

Lecture 3 - Maximum Margin Classifiers and Support Vector Machines

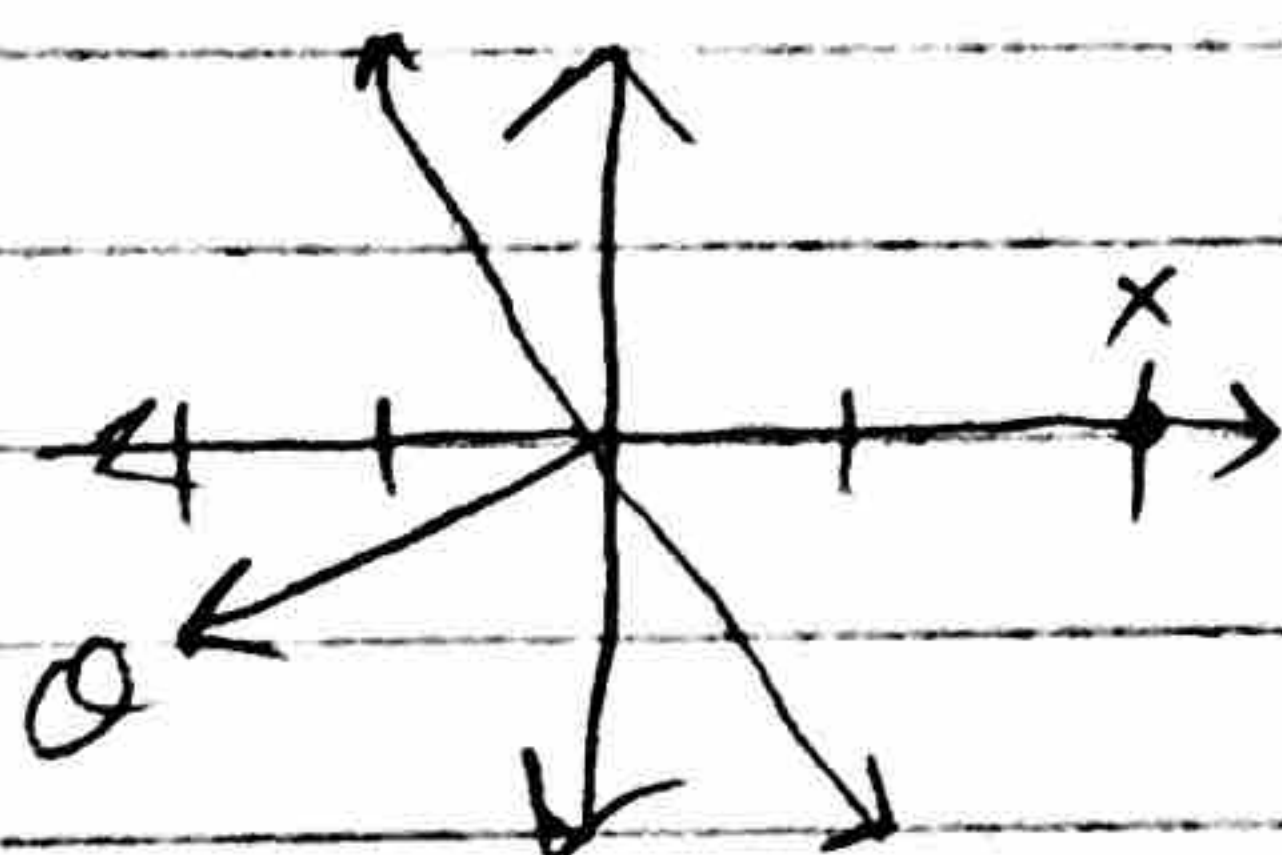
Review:

(1) Let $\theta = \begin{pmatrix} -2 \\ -1 \end{pmatrix}$, $\theta_0 = 0$.

a) Draw θ .

b) Draw the decision boundary.

c) Let $x = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$. What is $h(x; \theta)$? (You can either compute it or determine it by looking at your graph.)



$$h(x; \theta) = \text{sign}(\theta \cdot x) = \text{sign}\left(\begin{pmatrix} -2 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} 2 \\ 0 \end{pmatrix}\right) \\ = \text{sign}(-2 \cdot 2 + (-1) \cdot 0) = \text{sign}(-4) = -1 \\ h(x; \theta) = -1$$

(2) Let $\theta = \begin{pmatrix} 1 \\ 3 \\ -2 \end{pmatrix}$, $\theta_0 = -1$, $x = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}$, $y = -1$.

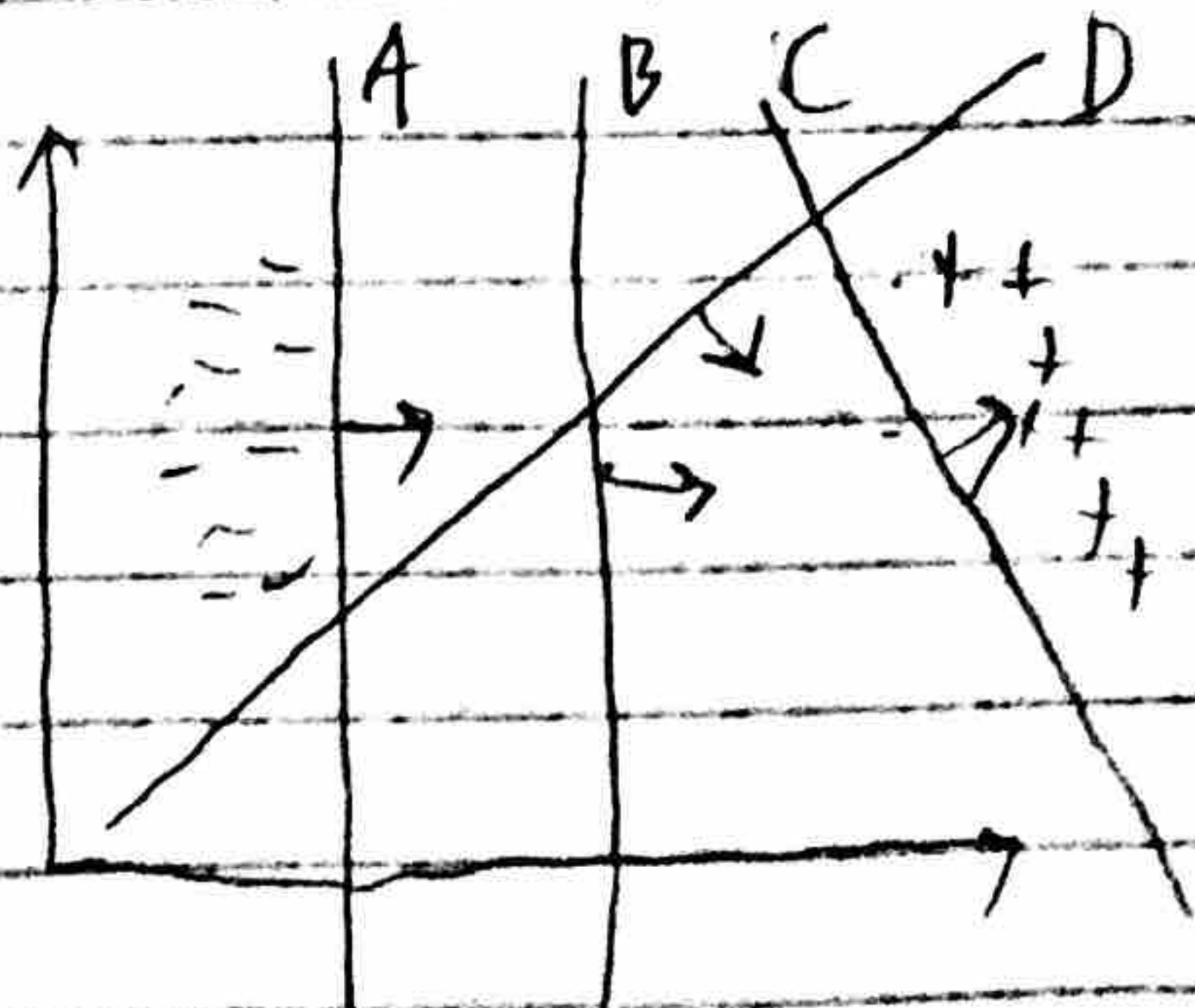
a) What is the agreement?

