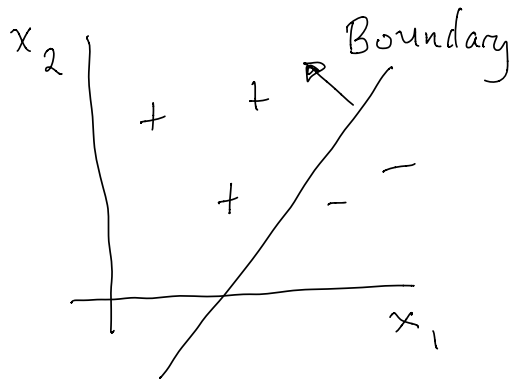CS 181

# Linear Classification



$$D = \{ (\underline{x}_1, y_1), \cdots, (\underline{x}_n, y_n) \}$$

$$\underline{x}_n \in \mathbb{R}^D$$

$$y_n \in \{+1, -1\}$$

Goal: predict $\hat{y}$, on a new example $\underline{x}$

Method 1    Non-parametric

k-Nearest Neighbor            kernel methods

$\rightarrow$ majority vote on the k-NNs

✓ flexible       ✗ Slow on big data set
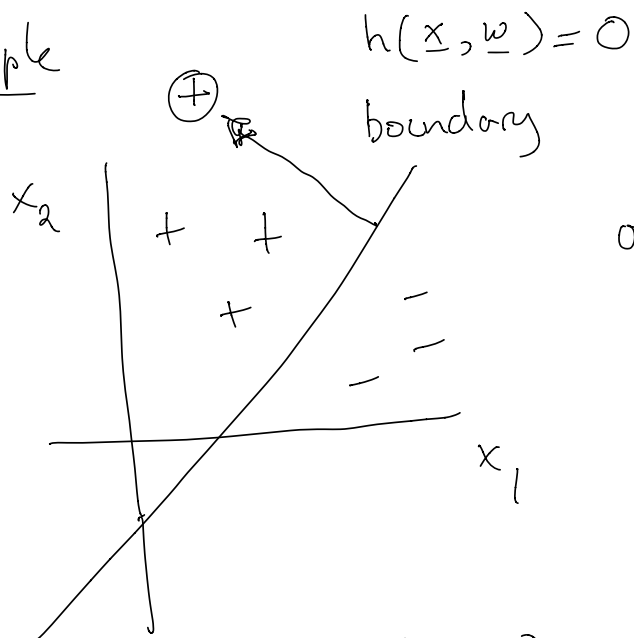                 ✗ Not interpretable

## Method 2  Linear classification

discriminant function
$$h(\underline{x}, \underline{w}) = \underline{w}^T \underline{x} + w_0$$

predict $\hat{y} = \begin{cases} +1 & \text{if } h(\underline{x}, \underline{w}) > 0 \\ -1 & \text{otherwise} \end{cases}$

Example

$h(\underline{x}, \underline{w}) = 0$
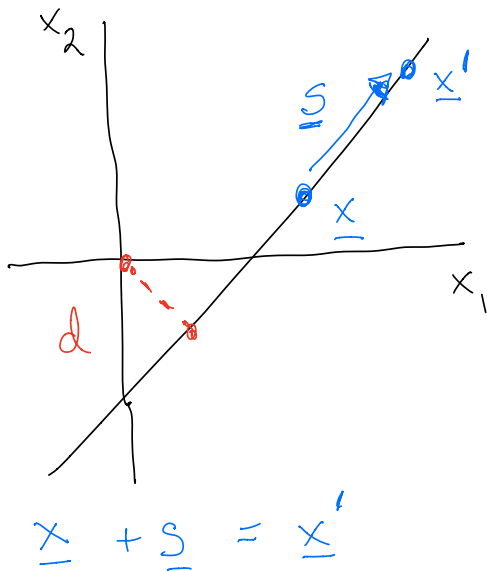
boundary

orthogonal vector

$$\underline{w} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$h(\underline{x}, \underline{w}) = -x_1 + x_2 + 1$$

# Understanding boundary

Generally, hyperplane

$$S = \{ \underline{x} : \underline{\omega}^T \underline{x} + \omega_0 = 0 \}$$

$$\underline{\omega}^T \underline{s} = \underline{\omega}^T \underline{x}' - \underline{\omega}^T \underline{x}$$

$$= \left( \underline{\omega}^T \underline{x}' + \omega_0 \right) - \left( \underline{\omega}^T \underline{x} + \omega_0 \right)$$

$$= 0 - 0 = 0$$

So, $\underline{\omega}$ is orthogonal to the boundary.

$\hookrightarrow$ $\underline{\omega}$ sets the orientation.

