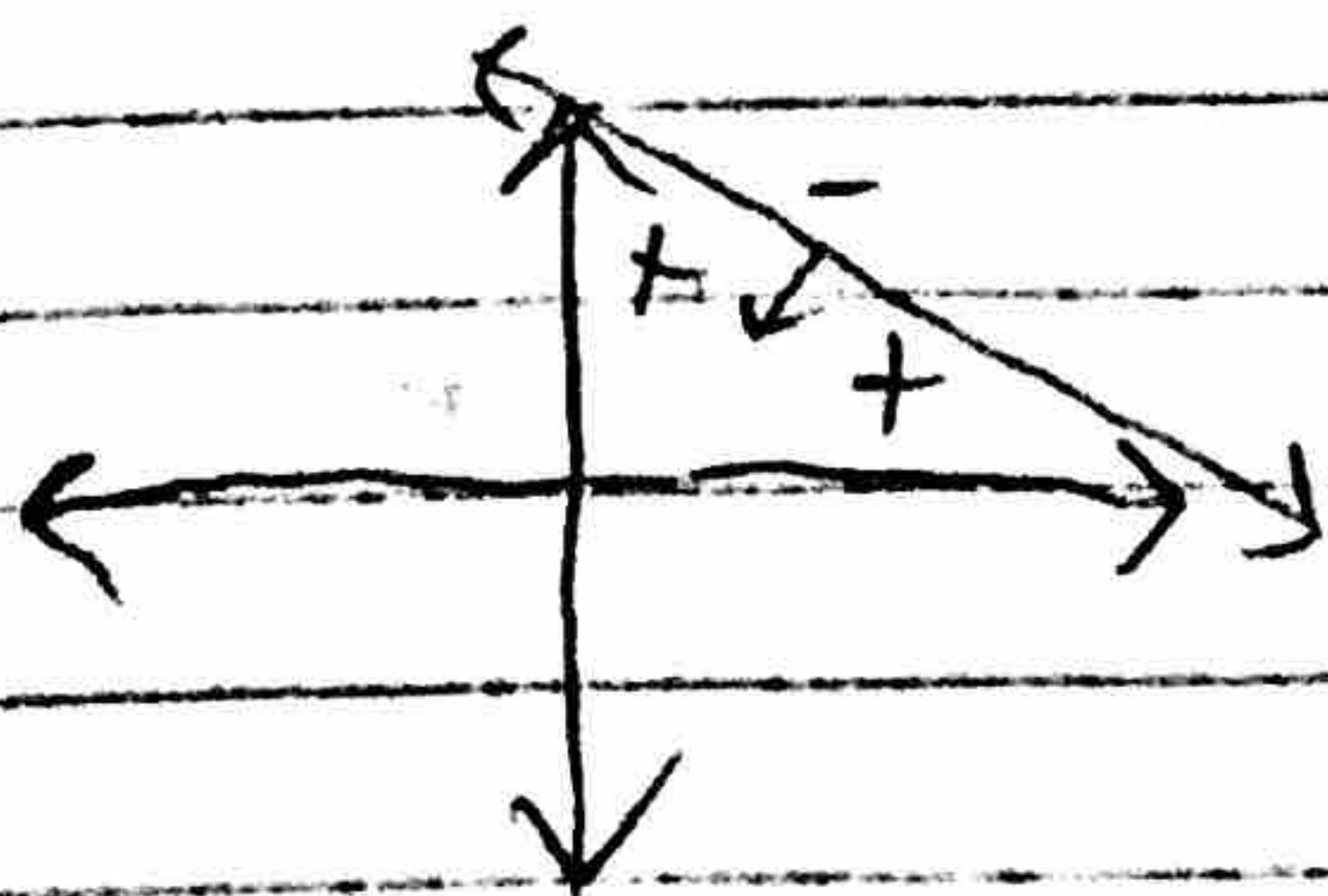
(Note: Since $\theta_0 = 0$, $x = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ will always satisfy $\theta \cdot x = 0$.)

Conclusion: Linear classifiers without offset always go through the origin.)

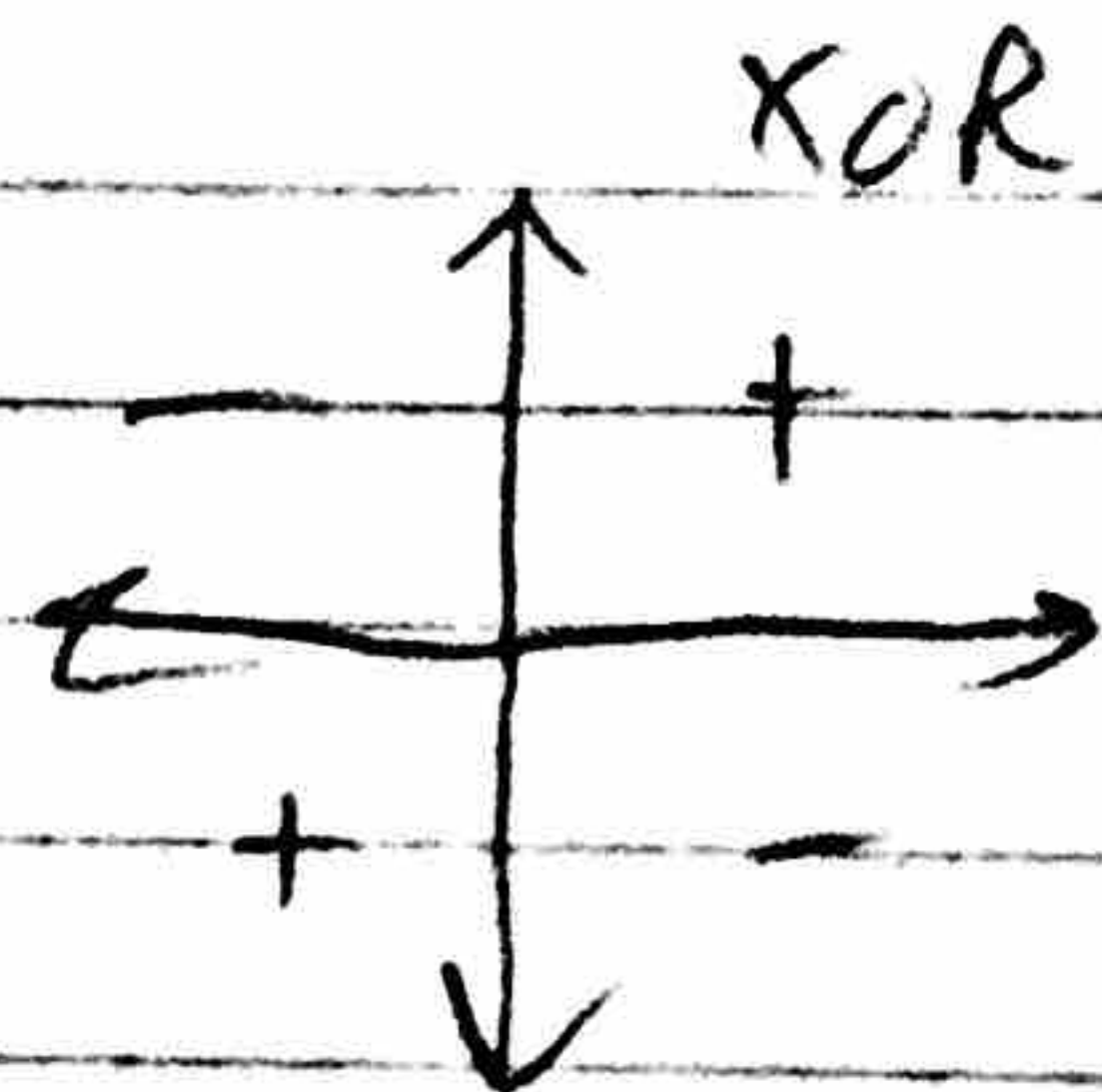
Pro: Easier to reason about.