$$\text{agreement} = y(\theta \cdot x + \theta_0) = -1 \left[\begin{pmatrix} 1 \\ 3 \\ -2 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} - 1 \right] \\ = -[1 + 3 - 4 - 1] = -1 \cdot (-1) = 1$$

b) Is x classified correctly?

Yes - agreement > 0

(3) Which is the best linear classifier for the following set of data?



Answer depends on "best".
All correctly classify all the points so all are perfect wrt training points. But B is best in terms of ability to generalize. What is special about B? Its margin.

Today

- Margins
 - Motivation
 - Definition
 - Examples
- Support Vector machines
 - Objective function
 - Hinge loss
 - Regularization
 - Offline vs. online algorithms
 - Online SVM
 - Gradient descent
 - Pegasos algorithm

Margins

Motivation: Perceptron has many solutions, which is best?

Goal: Generalize to the test set.

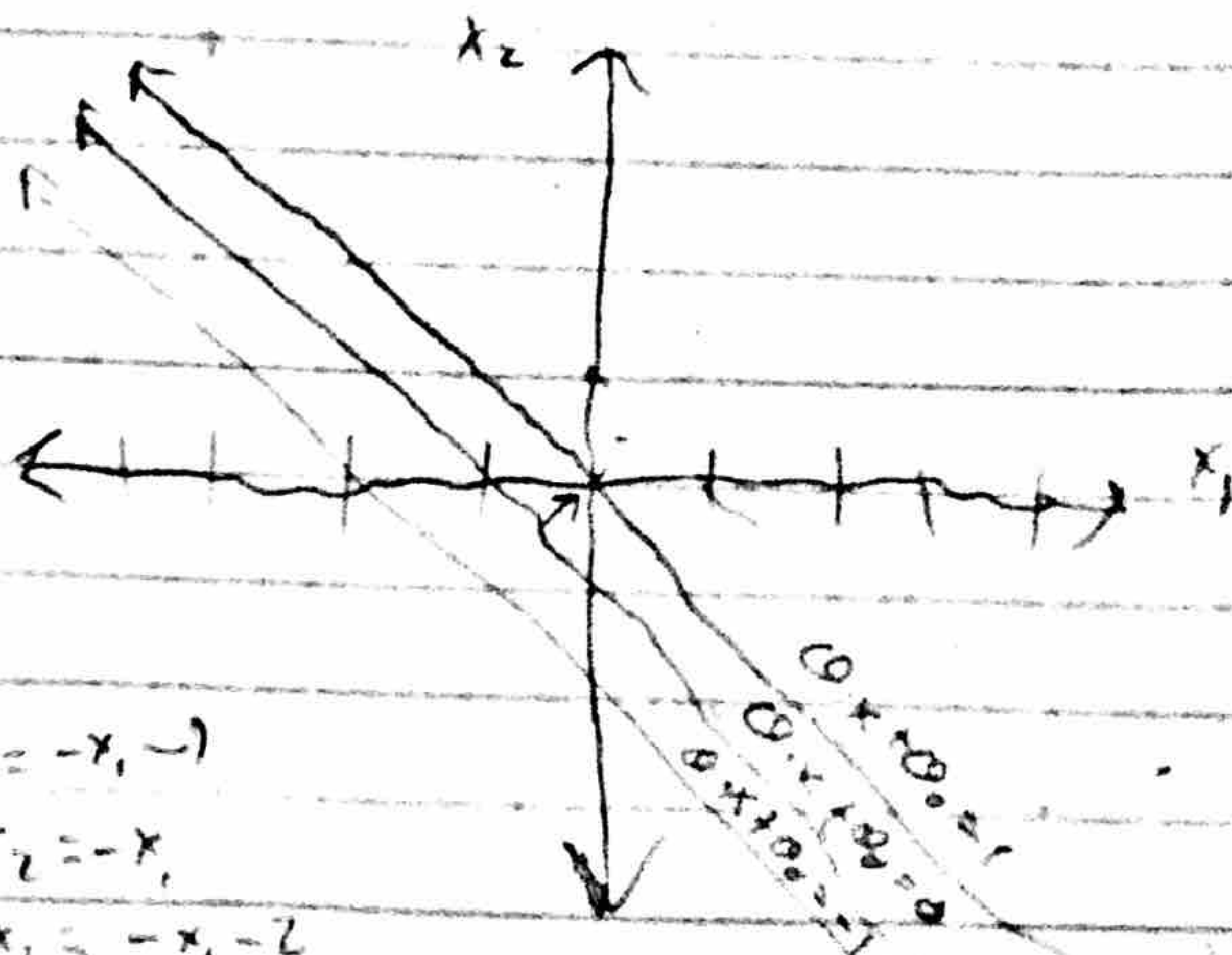
How? Maximize the margin (between the decision boundary and the data points)

Def: The margin of a linear classifier $h(x; \theta, \theta_0)$ is the set of all points where $\theta \cdot x + \theta_0 = \pm 1$.

