

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

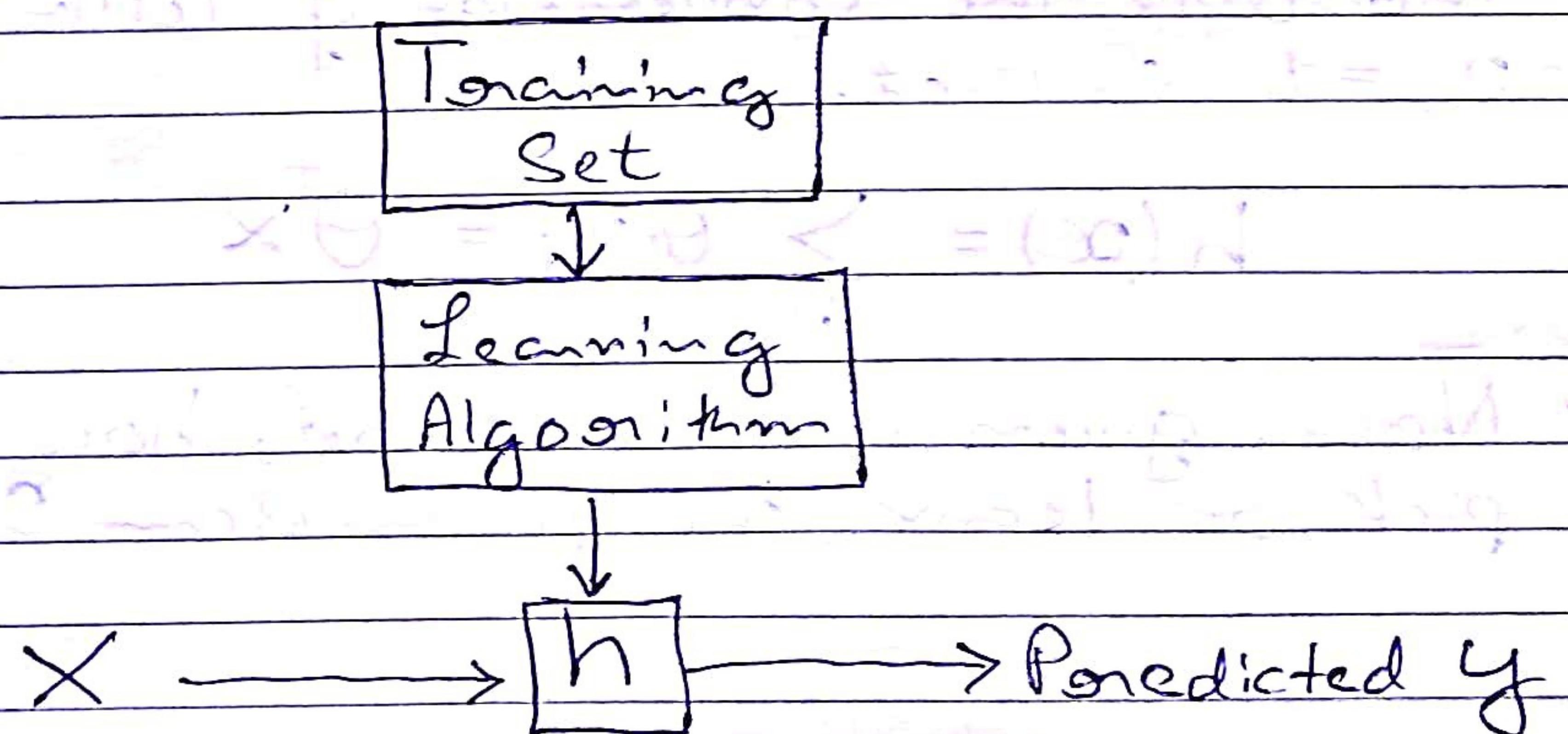
$Y \Rightarrow$ denotes the Space of Output

$m \Rightarrow$ # training example

* Supervised learning

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

↳ 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 parameters θ ?

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

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

★ 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 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

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

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

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

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

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

⇒ 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 {

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

}

⇒ 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.

* Normal equation (Least Square Solution)

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

Let $X = \begin{bmatrix} x^{(1)\top} \\ x^{(2)\top} \\ \vdots \\ x^{(m)\top} \end{bmatrix}$ & $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$

(m × n+1) , (n+1 × 1)

Then, $\begin{bmatrix} h(x^{(1)}) \\ h(x^{(2)}) \\ \vdots \\ h(x^{(m)}) \end{bmatrix} = X\theta$

Let $y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$

Then, $J(\theta) = \frac{1}{2} (X\theta - y)^T (X\theta - y)$

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (x\theta - y)^T (x\theta - y) \\ &= \frac{\partial}{2} (\theta^T x^T x \theta - 2y^T x \theta + y^T y) \\ &= x^T x \theta - 2x^T y = 0\end{aligned}$$

$\Rightarrow \theta = (x^T x)^{-1} x^T y$

~~X~~ ~~X~~