

# Lecture 18: Regression Basics

## 1. Regression tasks

1) Data:  $X = \{x_1, x_2, \dots, x_n\} = \{x_i\}_{i=1}^n$  a set of  $n$  data samples

$$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{im} \end{pmatrix} \text{ feature vector (ith)}$$

$y = \{y_i\}_{i=1}^n$  targets / ground truth

$y_i$ : continuous value, e.g.,  $y_i \in \underline{\mathbb{R}}$   $(x^2)' = 2x$   
 $(\frac{1}{2}x^2)' = x$

$$x_i \rightarrow \boxed{f} \rightarrow \hat{y}_i \text{ (estimation)} \longleftrightarrow y_i$$

$\Rightarrow$  model:  $f$ : if  $f$  is a linear function,  $f$  is a linear model.

3) Loss function,

$$\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2, \quad \frac{1}{2n} \sum_{i=1}^n (\hat{y}_i - y_i)^2, \quad \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$



## 2. Linear Regression

Model:  $f = w^T \cdot x_i + \underbrace{w_0}_{\text{bias}} = 1 \cdot w_0 + w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3} + \dots + w_m x_{im}$

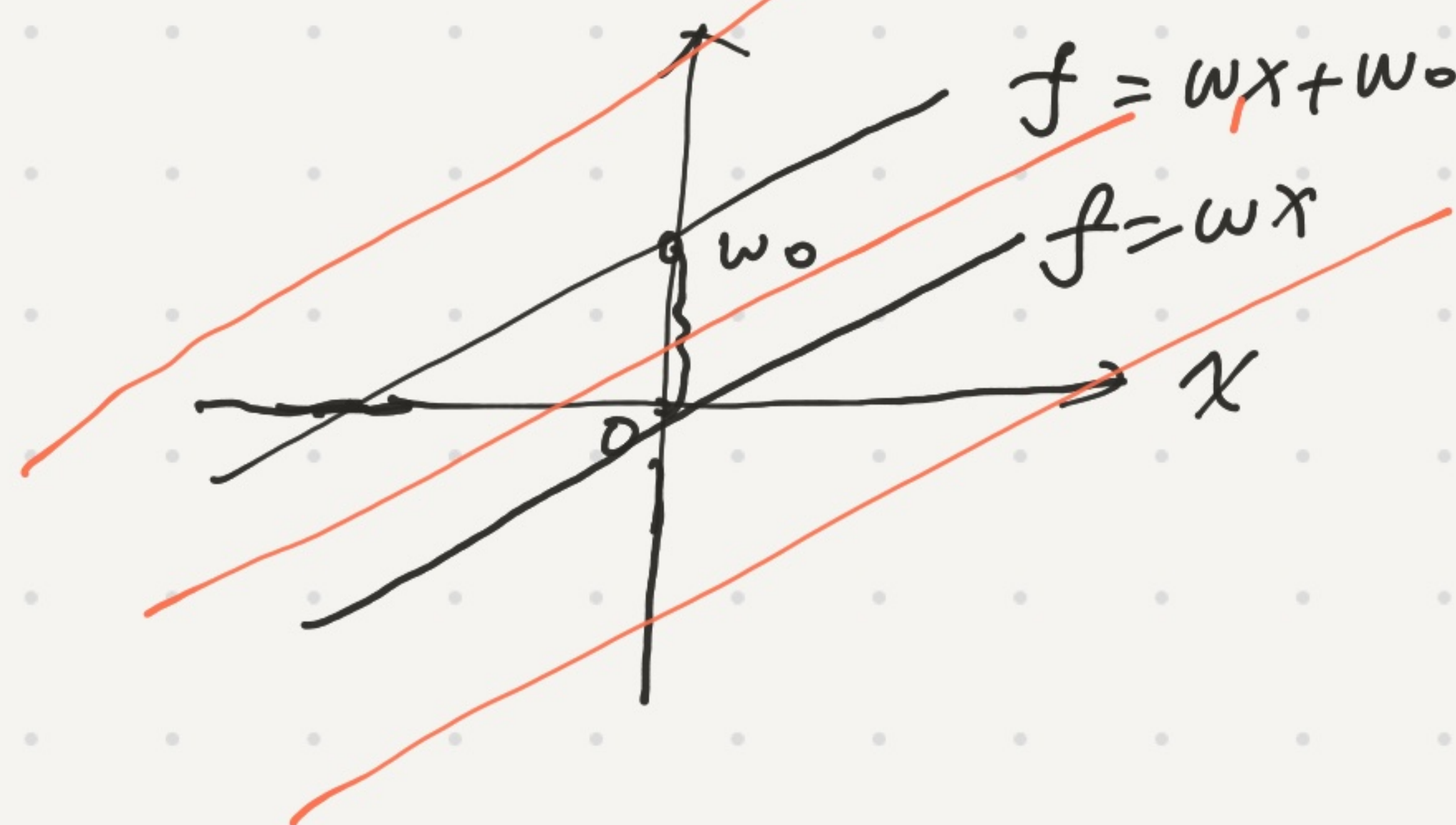
$$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{im} \end{pmatrix} \quad w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{pmatrix} \Rightarrow \text{defines the contribution of each feature.}$$