The margin is $\left[\begin{array}{l} \text{a pair of lines in } \mathbb{R}^2 \\ \text{a pair of hyperplanes in } \mathbb{R}^d \end{array} \right]$ parallel to the decision boundary.

Example

$$\theta = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \theta_0 = 1$$



(?) Draw decision boundary.

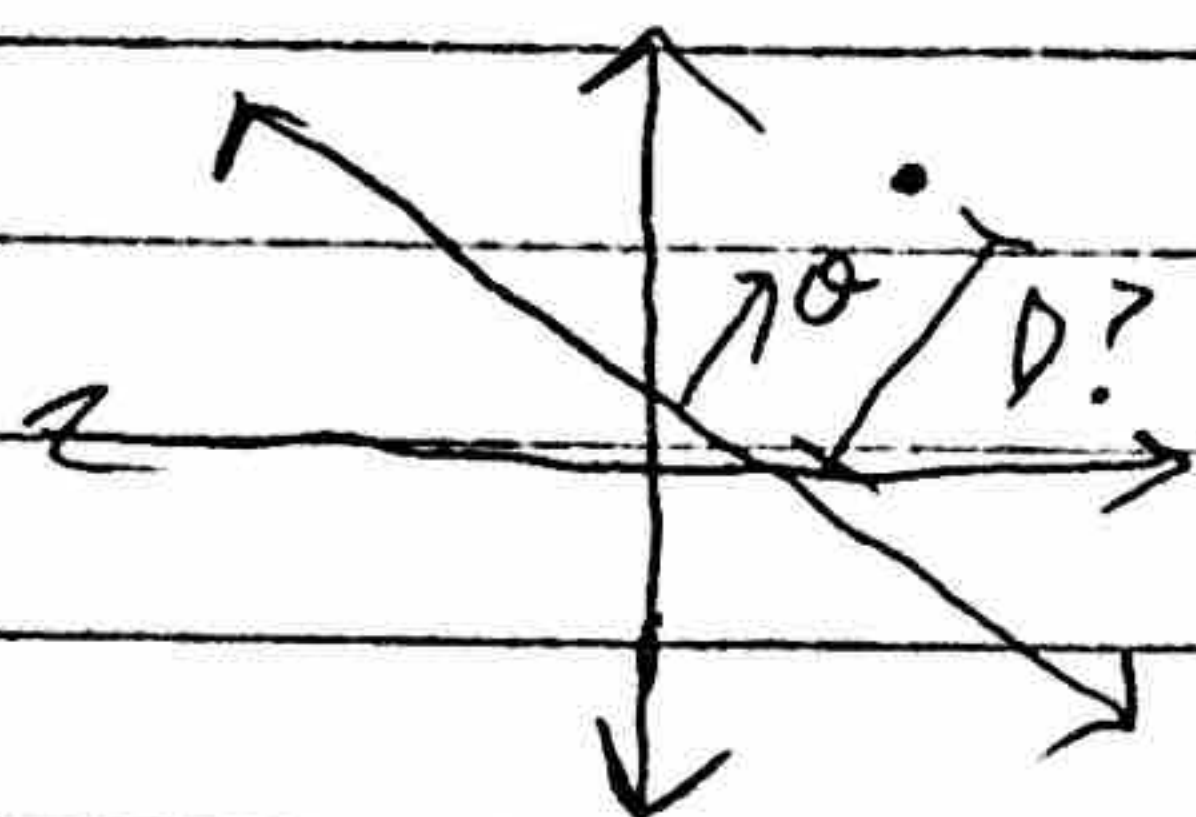
(?) Draw margins.

$$\theta \cdot x + \theta_0 = 0 \rightarrow x_1 + x_2 + 1 = 0 \rightarrow x_2 = -x_1 - 1$$

$$\theta \cdot x + \theta_0 = 1 \rightarrow x_1 + x_2 + 1 = 1 \rightarrow x_2 = -x_1$$

$$\theta \cdot x + \theta_0 = -1 \rightarrow x_1 + x_2 + 1 = -1 \rightarrow x_2 = -x_1 - 2$$

(?) Distance to the margin
Distance from a line to a point



$$D = \frac{|Q \cdot x + Q_0|}{\|Q\|}$$

Reminder: $\|Q\| = \sqrt{Q_1^2 + Q_2^2 + \dots + Q_d^2}$

= magnitude of the vector Q

= distance from the origin to the point Q

Distance to the margin

$$D = \frac{|Q \cdot x + Q_0|}{\|Q\|} = \frac{|\pm 1|}{\|Q\|} = \frac{1}{\|Q\|}$$

$$D = \frac{1}{\|Q\|}$$

Note: Distance to the margin is inversely proportional to the magnitude of Q .