Con: Weaker classifier

(?) Can you think of a set of labelled points which can only be classified with an offset?



- (?) On the same topic, can you think of a set of points which even a linear classifier with offset can't classify?



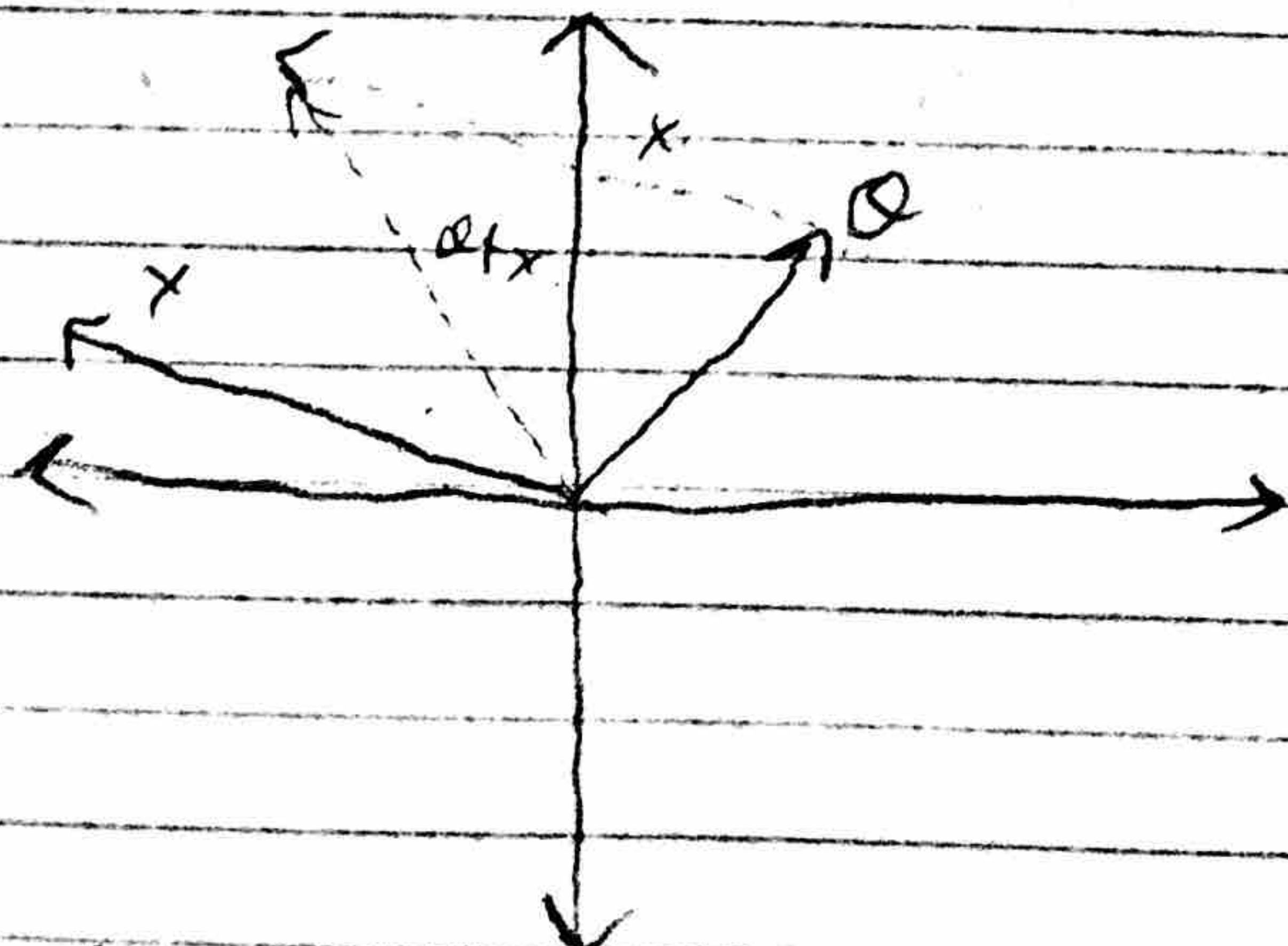
(Remember from last class? Looking ahead: The perceptron algorithm, which we're about to learn, uses linear classifiers. Therefore, perceptron is unable to classify functions such as XOR.)

(Linear classifiers, especially without offset, clearly have limitations, but they are easy to work with and to understand so still worth learning. Furthermore, next week we will see how to make small modifications to make them much more powerful using a technique called kernels.)

Vector Addition

(Last bit of math necessary to understand perceptron)

Graphically



- (?) Can someone draw $Q + x$? (as a vector from the origin)

Algebraically

$$Q = \begin{pmatrix} Q_1 \\ Q_2 \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

- (?) Add element-wise! $Q + x = \begin{pmatrix} Q_1 \\ Q_2 \end{pmatrix} + \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} Q_1 + x_1 \\ Q_2 + x_2 \end{pmatrix}$

