# FoDA    L25

## (Linear) Classification

↪ Loss functions

# Classification core ML

Input $(X, y)$    $X \subset \mathbb{R}^d$    $y \in \{-1, +1\}$

*supervised problem*    $X = \{x_1, x_2, \ldots x_n\}$    two classes

$X$    attributes of data

$y$    outcome

$$g(x) = \text{sign}(f(x))$$
$$= \begin{cases} +1 & \text{if } f(x) \geq 0 \\ -1 & \text{if } f(x) < 0 \end{cases}$$

Goal: function $f : \mathbb{R}^d \to \{-1, +1\}$    $g(x_i) = y_i$

s.t. on data $(x_i, y_i) \in (X, y)$    $\underline{f(x_i) = y_i}$

on as many as possible

assume $(x_i, y_i) \overset{iid}{\sim} \mu$    *unknown dist*    build $f$ s.t. on

new data $(x, y) \sim \mu$    $f(x) = y$
$$g(x) = y$$

# OCCAMS RAZOR

Simple model tend to generalize better.

---

Restrict $f$ is _linear_ _function_

$$f(x) = \underset{b + w_1}{\underset{\mid}{w_0 + w_1}} x^{(1)} + \underset{w_2}{\underset{\mid}{w_2}} x^{(2)} + \cdots + \underset{w_d}{\underset{\mid}{w_d}} x^{(d)}$$

$$x = (x^{(1)}, x^{(2)}, \ldots x^{(d)})$$
$$\in \mathbb{R}^{d}$$

$$\left. \begin{array}{l} b \in \mathbb{R} \\ w \in \mathbb{R}^d \end{array} \right\} \text{ parameters of the model}$$
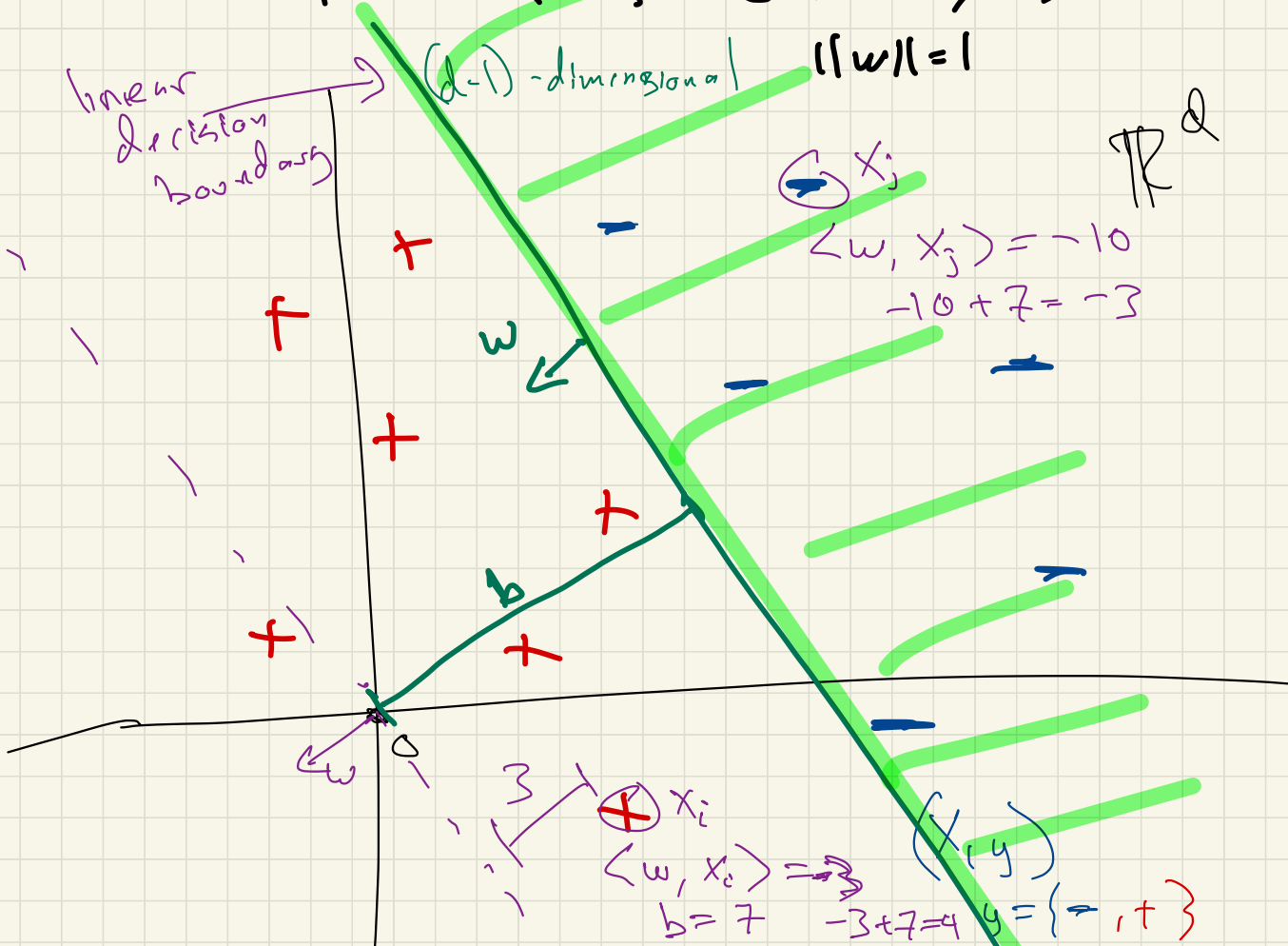
typically $\|w\| = 1$

