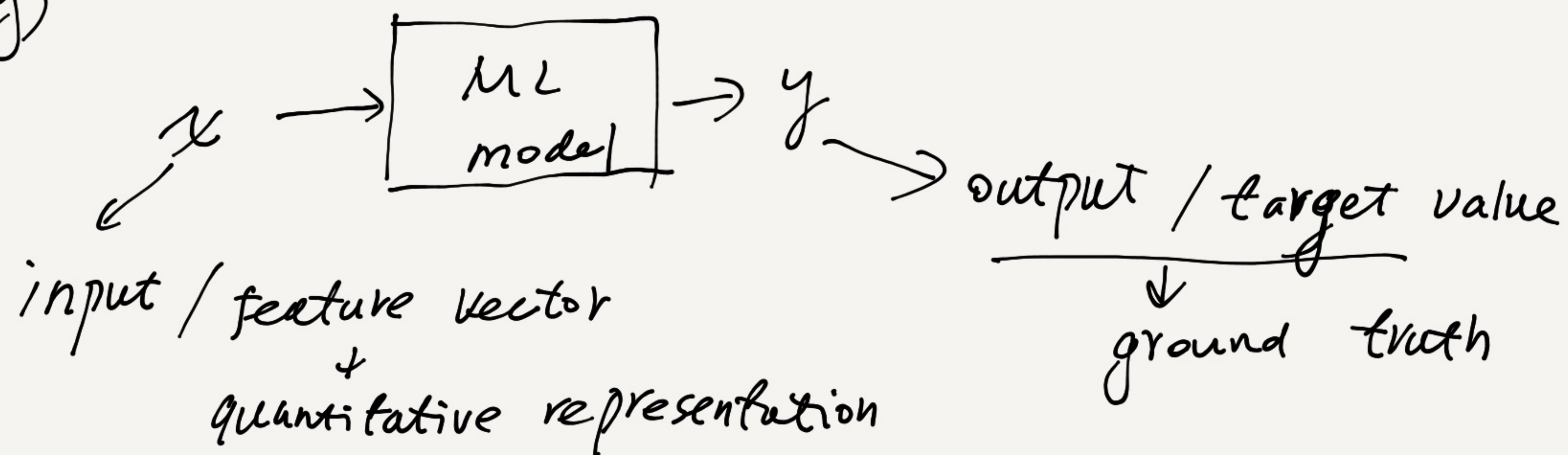# Lecture 5. ML Algorithm Basics

## 1. Two Categories of ML algorithms

1) supervised learning: we learn ML models to associate the input with output

$(x)$ $(y)$

$$x \longrightarrow \boxed{\begin{array}{c} ML \\ model \end{array}} \longrightarrow y$$

input / feature vector
↓
quantitative representation

output / target value
↓
ground truth

Linear regression
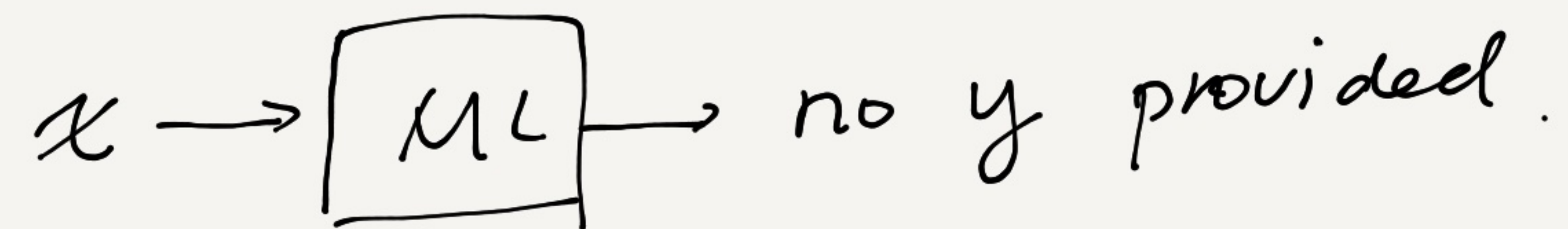Logistic regression

Naive Bayes

Support Vector machines (SVMs)

Decision trees

K-nearest neighbor (KNN)

* Artificial Neural Networks (ANNs)

2) Unsupervised learning. We only have input feature vectors

$$x \longrightarrow \boxed{ML} \longrightarrow \text{no } y \text{ provided.}$$

① Dimensionality reduction (DR)

high-dimensional data

$$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{im} \end{pmatrix} \longrightarrow \boxed{DR} \longrightarrow \widetilde{x_i} = \begin{pmatrix} \widetilde{x}_{i1} \\ \widetilde{x}_{i2} \\ \vdots \\ \widetilde{x}_{iK} \end{pmatrix} \quad \underline{K << m}$$