Large $Q \iff$ small margin

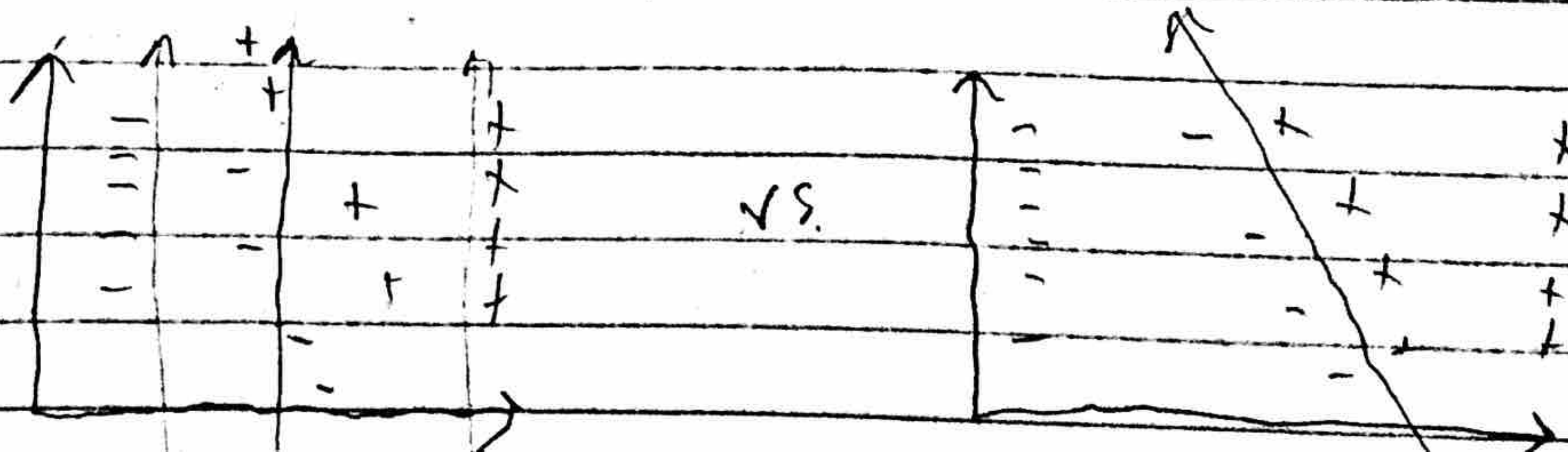
Small $Q \iff$ large margin

(*) The direction of Q does not affect the distance to the margin, only the magnitude matters.

Takeaway: Better generalization \iff large margin \iff small Q

(7) Question: What happens if we are too much about maximizing the margin?
Answer: Potential loss of accuracy.

Example:



(?) Question: How can we balance the competing interests of large margin and high accuracy?

3-4

Answer: SVM

Support Vector Machines (Maximum margin linear classifier)

Objective: High accuracy and large margin
Low loss and small ϕ

Objective function:
$$\min_{\phi} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y^{(i)} \phi \cdot x^{(i)}) + \frac{\gamma}{2} \|\phi\|^2$$

Taking it apart:

1) $\min_{\phi} \rightarrow$ Find the ϕ which minimizes the expression.

2) $\frac{1}{n} \sum_{i=1}^n \rightarrow$ Averages the loss across all training examples.

3) $\text{Loss}(y^{(i)} \phi \cdot x^{(i)}) \rightarrow$ The "loss" of a single training example.
(will explain shortly.)

4) $\|\phi\|^2 \rightarrow$ The magnitude of ϕ (squared). Since we're doing a min, this will give us small $\|\phi\|$ so large margin.

This is called a regularization term because it regulates (controls) the complexity of the model and helps prevent overfitting and encourages generalization.

5) $\frac{\gamma}{2} \rightarrow$ is a hyperparameter which allows us to control which is more important: accuracy or margin.
($\frac{1}{2}$ is arbitrary, makes math easier.)

large $\gamma \rightarrow \|\phi\|$ dominates and small ϕ / large margin is most important
small $\gamma \rightarrow$ loss dominates and high accuracy is most important

(As in lab2, we can test different values of γ to see what performs best.)

(check if all parts make sense.)

(So now that we understand how SVM makes its tradeoffs, what is Loss?)

Loss functions

Loss functions are a generalization of accuracy (like agreement).

Loss functions are a measure of how far away we are from our objective.

Smaller losses better, goal is to minimize loss.

Loss is defined for individual examples. Total loss is sum of loss for each example.

0-1 loss

$$\text{Loss}_{0-1}(y^{(i)} @ x^{(i)}) = \begin{cases} 1 & y^{(i)} @ x^{(i)} \leq 0 \\ 0 & y^{(i)} @ x^{(i)} > 0 \end{cases}$$

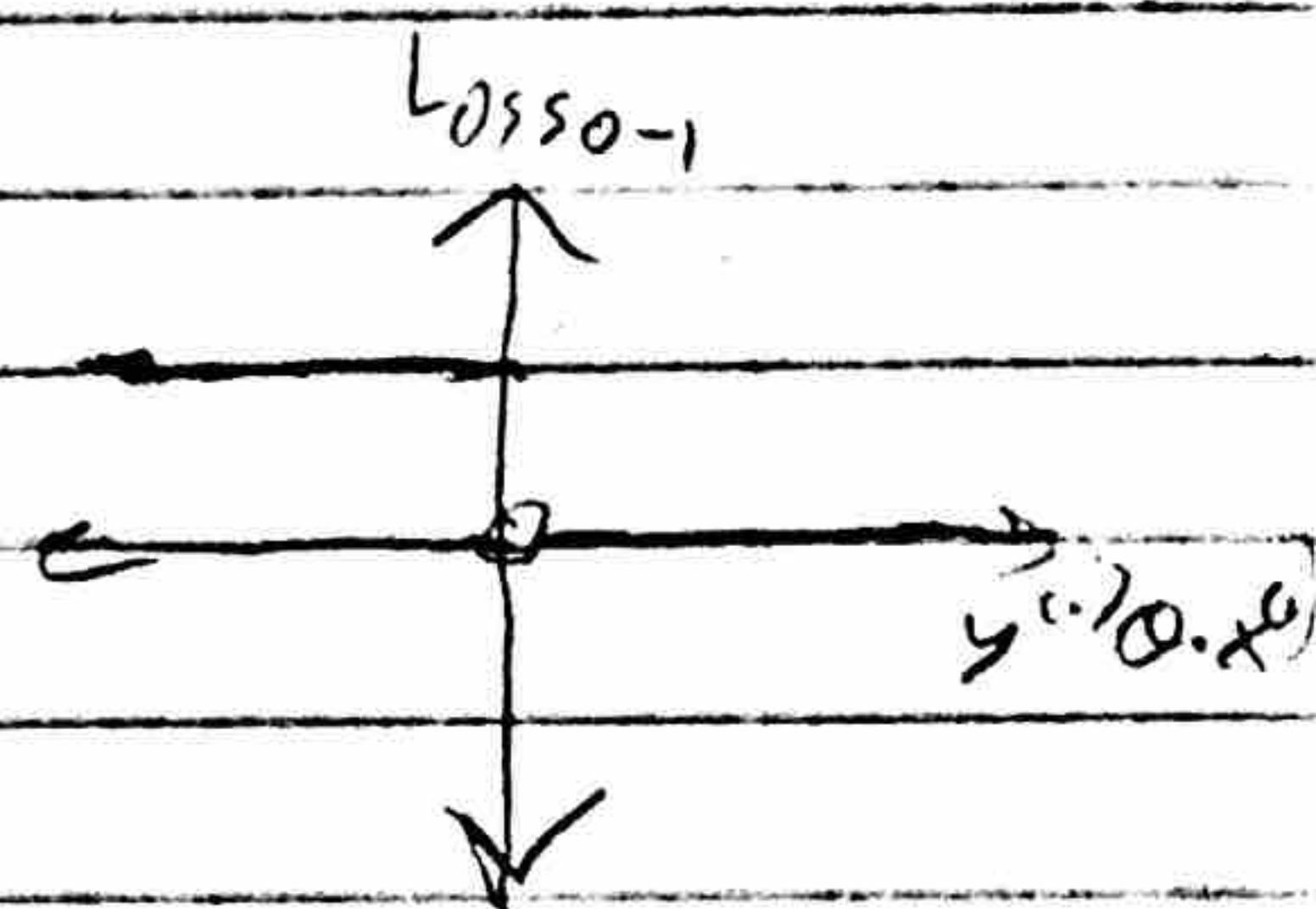
(?) Graph loss

0-1 loss is essentially the inverse of accuracy.