$$g(x) = \text{sign}(f(x)) \qquad f(x) = b + \langle w, x \rangle$$

$(d-1)$-dimensional

$\|w\| = 1$

linear decision boundary

$\mathbb{R}^d$

$x_j$

$\langle w, x_j \rangle = -10$

$-10 + 7 = -3$

$w$

$b$

$w$

$0$

$3$

$x_i$

$\langle w, x_i \rangle = -3$
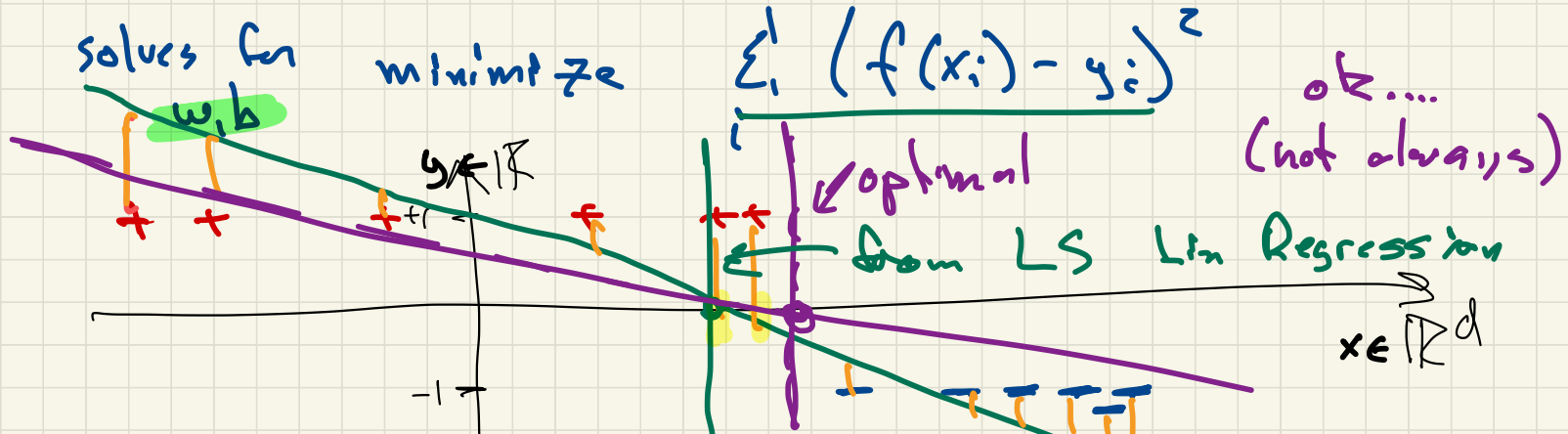
$b = 7 \qquad -3 + 7 = 4$