General form:  $w = \begin{pmatrix} w_0 \\ \vdots \\ w_m \end{pmatrix}_{(m+1)}$   $x_i = \begin{pmatrix} x_{i0} = 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{im} \end{pmatrix}_{(m+1)}$

$$f = w^T \cdot x_i$$

$w_0$ : bias / intercept

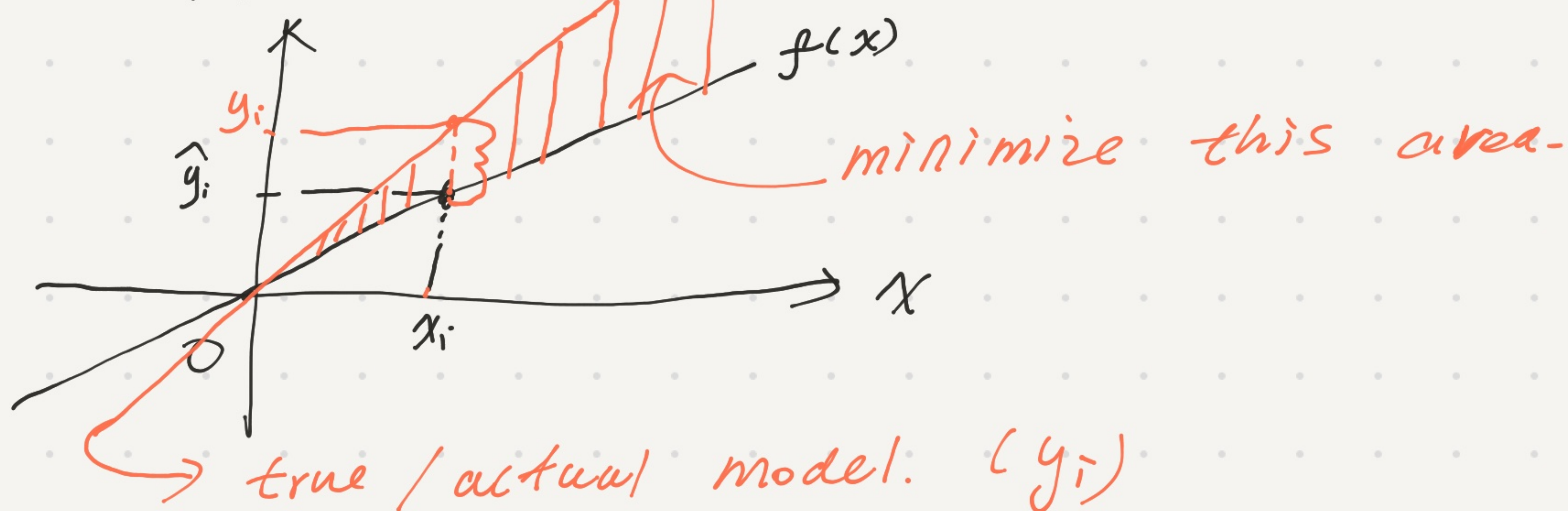
$w_1, w_2, \dots, w_m$ : Coef





Loss:

$$L = \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \sum_{i=1}^n (f(x_i) - y_i)^2$$



