

2*

Linear Regression

[Lecture Notes]

$x^{(i)}$ \Rightarrow i^{th} input feature in training set

$y^{(i)}$ \Rightarrow i^{th} output in training set

$(x^{(i)}, y^{(i)})$ \Rightarrow i^{th} training example

$\{(x^{(i)}, y^{(i)}) \mid 1 \leq i \leq m\}$ \Rightarrow training set

X \Rightarrow denotes the Space of input features.

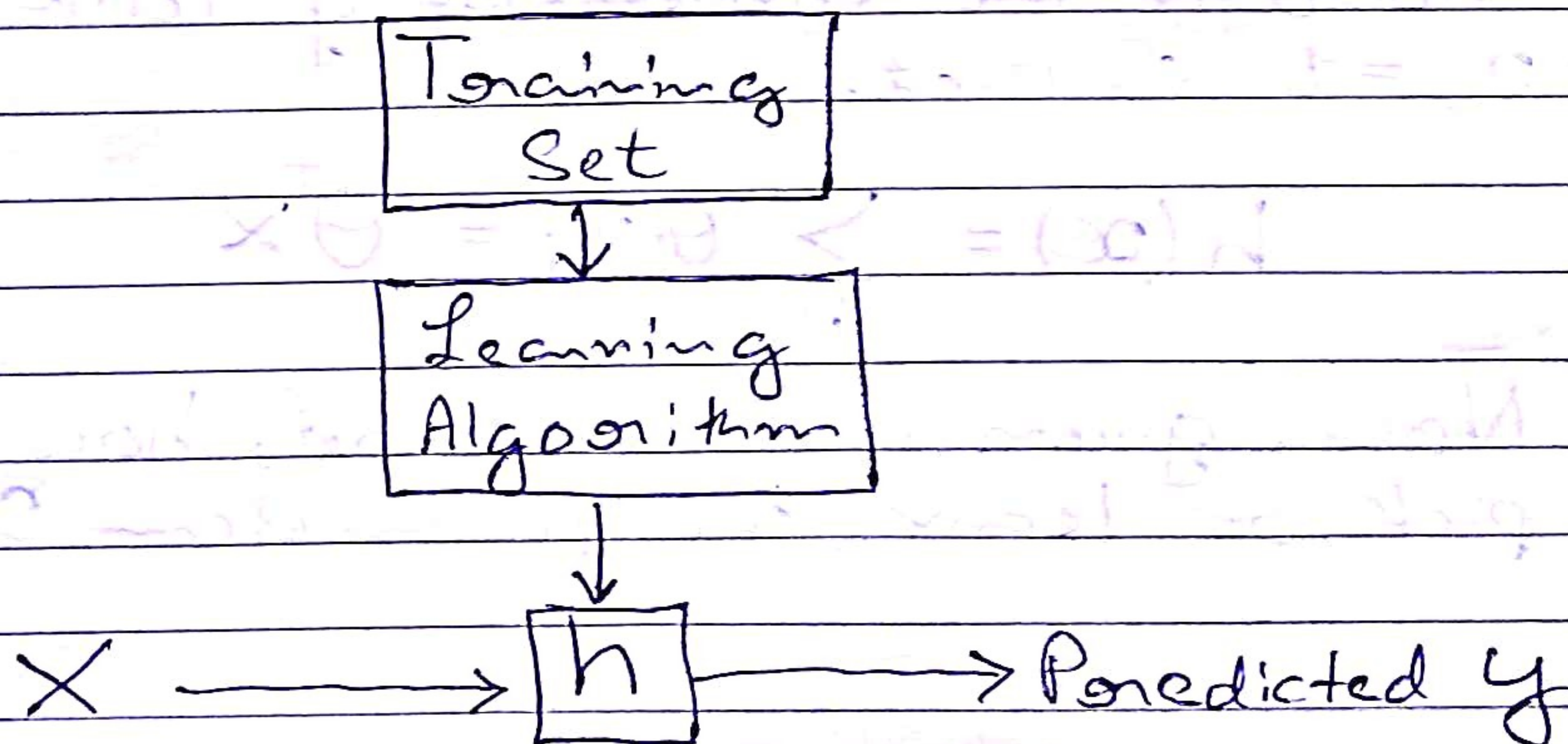
Y \Rightarrow denotes the Space of output

m \Rightarrow # training example

* Supervised Learning

"Given a training set, learn a function $h: X \mapsto Y$, so that $h(x)$ is a "good" prediction for the corresponding value of y "

\hookrightarrow For historical reasons, this function h is called a hypothesis.



Linear Regression

⇒ To perform supervised learning, we must decide how we're going to represent functions/hypotheses h in a Computer.

⇒ As an initial choice, let's say we decide to approximate y as a linear function of x :

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$\forall x \in \mathbb{R}^2$$

⇒ Here, the θ_i 's are the parameters parameterizing the space of linear functions mapping X to Y .

⇒ To simplify our notation, we also introduce the convention of letting $x_0 = 1$ so that,

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

⇒ Now, given a training set, how do we pick, or learn, the parameter θ ?

⇒ One reasonable method seems to be to make $h(x)$ close to y , at least for the training examples we have.

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

★ LMS algorithm

⇒ We want to choose θ so as to minimize $J(\theta)$.

⇒ To do so, let's use a search algorithm that starts with some initial "guess" for θ and that repeatedly changes θ to make $J(\theta)$ smaller, until hopefully we converge to a value of θ that minimizes $J(\theta)$.

⇒ Let us consider the gradient descent algorithm.

⇒ Start with some initial θ , and repeatedly perform the update

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

for $j = 0, \dots, n$

$\alpha \Rightarrow$ Learning rate

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_j} h_{\theta}(x^{(i)})$$

$$\frac{\partial}{\partial \theta_j} \sum_{k=0}^m \theta_k x_k^{(i)} \rightarrow x_j^{(i)}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

⇒ This rule is called the LMS update rule and is also known as the **Widrow-Hoff learning rule**.

⇒ This method looks at every example in the entire training set on every step, and is called **batch gradient descent**.

⇒ The cost function $J(\theta)$ is a **convex quadratic function**, so it only has one global optimum.

↳ Hence gradient descent will always converge to global optimum.

⇒ There is an alternative to batch gradient descent that also works very well.

for $i=1$ to m }

$$\theta_j := \theta_j - \alpha (h_0(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \left\{ \begin{array}{l} \text{FOR} \\ \text{EVERY} \\ j \end{array} \right\}$$

}

⇒ This algorithm is called **stochastic gradient descent**.

⇒ Batch gradient descent has to scan through the entire training set before taking a single step.

{ A costly operation if m is large }

⇒ Stochastic gradient descent can start making progress right away, and continue to make progress with each example it looks at.

⇒ On the down side, Stochastic gradient descent may never converge to the minimum, and the parameters θ will keep oscillating around the minimum of $J(\theta)$.

↳ But in practice most of the values near the minimum will be reasonable good approximations to the true minimum.

★ The normal equations

SAME