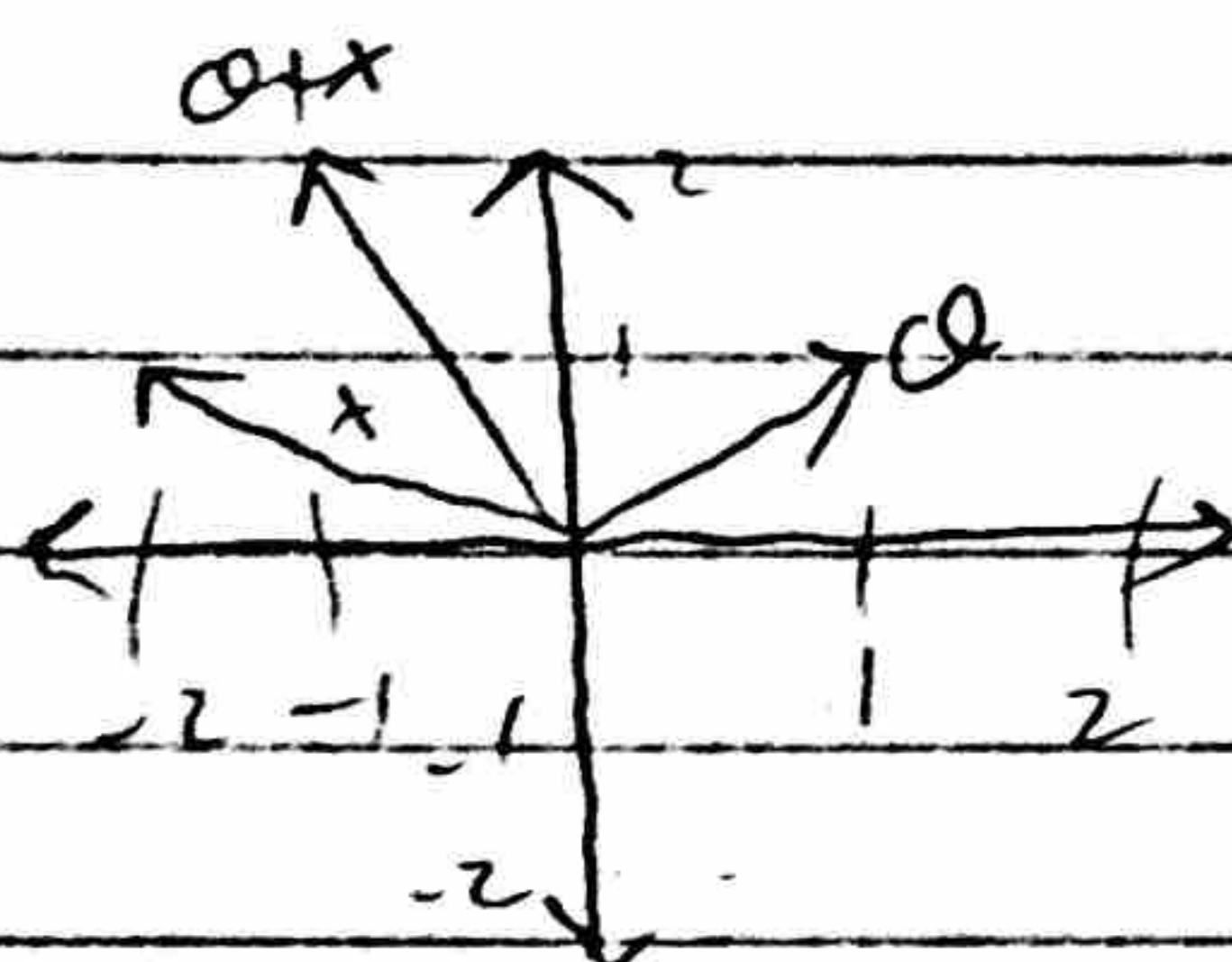
Example

$$\theta = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad x = \begin{pmatrix} -2 \\ 1 \end{pmatrix}$$

(1) compute

$$\theta + x = \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \begin{pmatrix} -2 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ 2 \end{pmatrix}$$

(2) draw



Perceptron algorithm (without offset)

Idea: Mistake-driven algorithm using a linear classifier.

1) Guess a linear classifier θ .

2) For each training example x :

- Check if the linear classifier classifies x correctly.

- If correct, continue.

- If mistake, update θ based on x .

- If x has label $+1$, point θ towards x .

- If x has label -1 , point θ away from x .

3) Repeat step 2 T times.

(What should T be? You'll explore this more in lab today.)

Algorithm:

$$\theta = \vec{0}$$

repeat T times:

for $i \in \{1, 2, \dots, n\}$:

if $y^{(i)} \neq h(x^{(i)}, \theta)$:

$$\theta = \theta + y^{(i)} x^{(i)}$$

$$\left(\begin{array}{l} n = \# \text{ training examples} \\ h(x^{(i)}, \theta) = \text{sign}(x^{(i)} \cdot \theta) \end{array} \right)$$

return θ

(?) Why adding $y^{(i)} x^{(i)}$?

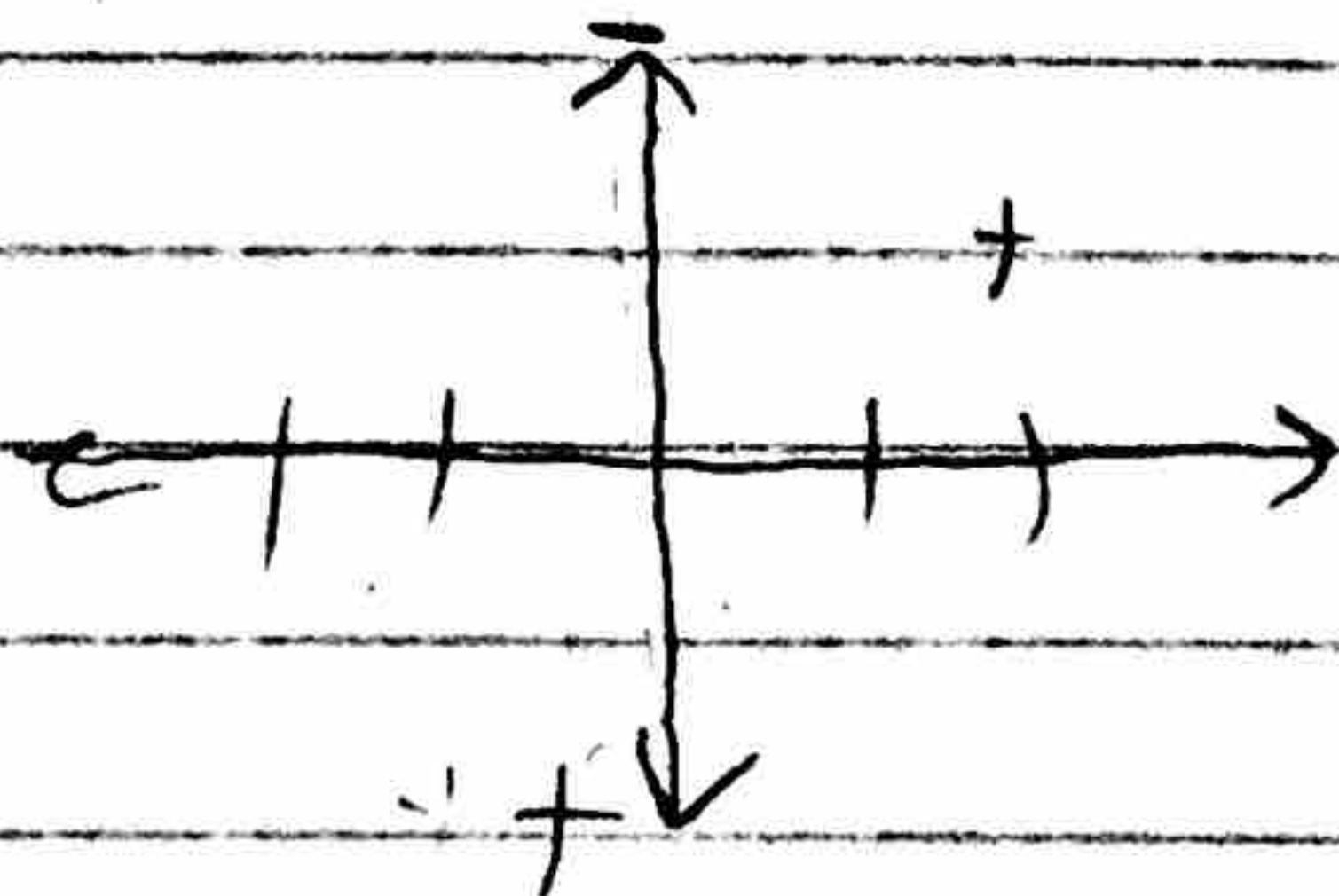
($y^{(i)}$ says whether to move toward or away from x . x tells direction.)