Correct = 0

Incorrect = 1

Zero loss = 100% accuracy



0-1 loss is limited because it doesn't say how close we are to being correct; it only says if we are correct or not (on the right side of the decision boundary or not).

(?) How could we improve the 0-1 loss? Could we include distance information?

Hinge loss

Idea: Make the loss proportional to the distance from a correct classification.

(?) Graph loss

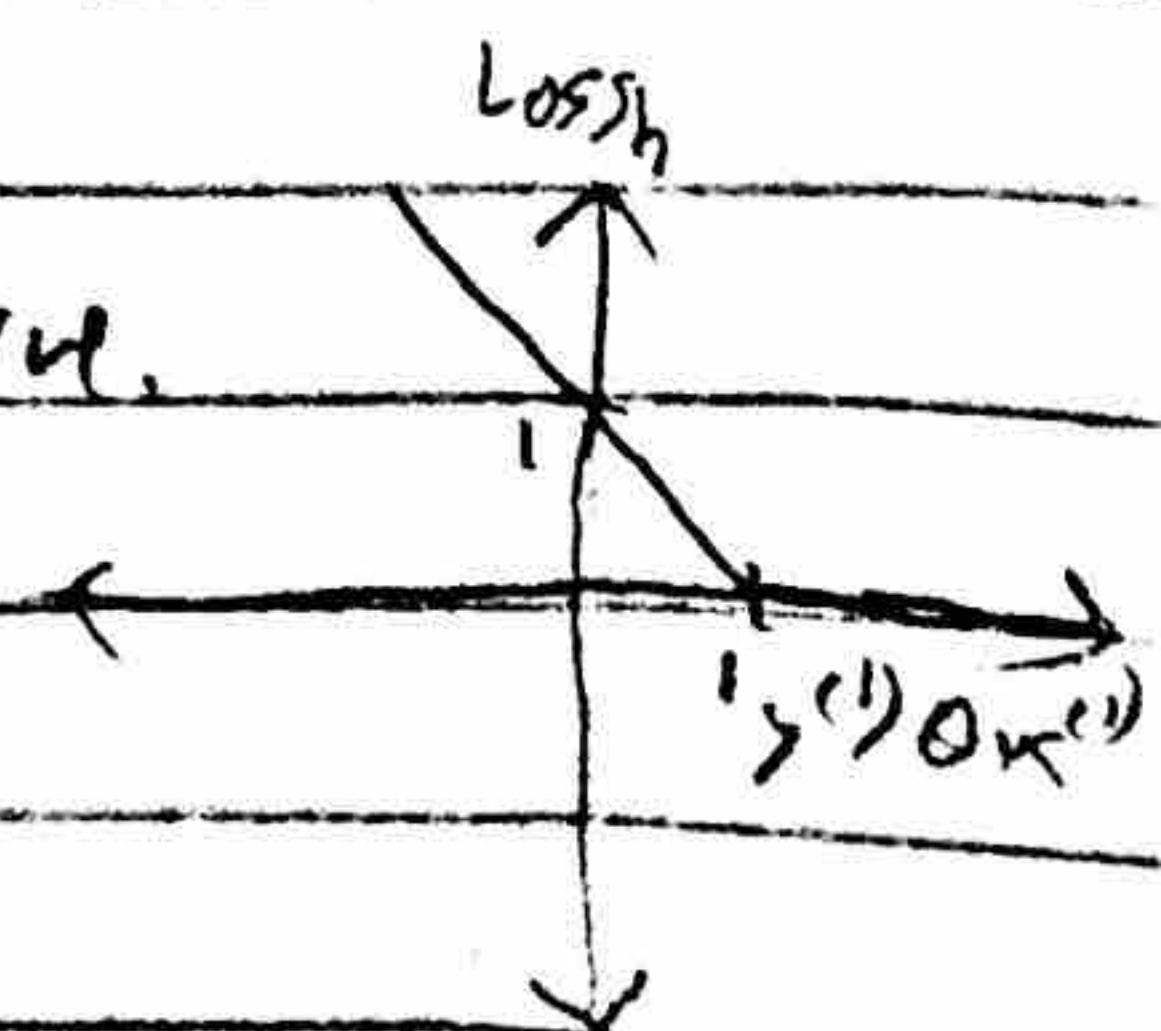
$$\text{Loss}_h(y^{(i)} @ x^{(i)}) = \begin{cases} 1 - y^{(i)} @ x^{(i)} & y^{(i)} @ x^{(i)} \leq 1 \\ 0 & y^{(i)} @ x^{(i)} > 1 \end{cases}$$

$$\text{Loss}_h(y^{(i)} @ x^{(i)}) = \max(0, 1 - y^{(i)} @ x^{(i)})$$

(?) Why 1 instead of 0?

We're using margins rather than the decision boundary.

Enforces stronger separation between positive and negative.