$(x, y)$

$y = \{-, +\}$

Input $(X, y) \subset \mathbb{R}^d \times \{-1, +1\}$

Goal: Find linear $f_{w,b} : \mathbb{R}^d \to \mathbb{R}$   so $\text{sign}(f(x_i)) = y_i$

How to solve for $w, b$?   $f(x) = b + \langle w, x \rangle$

Apply $\overset{\text{Least Squares}}{\wedge}$ Linear Regression Algorithm.

solves for $\boxed{w, b}$   minimize   $\sum_i \left( f(x_i) - y_i \right)^2$   ok....
(not always)

optimal

from LS Lin Regression



$y \in \mathbb{R}$

$x \in \mathbb{R}^d$

# Can we use GD to minimize?

$$g_\alpha = \text{sign}(f_\alpha(-))$$

$$\Delta\left(g_\alpha(x_{i,5})\right) = \sum_{i=1}^{n}\left(1 - \mathbb{1}\left(\text{sign}(g_i) = \text{sign}(f_\alpha(x_i))\right)\right)$$

$$g_i := g_\alpha(x_i)$$

$$= \# \text{ misclassified points.}$$

identity function

$$\mathbb{1} : \{\text{true, false}\} \to \{0,1\}$$

$$\mathbb{1} = \begin{cases} 1 & \text{if True} \\ 0 & \text{if false} \end{cases}$$

## No: can't use GD

— not convex

— no gradient!

$$\Delta\left(g_\alpha(x_{i,y})\right)[\alpha] \quad \text{piece-wise} \quad \text{constant}$$

# Loss Functions  (approximate $\Delta$)

$$f(\alpha) = \mathcal{L}\left(g_{\alpha,}(x_{i,5})\right) = \sum_{i=1}^{n} \ell\left(g_{\alpha,}(x_{i,}, 5_i)\right)$$

bivariate

$$= \sum_{i=1}^{n} \ell_\alpha(z_i)$$

univariate

$$\boxed{z_i = y_i \cdot g_\alpha(x_i)}$$

$$y_i \, g_\alpha(x_i) = \begin{cases} > 0 & \text{if } y_i < 0 \\ & g_\alpha(x_i) < 0 \\[4pt] > 0 & \text{if } y_i > 0 \\ & g_\alpha(x_i) > 0 \\[8pt] < 0 & \text{otherwise} \\ & (\text{bad prediction}) \end{cases}$$

linear

$$= \sum_{i=1}^{n} f_i(\alpha)$$

$f$ is decomposible

$$\boxed{f_i(\alpha) = \ell_\alpha(z_i = y_i \cdot g_\alpha(x_i))}$$
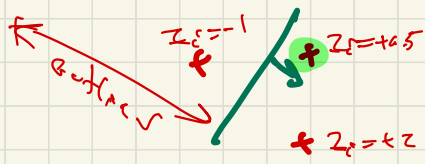
$\ell_\alpha$ is good if
① convex   ② has gradient   ③ approximate $\Delta$

$\Delta(z) = \begin{cases} 0 & \text{if } z \geq 0 \\ 1 & \text{if } z < 0 \end{cases}$

surface $z_i = -1$    + $z_i = +.5$    $z_i = y_i \cdot g_\alpha(x_i)$

+ $z_i = +2$

good if $> 0$, ok bad

## loss functions

hinge loss    $\ell(z) = \max\{0, 1-z\}$

hinge

squared hinge    $\ell(z) = \left(\max\{0, 1-z\}\right)^2$

## logistic loss

$\ell(z) = \ln\left(1 + \exp(-z)\right)$

smoothed hinge

$\Delta$ ReLU

smoothed hinge
$\ell(-z) = \begin{cases} 0 & z \geq 1 \\ (1-z)^3/2 & \text{if } 0 \leq z < 1 \\ \frac{1}{2} - z & \text{if } z < 0 \end{cases}$

logistic    $z_c > 0$

logistic    $z_c > 0$

classified incorrectly    ← → classified correctly    $z_i \in \mathbb{R}$