$m$ is large

$\Big\{$ Principal component analysis (PCA) $\rightarrow$ linear DR

manifold learning $\longrightarrow$ a set of algorithms for non-learn DR.
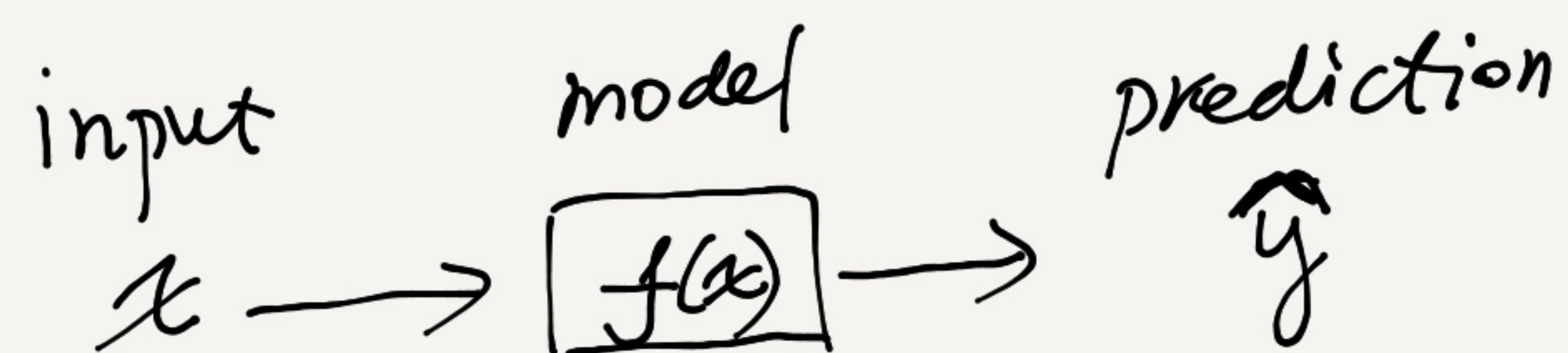
② Clustering. group data. e.g., k-means

## 2. Four Components in ML

① Data sets : $\underset{\text{find best model}}{\underbrace{\overset{(60\%-70\%)}{\text{training set}}}}$ + $\underset{\substack{\text{determine the} \\ \text{hyperparameters}}}{\underbrace{\overset{(15\%-20\%)}{\text{validation set}}}}$ ← training stage.

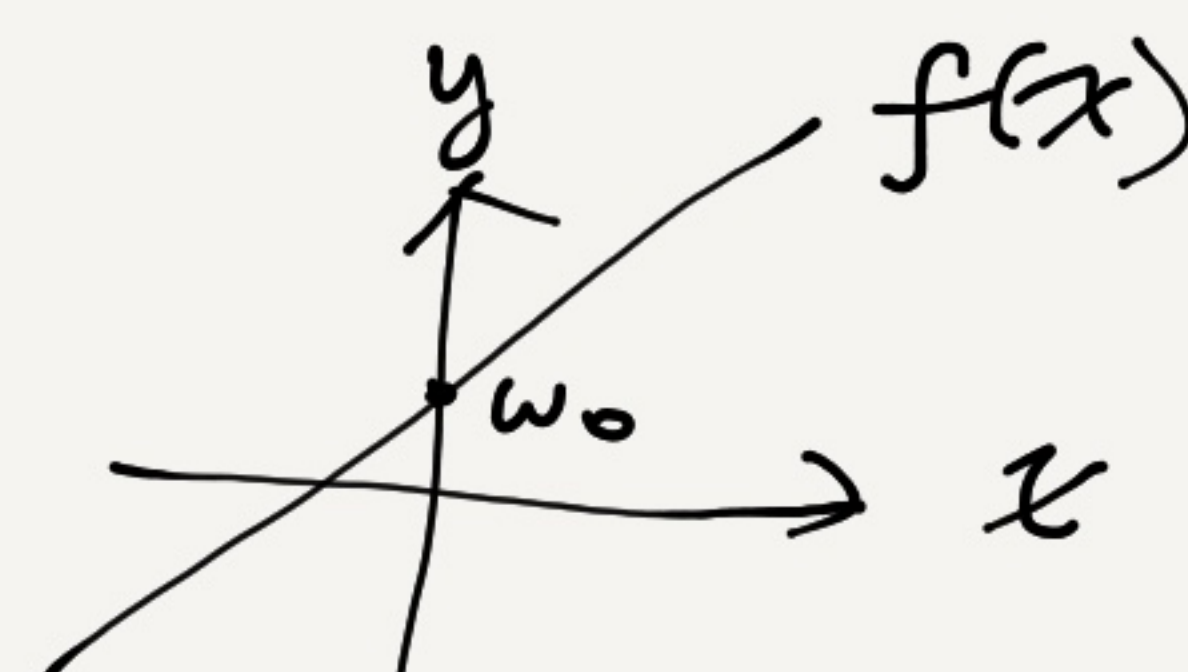$\underset{}{\underline{\overset{(15\%-20\%)}{\text{test set}}}}$ : evaluate generalization performance,

i.e., performance on new data

② Model : defines a function to compute prediction from input

input    model    prediction

$x \longrightarrow \boxed{f(x)} \longrightarrow \hat{y}$

simple model : linear model.