3-6

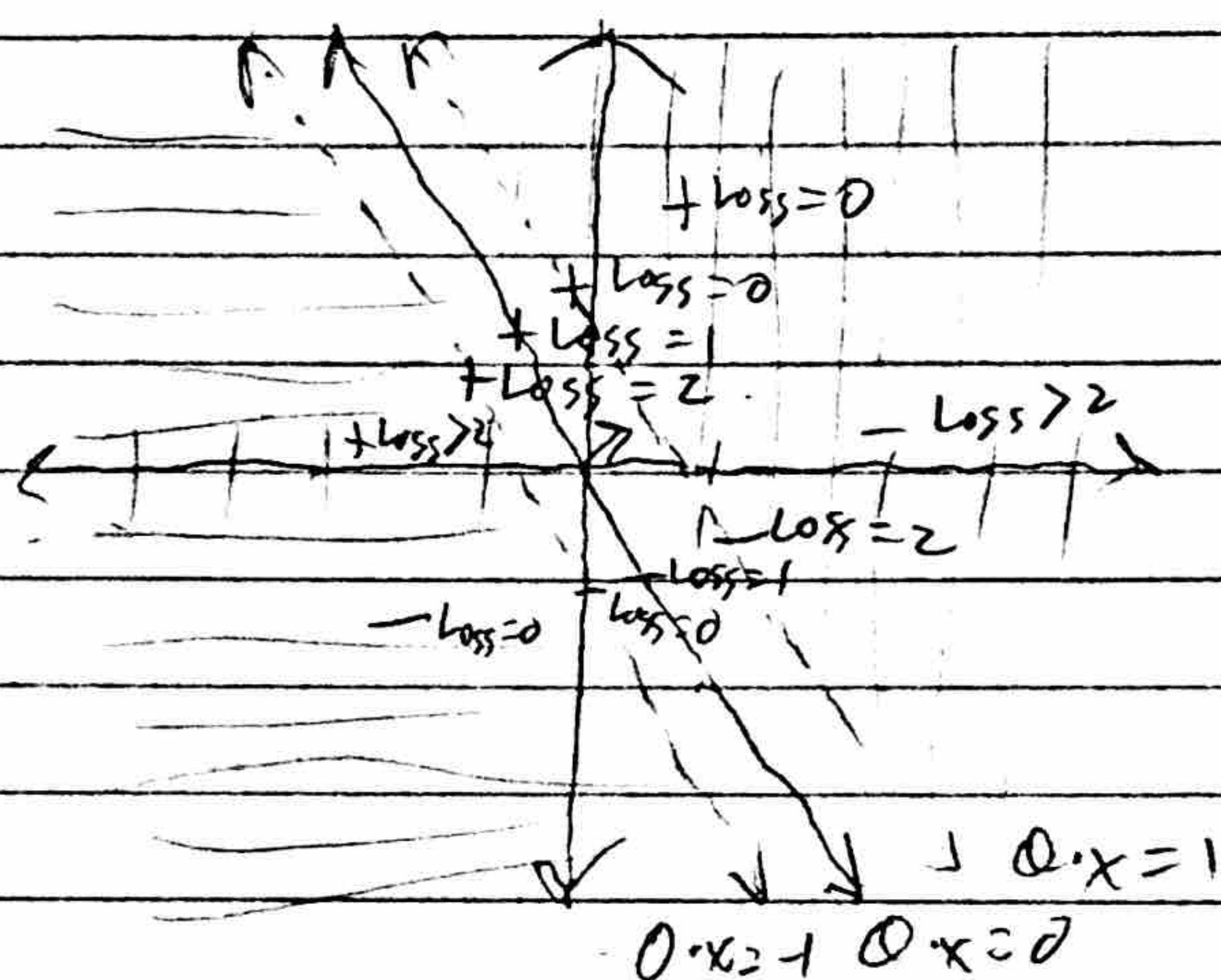
Example

$$Q = \begin{pmatrix} 1 \\ 1/2 \end{pmatrix}$$

$$D = \frac{1}{\|Q\|} = \frac{1}{\sqrt{1^2 + (\frac{1}{2})^2}} = \frac{1}{\sqrt{1 + \frac{1}{4}}} = \frac{1}{\sqrt{\frac{5}{4}}} = \frac{2}{\sqrt{5}} \approx 0.9$$

① What is distance to margin?

② Draw $+$ and $-$ and ask for loss. Show of hands.



□ 0 loss for $+$ point

▢ 0 loss for $-$ point

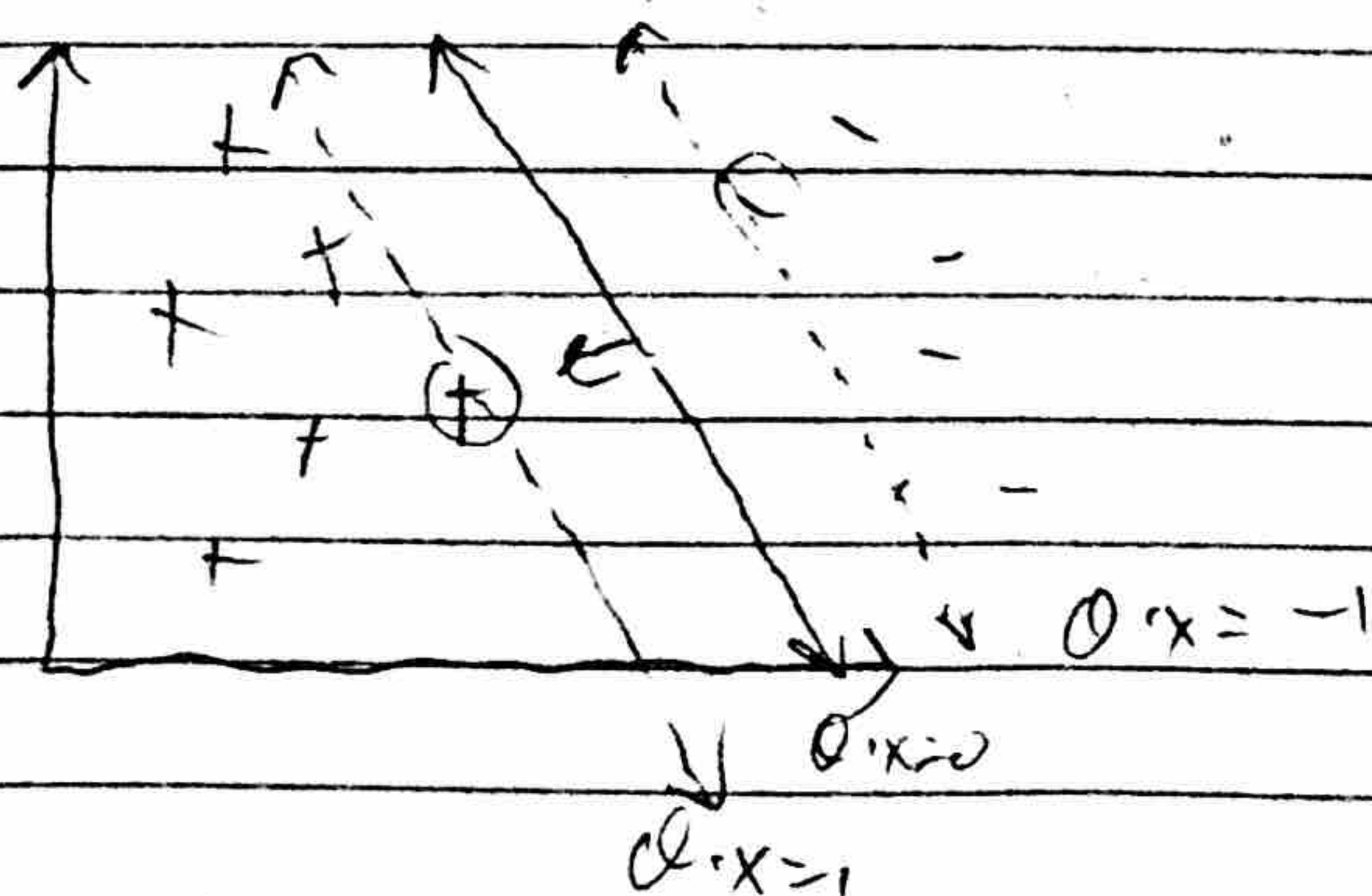
Loss increases linearly from 0 as the point moves away from the margin

③ Realizable

Realizable Case

When the points are linearly separable (realizable), there is a unique solution where hinge loss is 0 and margin is maximized.

④ Ask someone to draw solution.



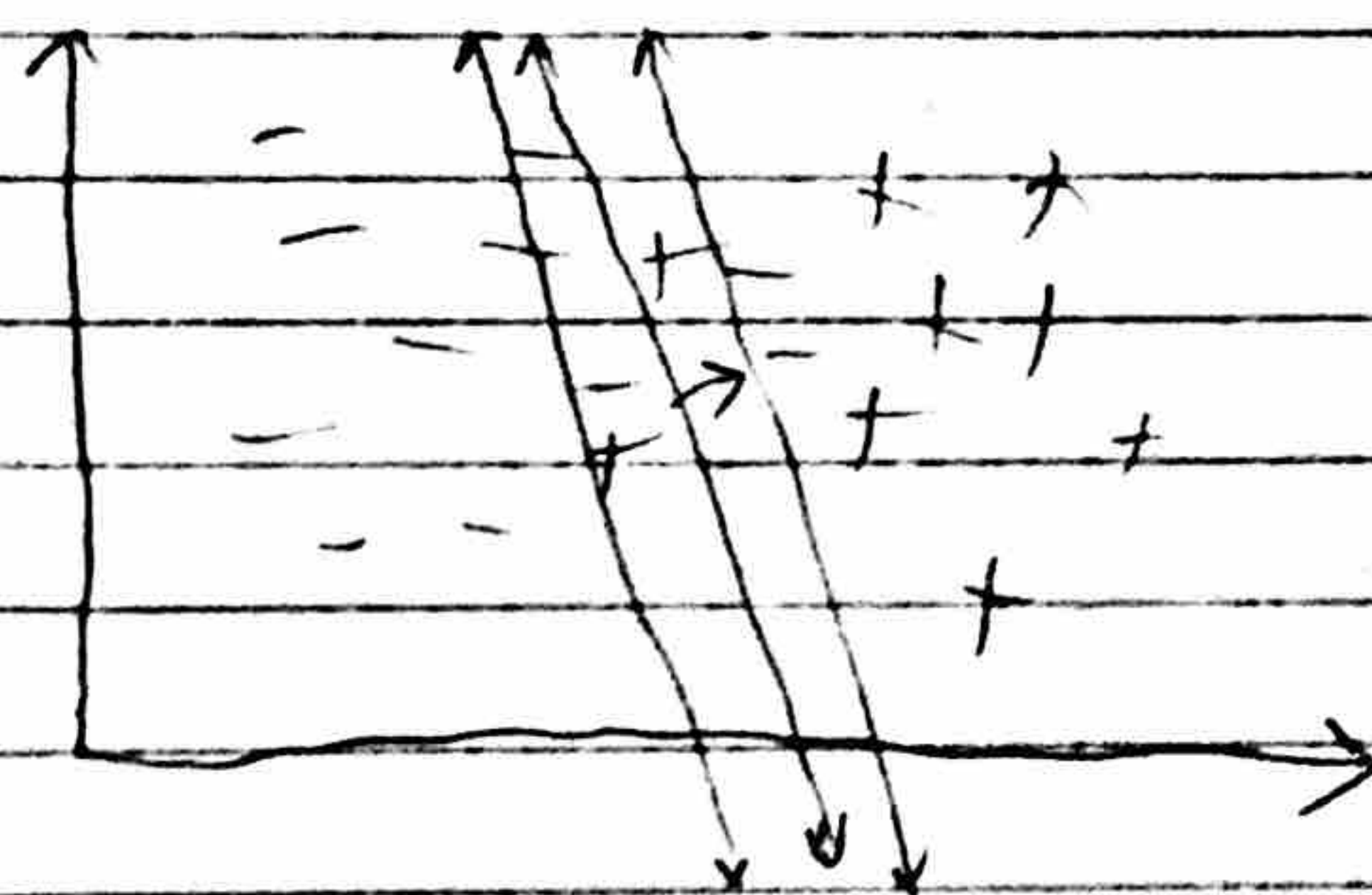
- 0 hinge loss because all points are beyond margin
- Max margin because decision boundary is in the middle between closest positive and negative

(?) What happens if I remove a non-support vector?
What happens if I remove a support vector?

Nope! In the realizable case, the margin is determined entirely by the closest + and - which end up on the margin. These special points are called support vectors. Thus the name SVM.

Unrealizable case

Points can't be perfectly separated so simply balance loss and margin.



Solving SVM

Two methods:

1) offline - Directly optimize the objective function on all the points at once.

Pro: Achieves optimal solution

Con: Hard (involves Lagrangians)

2) online - Optimize the objective function point by point

Pro: Easier

Con: No guarantee of optimal solution

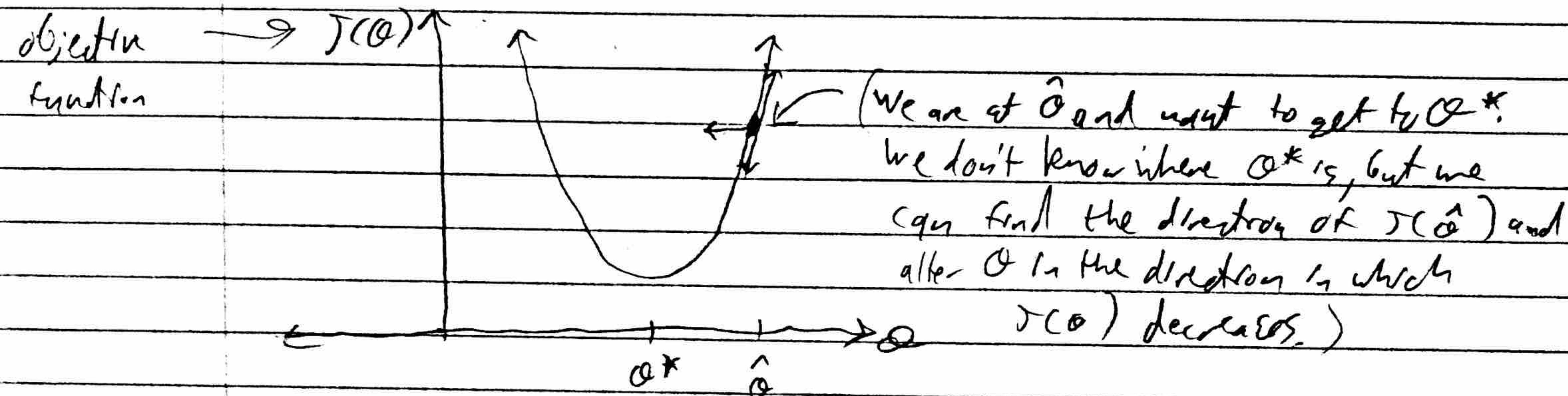
(Today we'll do online SVM because it's easier to understand, it still performs very well, and it introduces concepts which will be important later in the course.)

Online SVM

Idea: Optimize objective function for each point one at a time.

(?) How? If we choose ω just for one point, it's probably going to work very poorly for other points. How can we slightly adjust ω w/o totally overfitting?

(?) We're trying to find the minimum of the objective function, which is a function of θ , but we don't know where the minimum is (i.e. what θ produces minimum cost). Do we know anything about what direction the minimum is in?



Idea: The gradient (generalization of derivative) tells us the direction in which changing θ most increases $J(\theta)$. We will go in the opposite direction.

Gradient Descent

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \nabla_{\theta} J(\theta)$$

$\nabla_{\theta} J(\theta) \rightarrow$ gradient of $J(\theta)$ with respect to θ

$\eta \rightarrow$ learning rate - controls how much we change θ in the direction specified by the gradient

\rightarrow travel in the opposite direction of the gradient

Gradient Review

The gradient is a vector of derivatives.

If $\theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{pmatrix}$ then $\nabla_{\theta} J(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} J(\theta) \\ \frac{\partial}{\partial \theta_2} J(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_d} J(\theta) \end{pmatrix}$

Example

$$J(\theta) = \theta \cdot x \quad \text{where } \theta = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}$$

$$J(\theta) = \theta_1 x_1 + \theta_2 x_2$$

$$(?) \quad \frac{\partial}{\partial \theta_1} J(\theta) = \frac{\partial}{\partial \theta_1} (\theta_1 x_1 + \theta_2 x_2) = 1 \cdot x_1 + 0 \cdot x_2 = x_1$$

$$(?) \quad \frac{\partial}{\partial \theta_2} J(\theta) = \frac{\partial}{\partial \theta_2} (\theta_1 x_1 + \theta_2 x_2) = 0 \cdot x_1 + 1 \cdot x_2 = x_2$$

$$\nabla_{\theta} J(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} J(\theta) \\ \frac{\partial}{\partial \theta_2} J(\theta) \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

Gradient of SVM Objective Function

(ignore $\frac{1}{2}$)

$$\text{For a single point } x^{(i)}, J(\theta) = \text{loss}_h(y^{(i)} \theta \cdot x^{(i)}) + \frac{\lambda}{2} \|\theta\|^2$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} (\text{loss}_h(y^{(i)} \theta \cdot x^{(i)}) + \frac{\lambda}{2} \|\theta\|^2) = \nabla_{\theta} \text{loss}_h(y^{(i)} \theta \cdot x^{(i)}) + \nabla_{\theta} \left(\frac{\lambda}{2} \|\theta\|^2 \right)$$

chain rule

$$(?) \quad \nabla_{\theta} \frac{\lambda}{2} \|\theta\|^2 = 2 \cdot \frac{\lambda}{2} \theta \cdot \frac{\partial}{\partial \theta} \cdot \theta = \lambda \theta$$

$$(?) \quad \nabla_{\theta} \text{loss}_h(y^{(i)} \theta \cdot x^{(i)}) = \nabla_{\theta} \begin{cases} 1 - y^{(i)} \theta \cdot x^{(i)} & , \quad y^{(i)} \theta \cdot x^{(i)} \leq 1 \\ 0 & , \quad y^{(i)} \theta \cdot x^{(i)} > 1 \end{cases}$$

$$= \begin{cases} -y^{(i)} x^{(i)} \\ 0 \end{cases}$$

(combining the pieces)

$$\nabla_{\theta} J(\theta) = \begin{cases} -y^{(i)} x^{(i)} + \lambda \theta & , \quad y^{(i)} \theta \cdot x^{(i)} \leq 1 \\ \lambda \theta & , \quad y^{(i)} \theta \cdot x^{(i)} > 1 \end{cases}$$

Online SVM update step

$$\theta_{\text{new}} = \theta_{\text{old}} - \gamma \nabla_{\theta} J(\theta) = \theta_{\text{old}} - \gamma \begin{cases} -y^{(i)} x^{(i)} + \lambda \theta & , \quad y^{(i)} \theta \cdot x^{(i)} \leq 1 \\ \lambda \theta & , \quad y^{(i)} \theta \cdot x^{(i)} > 1 \end{cases}$$

Now after every point we will move θ in the direction that minimizes the objective function.

We can formalize this process as an algorithm

3-10

Pegasos Algorithm (Online SVM)

$\theta = 0$

repeat T times:

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

if $y^{(i)} \theta \cdot x^{(i)} \leq 1$:

$$\theta = \theta - \eta \cdot (-y^{(i)} x^{(i)} + 1) \theta$$

else:

$$\theta = \theta - \eta \cdot (1) \theta$$

return θ