## 2.1. Ordinary Least Squares (standard Linear Regression) (OLS)

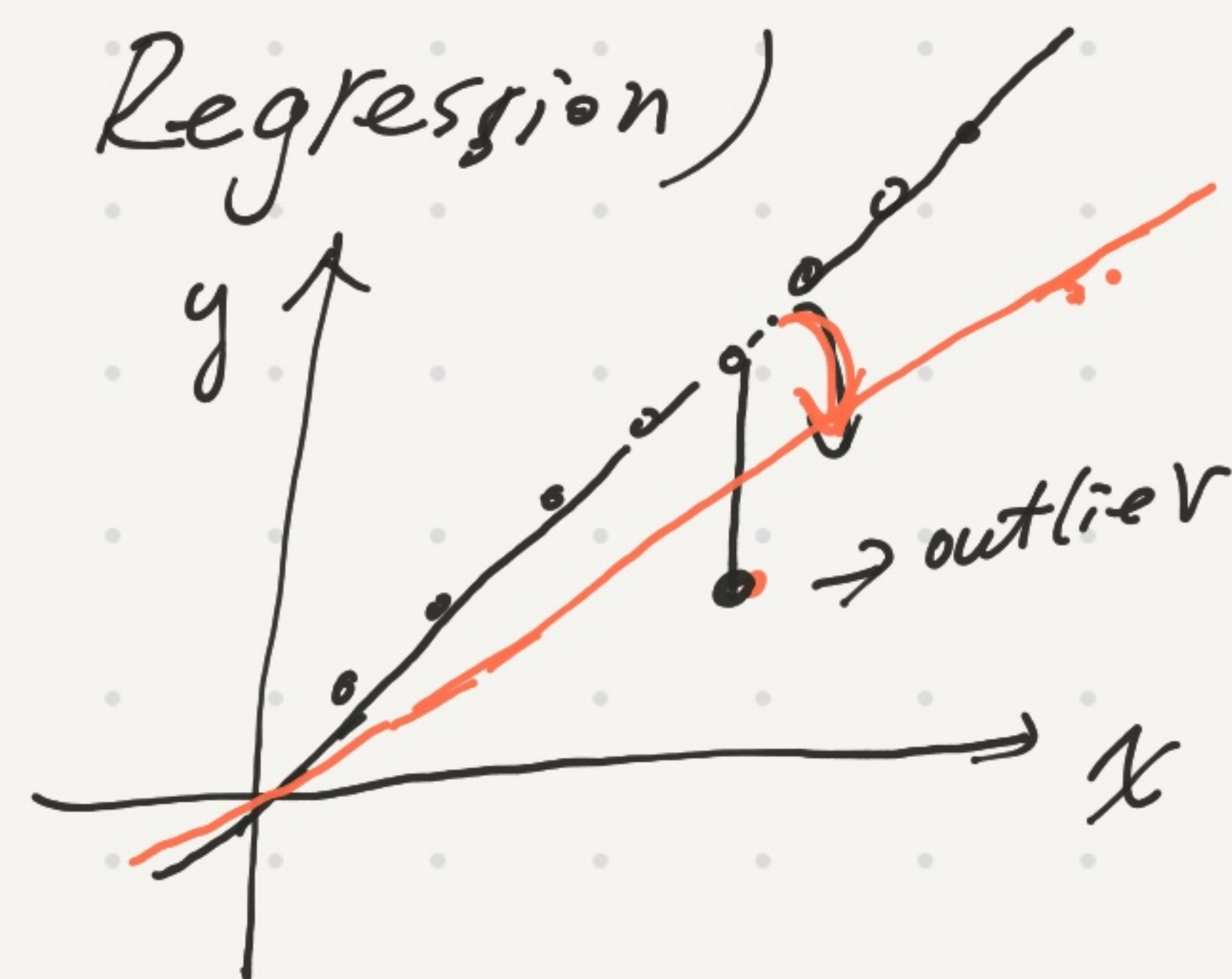
$$\text{Loss} : J_1 = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Optimization:  $\omega^* = \arg \min_{\omega} J_1$

All distances  $(\hat{y}_i - y_i)^2$  contribute equally.

→ Cause a major challenge in OLS: OLS sensitive to

outliers (noise) ⇒ { → remove outliers  
→ regularization { Ridge  
LASSO





### 3. Regularization: strategy used to model user preference

No-free-Lunch theorem: there is no machine learning algorithm that can generalize well to all application.  
good test performance.