1D : $f(x) = w_0 + w_1 x$

2D : $f(x) = w_0 + w_1 x_1 + w_2 x_2$   $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$
         plane

MD :

$f(x) = w_0 + w_1 x_1 + \cdots + w_m x_m$

$= w^T x + w_0$

$w = \begin{pmatrix} w_1 \\ \vdots \\ w_m \end{pmatrix}$   $x = \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}$

3D : $f(x) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_4$
     hyperplane         $x = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$
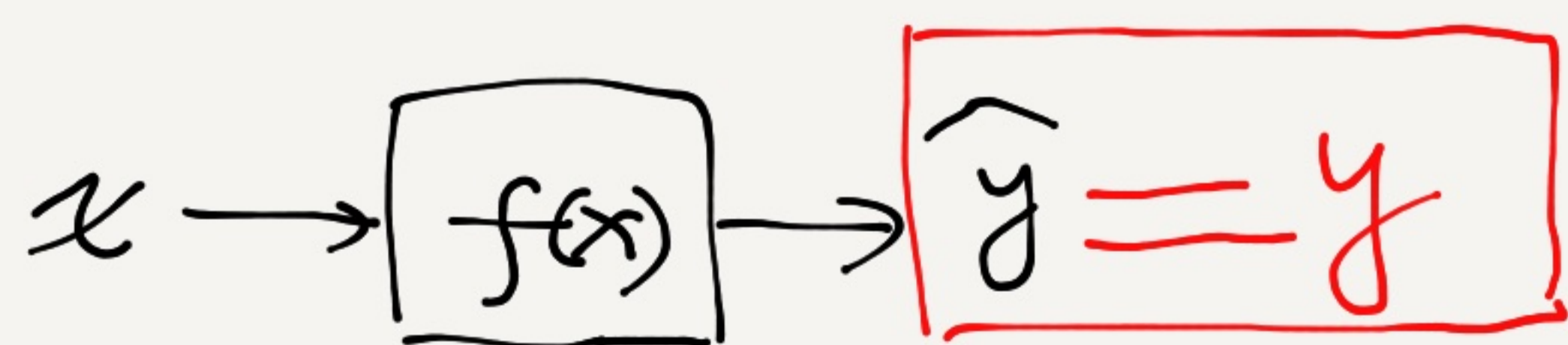
**generic form of linear function.**

$f(x) = w^T \textcircled{x}$

$w = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{pmatrix}_{m+1}$   $x = \begin{pmatrix} x_0 = 1 \\ x_1 \\ \vdots \\ x_m \end{pmatrix}$

3) Loss/cost/objective function : defines training goal, e.g., prediction should ($\hat{y}$) match the target value (y)

$$x \longrightarrow \boxed{f(x)} \longrightarrow \boxed{\hat{y} = y}$$

① Exact match $\quad L_0 = \begin{cases} 0 & \text{if } \hat{y} = y \\ +\infty & \text{if } \hat{y} \neq y \end{cases}$

residual : $y - \hat{y}$

② Residual-based Loss

$$L_1 = \sum_{i=1}^{n} (y_i - \hat{y_i})^2 .$$

$n$: # of data samples.
residual sum of squares.
RSS.

$$L_2 = \frac{1}{2} \sum_{i=1}^{n} (y_i - \hat{y_i})^2$$

$$L_3 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y_i})^2 \; : \; \text{mean square error (MSE)}$$

$$L_4 = \frac{1}{2n} \sum_{i=1}^{n} (y_i - \hat{y_i})^2$$

$$L_5 = \sum_{i=1}^{n} y_i \cdot \log \hat{y_i} \quad \longrightarrow \; \text{cross-entropy}$$

**4). Optimization algorithm / optimizer / solver**

defines steps to find $f(x)$ that minimize the loss function

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

$$\omega^* = \arg\min_{\omega} \frac{1}{2} \sum_{i=1}^{n} \overbrace{(\hat{y}_i - y_i)^2}$$   $\omega^*$ : best model parameters

$$\downarrow \hat{y} = f(x) = \omega^T x \quad \text{(linear model)}$$

$$= \arg\min_{\omega} \frac{1}{2} \sum_{i=1}^{n} (\boxed{\omega^T} x_i - y_i)^2$$

$\mathcal{L}$

<u>How ??</u>  The extrem points have zero gradient.

We can solve $\nabla \mathcal{L}|_{\omega=\omega^*} = 0$ to

get $\omega^*$.

$$\boxed{\nabla \mathcal{L}|_{\omega=\omega^*} = 0}$$



$w^*$

$\rightarrow \omega$

$$\mathcal{L}_0 = \omega^T x - y \longrightarrow \nabla_\omega \mathcal{L}_0 = x \longrightarrow \boxed{\text{Obtain } \omega^*}$$

$$\mathcal{L}_1 = (\omega^T x - y)^2 \longrightarrow$$