Example

$$x^{(1)} = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \quad y^{(1)} = +1$$

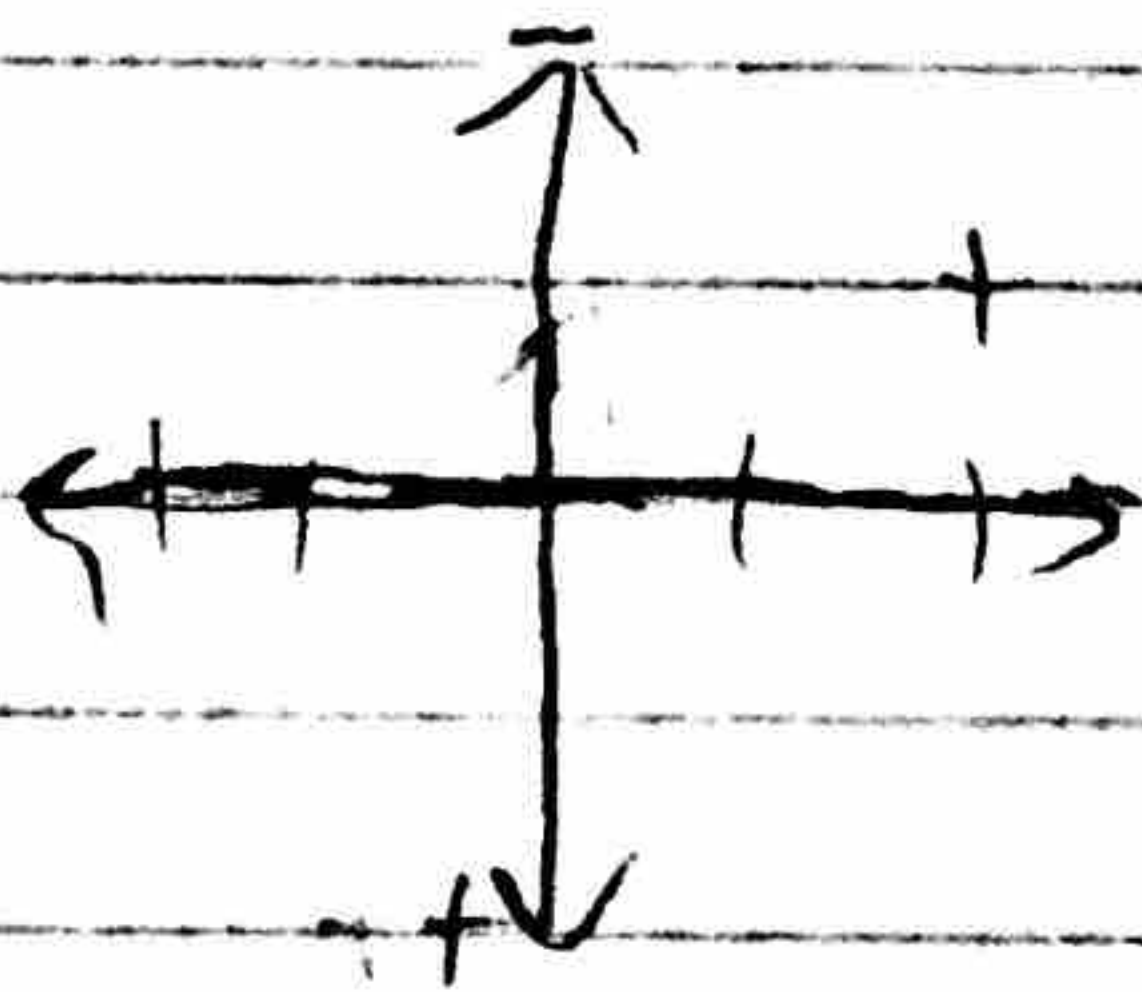
$$x^{(2)} = \begin{pmatrix} 0 \\ 2 \end{pmatrix} \quad y^{(2)} = -1$$

$$x^{(3)} = \begin{pmatrix} -0.5 \\ -2 \end{pmatrix} \quad y^{(3)} = +1$$

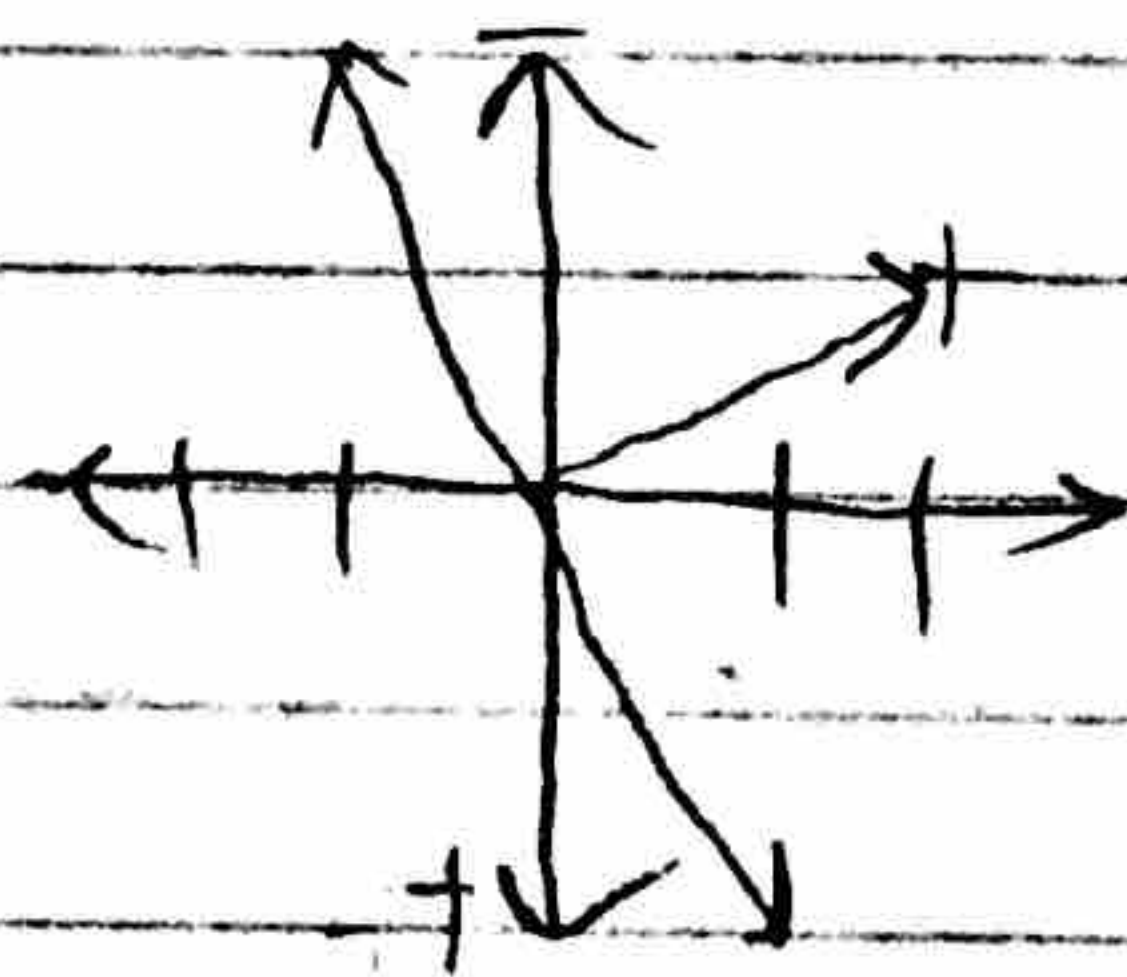


(?) Guess where decision boundary will be. Get help to do Perceptron

① $\theta = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$



② $\theta = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$



$$h(x^{(1)}; \theta) = \text{sign}(\theta \cdot x^{(1)})$$

$$= \text{sign}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 2 \\ 1 \end{pmatrix}\right) = \text{sign}(0) = -1$$

$$h(x^{(1)}; \theta) = -1$$

$$y^{(1)} = +1$$

$y^{(1)} \neq h(x^{(1)}; \theta) \rightarrow \text{mistake}$

update: $\theta = \theta + y^{(1)} x^{(1)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} + 1 \cdot \begin{pmatrix} 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$

$$h(x^{(2)}; \theta) = \text{sign}(\theta \cdot x^{(2)})$$

$$= \text{sign}\left(\begin{pmatrix} 2 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 2 \end{pmatrix}\right) = \text{sign}(2) = +1$$