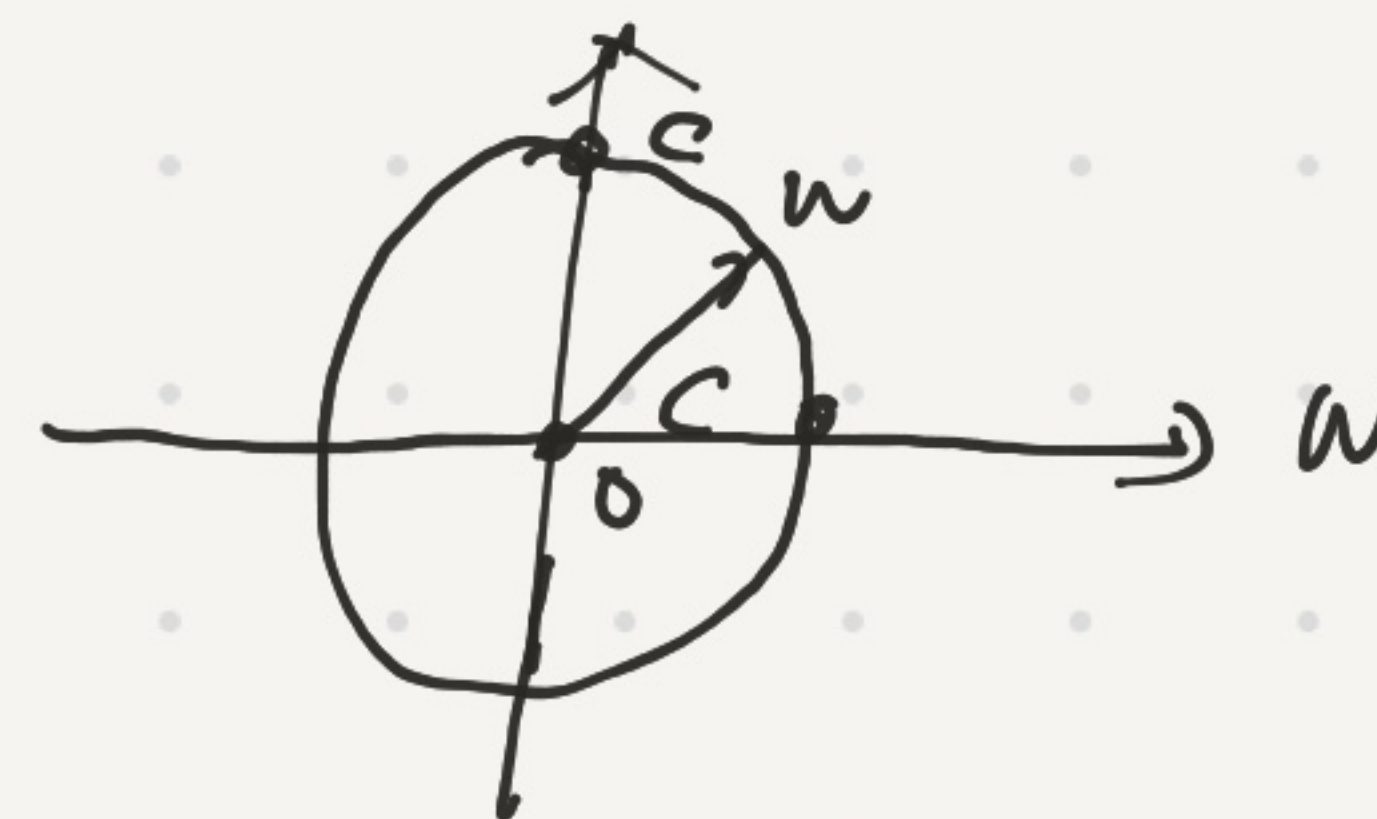
① LS:  $\min J_1 = \sum_{i=1}^n (\hat{y}_i - y_i)^2$

↓ define constraints to add preference.

↓  $\min J_1 = \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \sum_{i=1}^n (\underline{\underline{w}}^T x_i - y_i)^2$