$x_2$

$x_1$

$S$

$\underline{x}'$

$\underline{x}$

$d$

$\underline{x} + \underline{s} = \underline{x}'$

$d$ as the unsigned orthonormal distance between boundary + origin

$$d = \frac{|\omega_0|}{\| \omega \|_2}$$

$\hookrightarrow$ $\omega_0$ sets the distance to the origin

## 0/1 Loss function

Recall $\hat{y} = \begin{cases} +1 \\ -1 \end{cases}$ , if $\underline{w}^t \underline{x} + w_0 > 0$
otherwise

For a negative example $(\underline{x}_n, y_n)$

For a positive example



Define

$$\ell_{0/1}(z) = \begin{cases} 1, & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$d_D(\underline{w}) = \sum_n \ell_{0/1}\left(-y_n(\underline{w}^t \underline{x}_n + w_0)\right) = \sum_{n: y_n \neq \hat{y}_n} 1$$
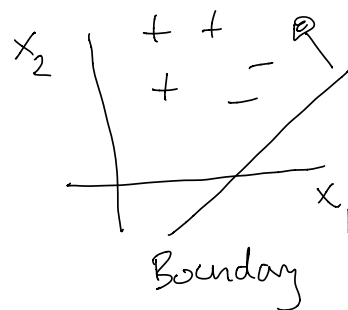
$\hookrightarrow$ scales with # mistakes

### Problem

Want to use gradient descent

$$\underline{w}_{t+1} \leftarrow \underline{w}_t - \eta \frac{\partial d(\underline{w})}{\partial(\underline{w})}$$
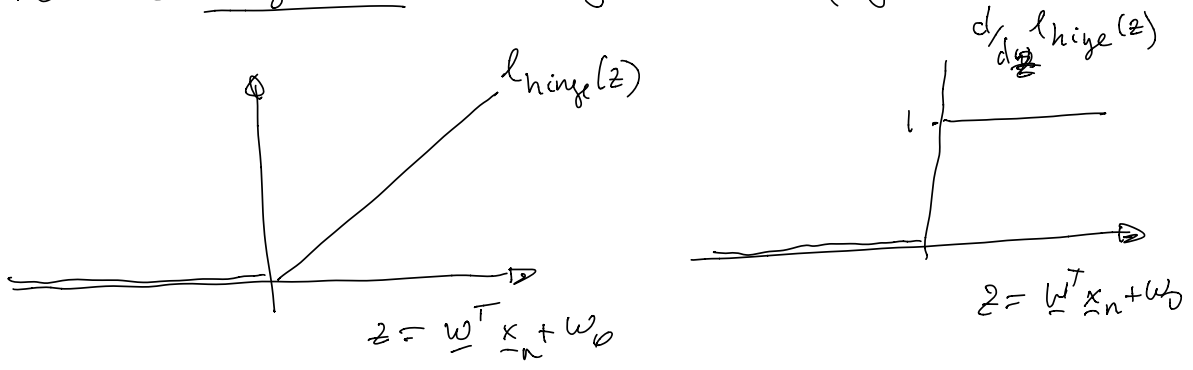
$\eta$ step size $(>0)$



Boundary

But, gradients of 0/1 loss functions are uninformative! Zero almost everywhere.

# Fix: Hinge loss

For a _negative_ training example ($y_n = -1$)



$\ell_{hinge}(z)$

$\frac{d}{dz}\ell_{hinge}(z)$

$z = \underline{w}^T \underline{x}_n + w_0$
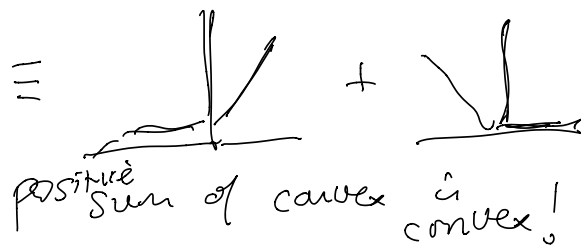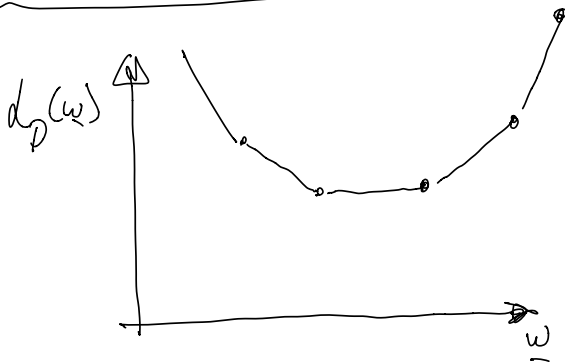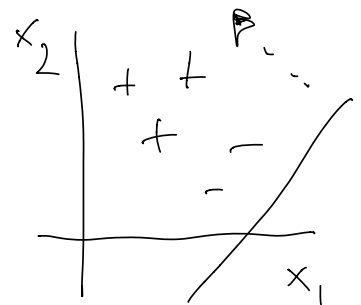
$z = \underline{w}^T \underline{x}_n + w_0$

$$\ell_{hinge}(z) = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$L_D(\underline{w}) = \sum_n \ell_{hinge}\left(-y_n(\underline{w}^t \underline{x} + w_0)\right) = \sum_{n: y_n \neq \hat{y}_n} -y_n(\underline{w}^t \underline{x}_n + w_0)$$

$L_D$ increases with more mistakes & worse mistakes

Gradient descent can now work

$\boxed{L_D \text{ is Convex}}$



$L_D(\underline{w})$

$\underline{w}$

$x_2$

$x_1$

positive sum of convex is convex!