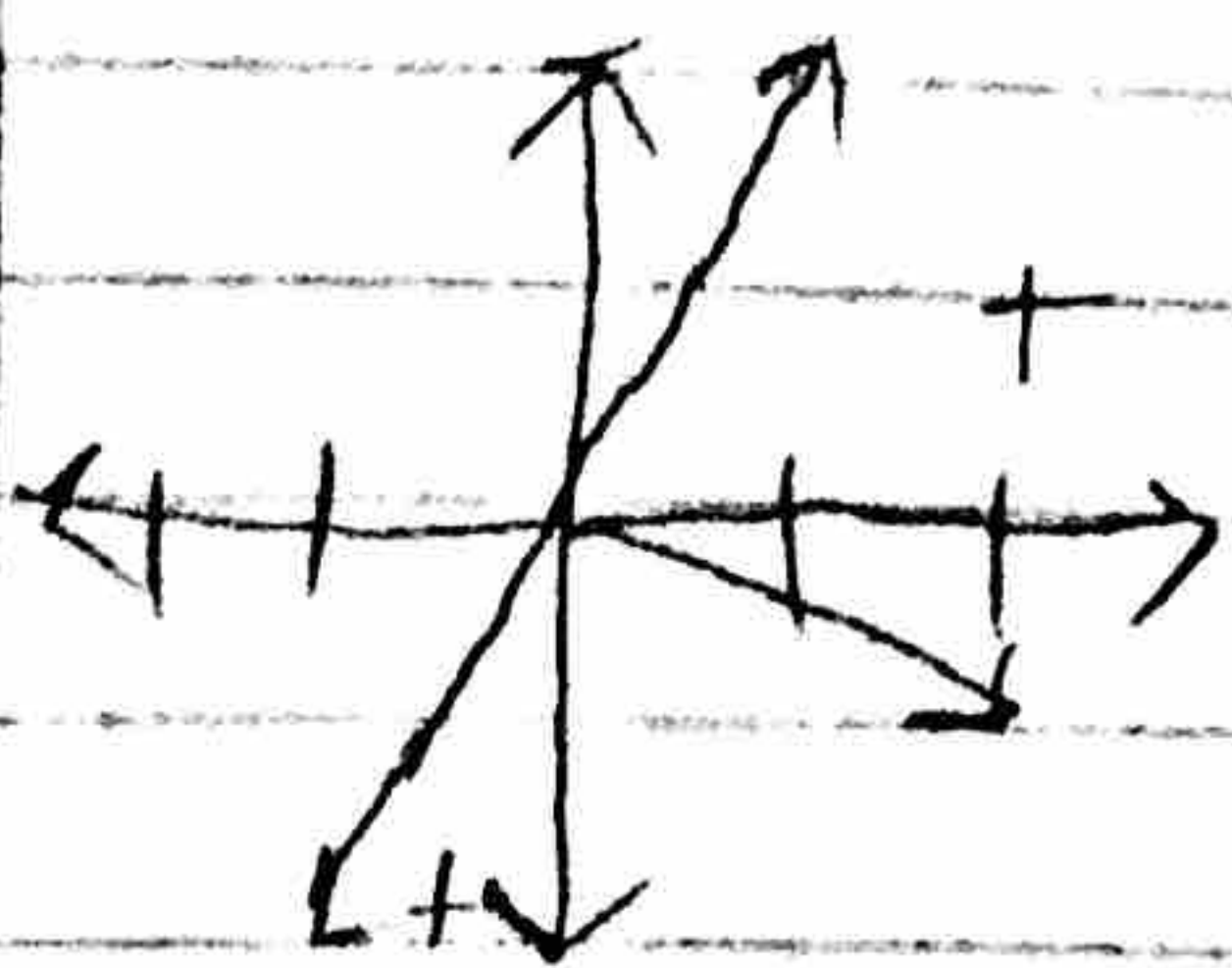
$$h(x^{(2)}; \theta) = +1$$

$$y^{(2)} = -1$$

$y^{(2)} \neq h(x^{(2)}; \theta) \rightarrow \text{mistake}$

update: $\theta = \theta + y^{(2)} x^{(2)} = \begin{pmatrix} 2 \\ 1 \end{pmatrix} - 1 \cdot \begin{pmatrix} 0 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$

③ $\theta = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$



$$h(x^{(3)}; \theta) = \text{sign}(\theta \cdot x^{(3)}) = \text{sign}\left(\begin{pmatrix} 2 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} -0.5 \\ -2 \end{pmatrix}\right)$$

$$= \text{sign}(-1 + 2) = \text{sign}(1) = +1$$

$$y^{(3)} = +1$$

$$y^{(3)} = h(x^{(3)}; \theta) \rightarrow \text{correct}$$

no update

(If we check $x^{(1)}$ and $x^{(2)}$, we will see that they are also correctly classified.)

(Since there are no more mistakes, ϕ will no longer change and we can say that the perceptron converged.)

Perceptron algorithm (with offset)

$$\phi = \vec{0}, \phi_0 = 0$$

repeat T times:

for $i = 1, 2, \dots, n$:

if $y^{(i)} \neq h(x^{(i)}; \phi, \phi_0)$; (note ϕ_0)

$$\phi = \phi + y^{(i)} x^{(i)}$$

$$\phi_0 = \phi_0 + y^{(i)}$$

return ϕ, ϕ_0

(*) Explain that we can think of ϕ_0 just as an extra dimension of x which is always 1. So we can actually think of a classifier with offset as a classifier with no offset in $d+1$ dimensions with $\phi = \begin{pmatrix} \phi_1 \\ \vdots \\ \phi_d \\ \phi_0 \end{pmatrix}$ and $x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \\ 1 \end{pmatrix}$

$$\text{b/c } \phi \cdot x = \phi_1 x_1 + \dots + \phi_d x_d + \phi_0 \cdot 1.$$

Questions:

(?) 1) Does the perceptron algorithm always converge?

(i.e. reach a state where all training points are classified correctly)

No - example is XOR where ϕ will keep spinning in circles

(?) 2) Are there any cases where the perceptron algorithm is guaranteed to converge?

Yes - whenever the points are linearly separable

Def: A training set S_n is linearly separable if there exists a $\hat{\phi}$ such that $h(x^{(i)}; \hat{\phi}) = y^{(i)}$ for all $i = 1, 2, \dots, n$.