subject to (s.t.):  $\|\underline{\underline{w}}\|_2^2 \leq C$

↓ constraints

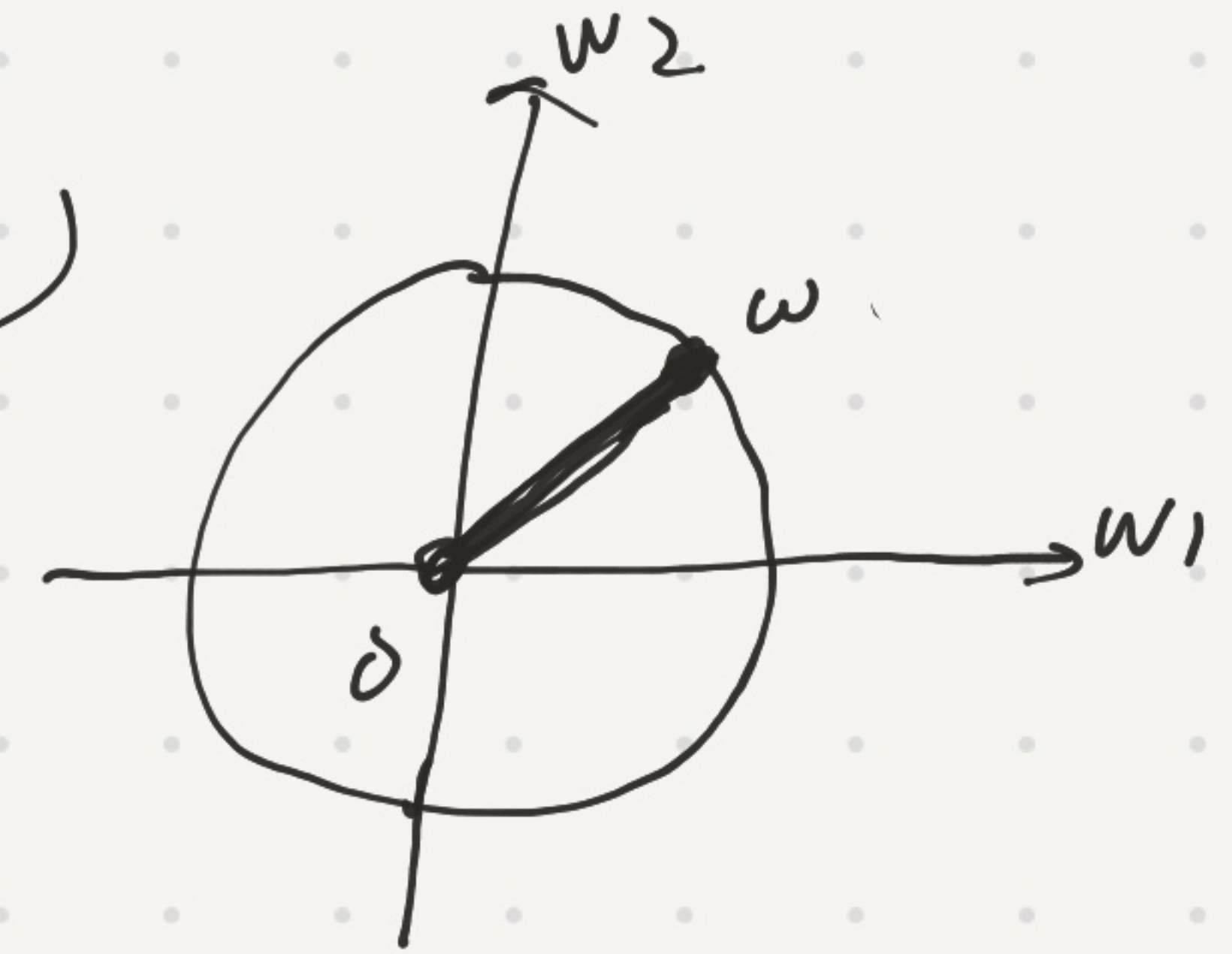




In ML, we transform all constraints to Regularizers which are terms in the loss function.

Ridge Regression

$$A) J_2 = J_1 + \underbrace{2 \|w\|_2^2}_{\substack{\text{regularizer} \\ \text{length of } w}} \rightarrow (L_2 \text{ norm})$$



Our preference in this example is small  $w$ .