Comments           + differentiable
 ① Convex function can be        } for small
   minimized via grad descent    } enough $\eta$

 ② Convex opt. is polynomial time
   solvable   ( o/w NP-hard )

| Gradient descent | Stochastic gradient descent |
|---|---|
| Repeat: | Repeat: |

Gradient descent:

Repeat:

Update $\underline{w}$ with step
size $\eta$ × gradient

$$\frac{1}{n} \sum_n \frac{\partial}{\partial \underline{w}} \ell_{hinge}(\underline{x}_n, \underline{w}_n)$$

☑ Stable ( fixed $\eta$ ) converges for

☒ Costly per step

Stochastic gradient descent:

Repeat:

Pick random minibatch
$B \subseteq D$ ( perhaps
           $|B| = 1$ )

Update $\underline{w}$ with step
size $\eta$ × gradient on $B$

$$\frac{1}{|B|} \sum_{n \in B} \frac{\partial}{\partial \underline{w}} \ell_{hinge}(\underline{x}_n, \underline{w}_n)$$

☑ Faster on large data

☒ Noisy gradient

(need an adaptive learning
rate ("step size") to
converge; e.g., $\eta_t \propto \frac{1}{t}$  )

Aside

$$\frac{\partial}{\partial \underline{w}} L_{\partial}(\underline{w}) = \frac{\partial}{\partial \underline{w}} \left( - \sum_{n: y_n \neq \hat{y}_n} y_n (\underline{w}^T \underline{x}_n) \right)$$

$$= - \sum_{n: y_n \neq \hat{y}_n} y_n \underline{x}_n$$

SGD on hinge loss $\quad (|B| = 1)$

**Perceptron**

Rosenblatt (1958)

$\eta = 1$

① Pick $(\underline{x}_n, y_n)$ at random

② If $y_n = \hat{y}$, do nothing

Else $\quad \underline{w}_{t+1} \leftarrow \underline{w}_t + \eta \, y_n \underline{x}_n$

Thm  Perceptron converges if and only
if data is linearly separable

## Remark ①

Multiclass classification

$$y_k \in \{ C_1, \ldots, C_K \}$$

"All vs one"

Train $k$ seperate binary classifiers.

Classify a new example as
$$\arg \max_k h(\underline{x}, \underline{w}_k)$$

## Remark 2

Basis function!

## Metrics

True

|  | 0 | 1 |
|---|---|---|
| Predict 0 | TN | FN |
| 1 | FP | TP |

TNR or SPEC

$$\frac{TN}{TN+FP}$$

"reject all negative?"

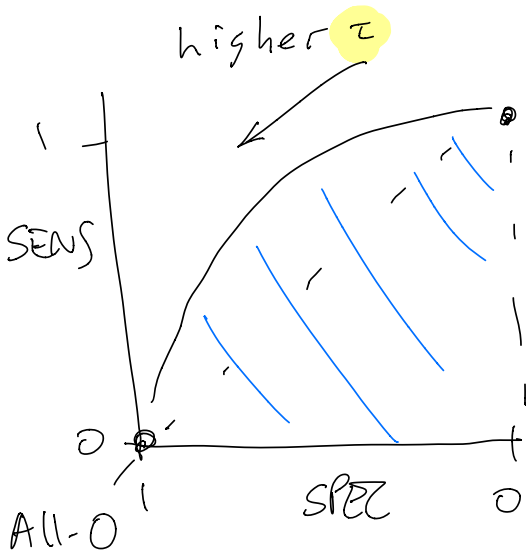TPR or SENS or Recall

$$\frac{TP}{FN+TP}$$

"find all positive?"

ROC curve / AUC

new parameter ↓

Classify $\hat{y} = \begin{cases} +1 & \text{if } h(\underline{w}, \underline{x}) > \text{threshold } \tau \\ -1 & \text{otherwise} \end{cases}$

Generally we've used $\tau = 0$

higher $\tau$



All-1

SENS

SPEC

All-0

"Area under curve"

$AUC \in [0.5, 1]$

Higher better

# F1-score / Prediction + Recall

$$\boxed{\text{PREC}} \quad \boxed{\text{REZALL}}$$

$$\frac{TP}{FP + TP} \qquad \frac{TP}{FN + TP}$$

$$\text{F1-score} = 2 \times \frac{P \times R}{P + R}$$

( harmonic mean )

| TN | FN |
|----|----|
| FP | TP |

Precision



lower threshold

"All 1"

Recall

Precision vs. recall