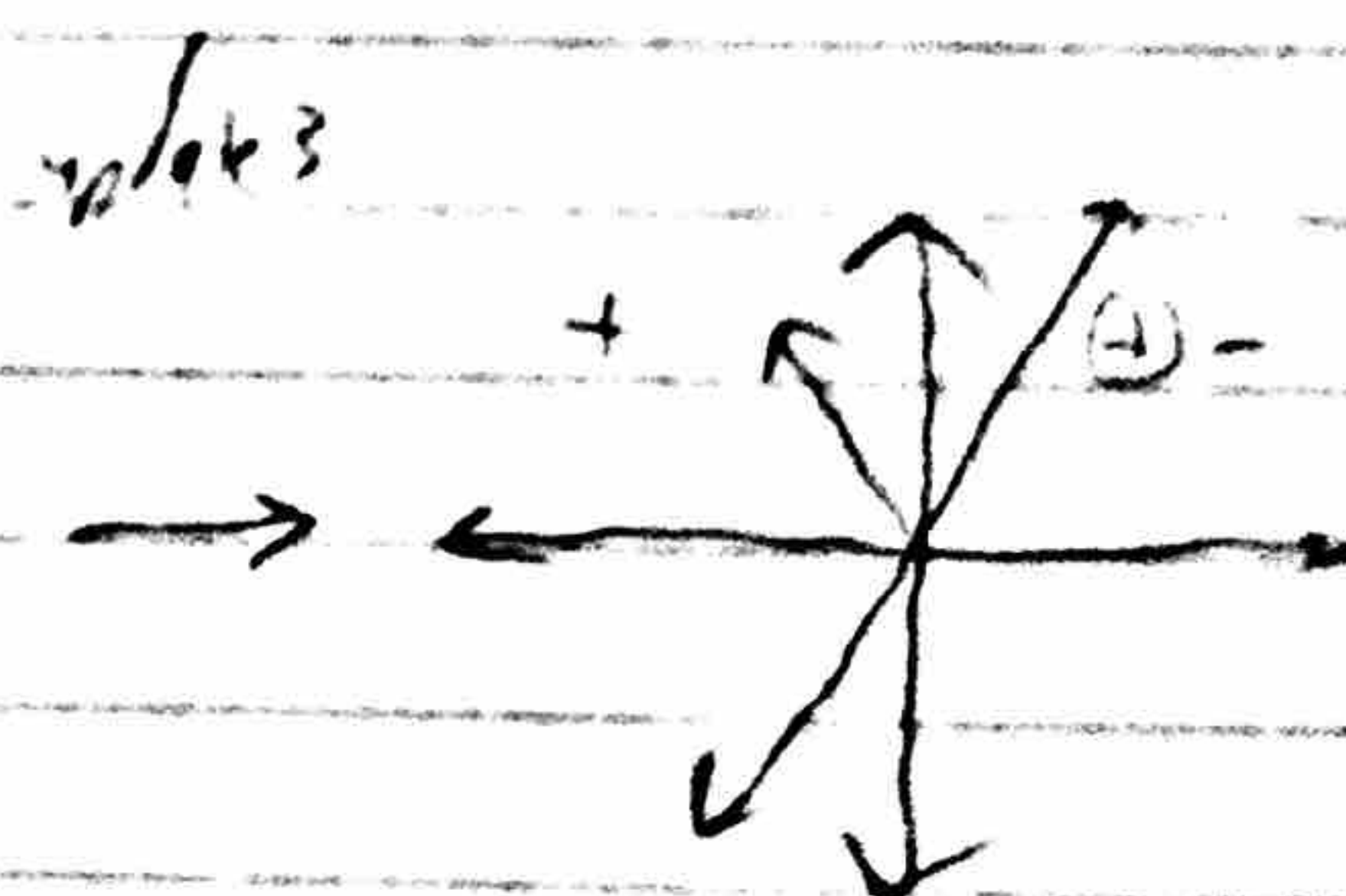
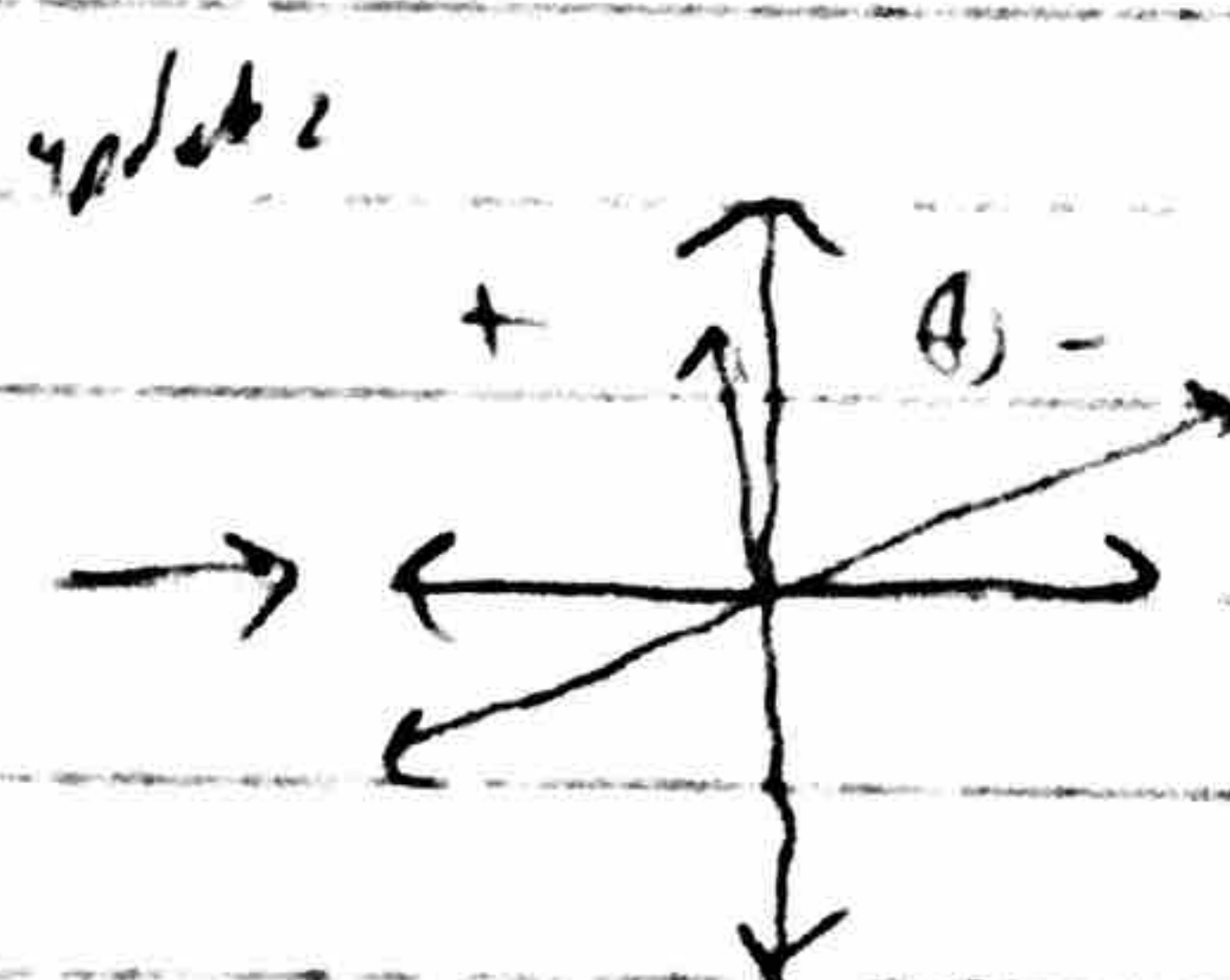
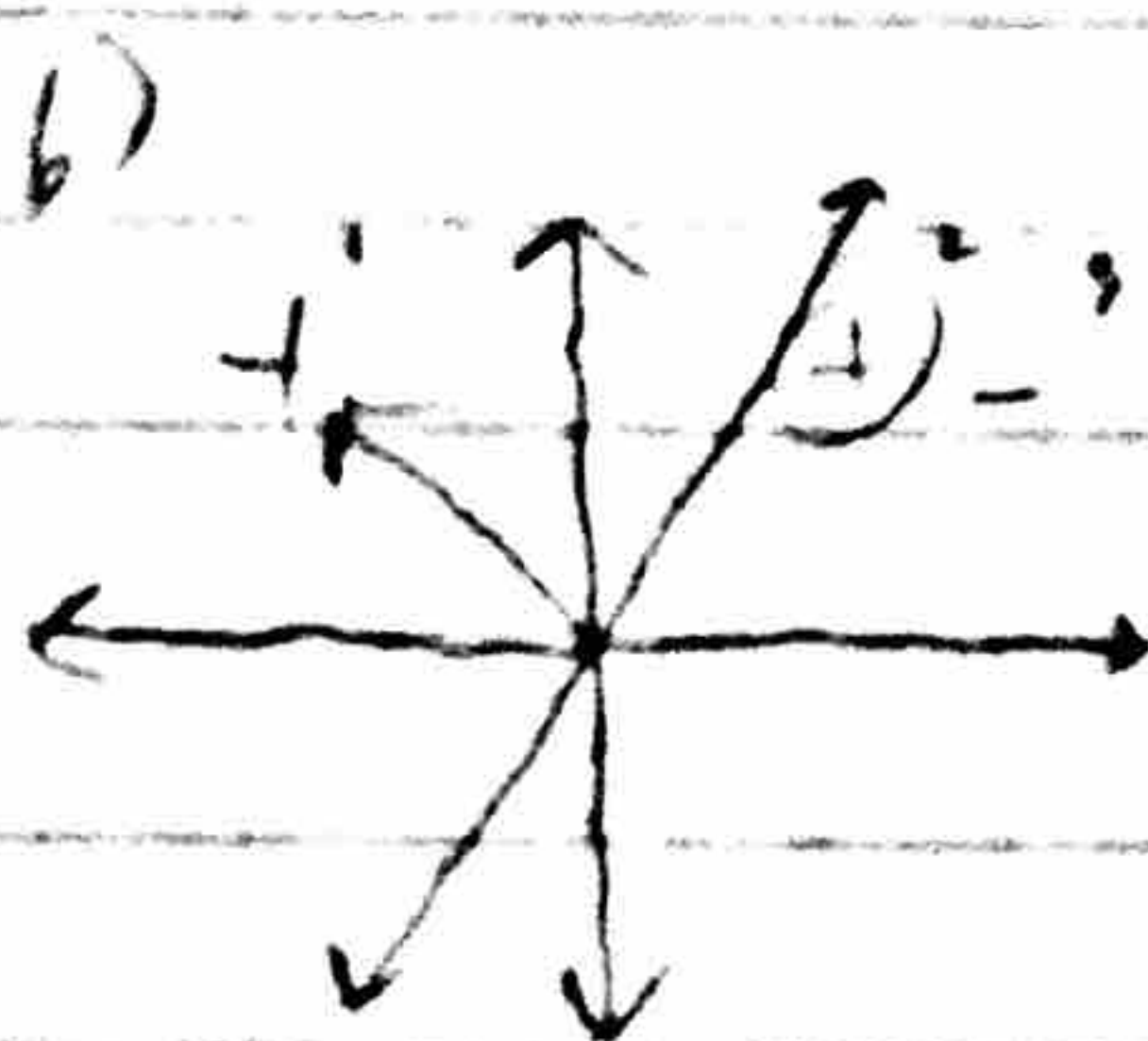
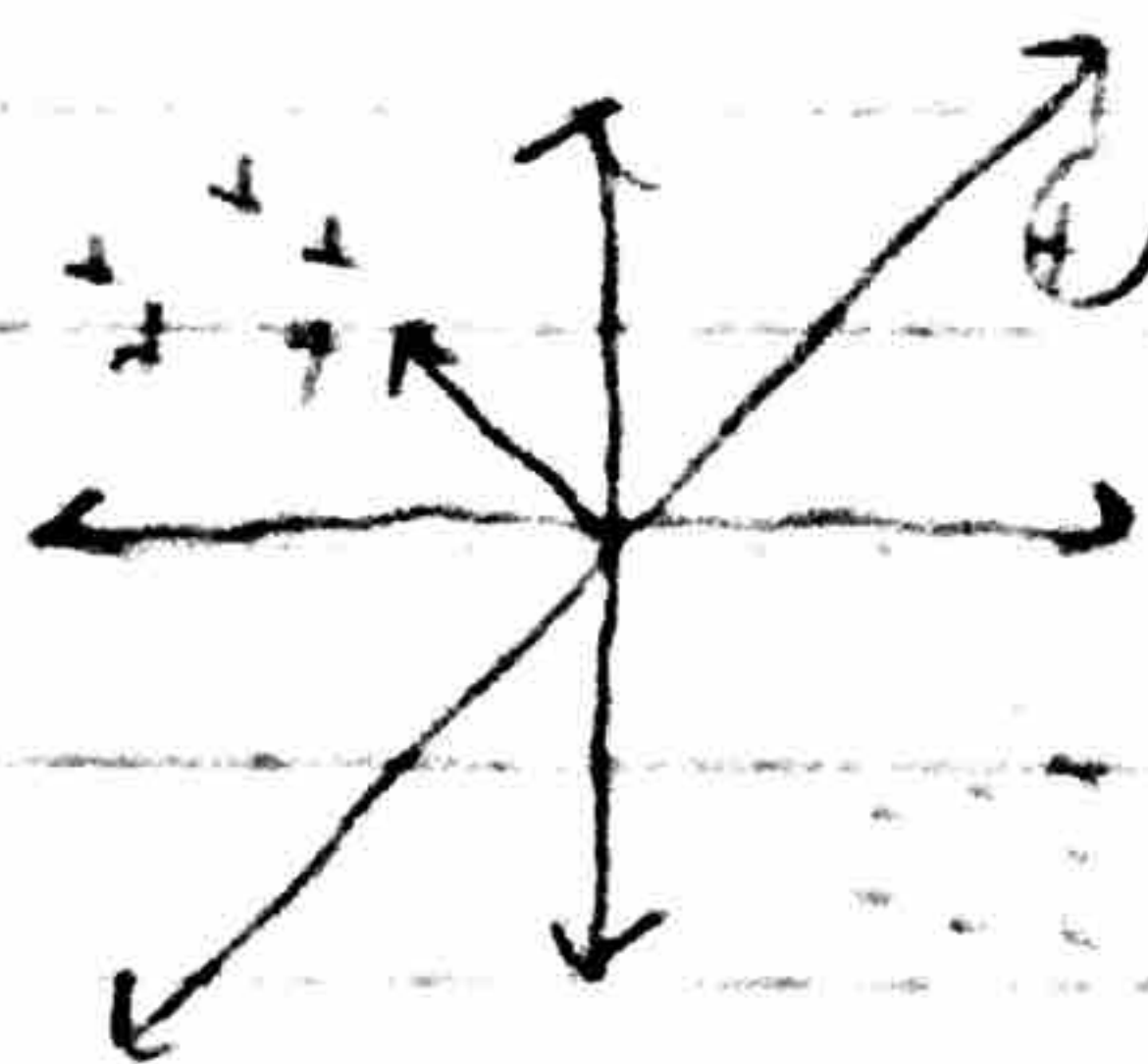
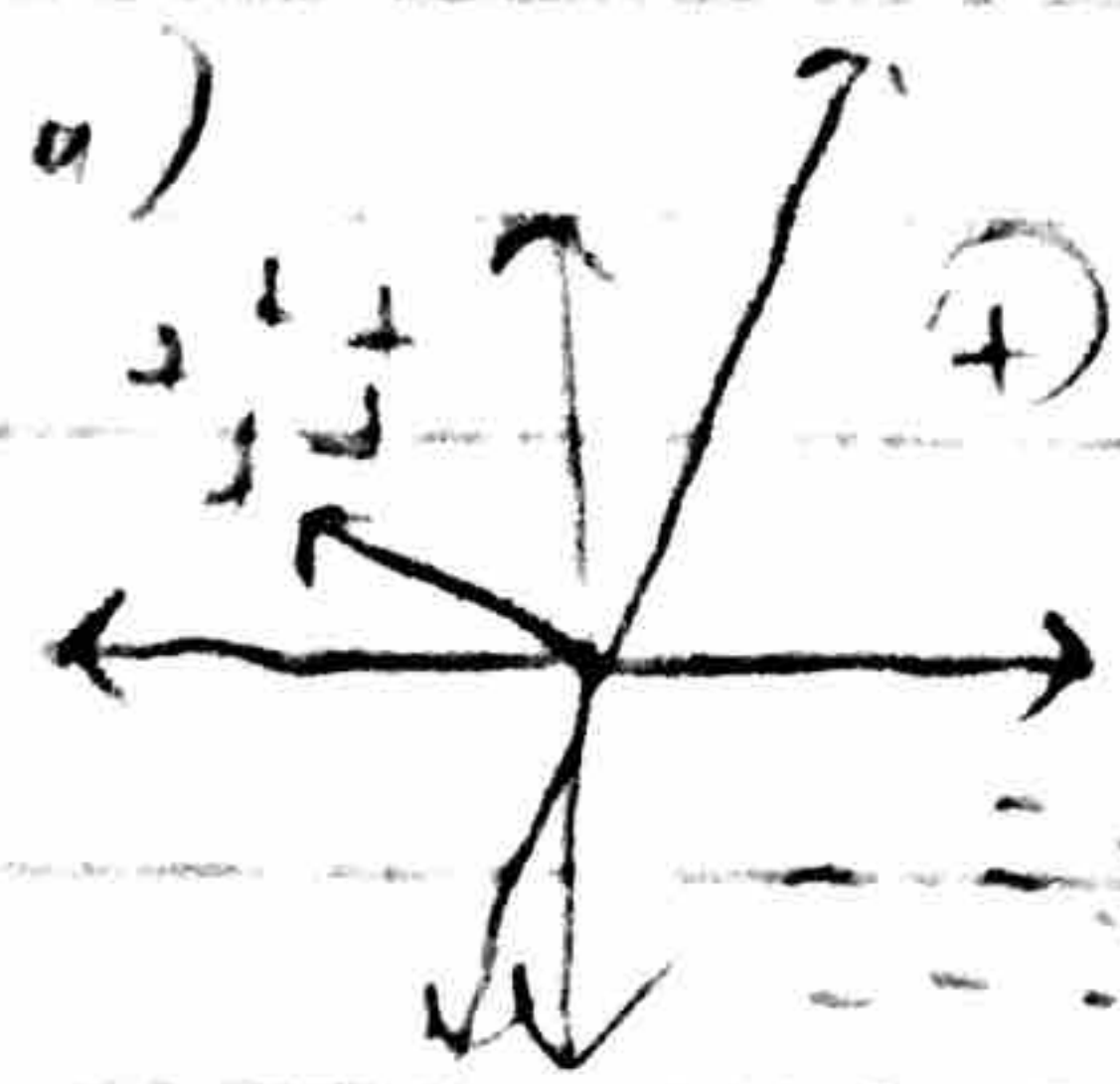
(Talk to me after class or in lab if you want to see proof of convergence)

- (?) 3) When the perceptron does converge, does it always converge to the same solution?
 No - depends on initialization and order in which data points are traversed

2-8

- (?) 4) Once the perceptron misclassifies a point, is it guaranteed to correctly classify the point in the future?

No - a) Even after the update to θ it can still misclassify the point.
 b) Also updates based on other points could change how it's classified



- 5) What can we say about how the update step affects the classification?
 The agreement increases

Def: Agreement = $y^{(i)}(\theta \cdot x^{(i)} + \theta_0)$ or just $y^{(i)}(\theta \cdot x^{(i)})$ if $\theta_0 = 0$

$$y^{(i)} \theta \cdot x^{(i)} > 0 \iff y^{(i)} = \text{sign}(\theta \cdot x^{(i)}) = h(x^{(i)}; \theta) \iff \text{correct}$$

$$y^{(i)} \theta \cdot x^{(i)} \leq 0 \iff y^{(i)} \neq \text{sign}(\theta \cdot x^{(i)}) = h(x^{(i)}; \theta) \iff \text{misclassified}$$

(Claim): Agreement on $x^{(i)}$ increases after updating θ .

Proof

Agreement before update = $p = y^{(i)} \theta \cdot x^{(i)}$

Update: $\theta' = \theta + y^{(i)} x^{(i)}$

Agreement after update = $p' = y^{(i)} \theta' \cdot x^{(i)} = y^{(i)} (\theta + y^{(i)} x^{(i)}) \cdot x^{(i)}$

$$= y^{(i)} (\theta \cdot x^{(i)} + y^{(i)} x^{(i)} \cdot x^{(i)}) = y^{(i)} \theta \cdot x^{(i)} + y^{(i)^2} \|x^{(i)}\|^2 = p + \|x^{(i)}\|^2 \geq p$$

$$\|x\|^2 = \sqrt{x \cdot x} = \sqrt{x_1^2 + \dots + x_n^2} \geq 0$$

Improvement: Average Perceptron

Problem: θ can jump around and final theta may be biased toward late-points.

Idea: Average θ across all time steps. Gives more weight to θ 's which worked for a large number of points without being updated.

Algorithm

$$\theta = \vec{0}$$

$$\theta_{sum} = \vec{0}$$

repeat T times:

for $i=1, 2, \dots, n$:

if $y^{(i)} \neq h(x^{(i)}, \theta)$:

$$\theta = \theta + y^{(i)} x^{(i)}$$

$$\theta_{sum} = \theta_{sum} + \theta$$

$$\text{return } \frac{\theta_{sum}}{nT}$$

(Note: not 1, the if so happens on every step)

$$\text{(works b/c } \bar{\theta} = \frac{\theta_1 + \theta_2 + \dots + \theta_{nT}}{nT})$$

Questions?