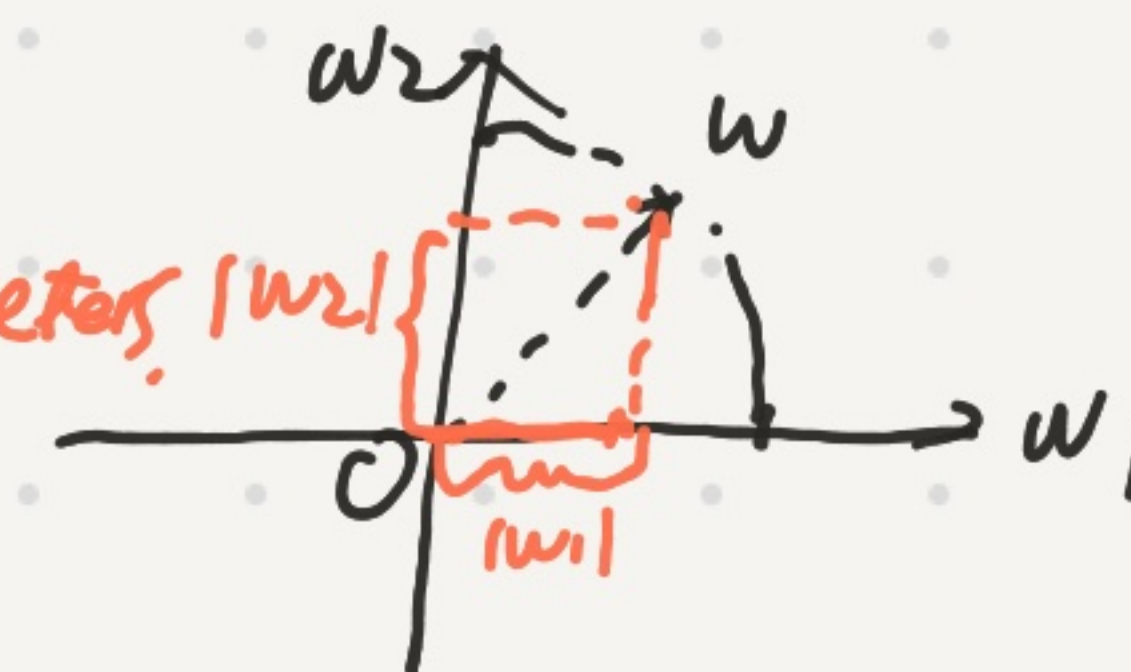
with this  $\|w\|_2^2$ , ridge regression will shrink the model parameters.

B) Least absolute shrinkage and selection operator (LASSO)

$$J_3 = J_1 + 2 \|w\|_1 \rightarrow L_1 \text{ norm.}$$

$$\|w\|_1 = |w_1| + |w_2| + \dots \quad (w_m) = \sum_{i=1}^m |w_i|$$

LASSO will produce more zero parameters





### C) Elastic Net

$$J_4 = J_1 + \underbrace{\alpha_1 \|W\|_2^2}_{L_2} + \underbrace{\alpha_2 \|W\|_1}_{L_1}$$

$\alpha_1, \alpha_2$  : constants, define the contribution of each term.

In a loss function, if we have  $M$  terms, we will  $M-1$  constants. (e.g.,  $\alpha_1, \alpha_2, \dots, \alpha_{M-1}$ )

$\left\{ \begin{array}{l} \text{Model parameters: } w_0, w_1, \dots, w_m \quad (f = w^T x) \\ \text{Hyperparameters: } \alpha_1, \alpha_2. \quad (\text{not related to } f) \end{array} \right.$



#### 4. Problems in Linear Regression.

$$f(x_i; w) = w_0 + w_1 \cdot x_{i1} + w_2 \cdot x_{i2} + \dots + w_m \cdot x_{im}$$

1) the independent assumption: All features are independent

$x_i = \begin{cases} x_{i1} = \text{height} \\ x_{i2} = \text{weight} \end{cases}$  the standard linear model ( $f$ ) cannot model the dependence.

We can add interaction terms in model to model dependency

$$f = w_0 + \underbrace{x_{i1} \cdot w_1 + w_{i2} \cdot x_{i2}}_{\text{interaction term}}$$

$$f_1 = w_0 + x_{i1} \cdot w_1 + w_{i2} \cdot x_{i2} + \boxed{w_3 \cdot x_{i1} \cdot x_{i2}}$$

If we have 3 features:  $x_{i1}, x_{i2}, x_{i3}$ .

$x_{i1} \cdot x_{i2}, x_{i2} \cdot x_{i3}, x_{i1} \cdot x_{i3}, x_{i1} \cdot x_{i2} \cdot x_{i3}$

interaction term between  $x_{i1}$  and  $x_{i2}$ .



2) Non-linear relationship between  $x$  and  $y$ .

add non-linear terms into model function.

Polynomial model:

Linear model:  $f = w_0 + w_1 \cdot \underline{x_{i1}} + w_2 \underline{x_{i2}}$ .



polynomial model:  $f = \underbrace{w_0 + w_1 x_{i1} + w_2 x_{i2}}_{\text{Linear terms}} + \underbrace{w_3 x_{i1}^2 + w_4 x_{i2}^2 + w_5 x_{i1} x_{i2}}_{\text{non-linear terms}}$  → with degree  $\geq$

Polynomial function / model with degree 3:

degree 3:  $x_{i1}^3, x_{i2}^3, x_{i1}^2 \cdot x_{i2}, x_{i1} \cdot x_{i2}^2$

degree 2:

degree 1 and 0:

