

Objectives:

- Write the training error as least squares criterion for linear regression
- Use stochastic gradient descent for fitting linear regression models.
- Solve closed-form linear regression solution
- Identify regularization term and how it changes the solution, generalization

Introduction:

Classification: $S_m = \{ (x^{(t)}, y^{(t)}) / t=1, \dots, m \}$
 $x^{(t)} \in \mathbb{R}^d, y^{(t)} \in \{-1, +1\}$

In many problems, it is not enough to say yes or no kind of questions.

Regression: $f(x, \theta, \theta_0) = \sum_{i=1}^d \theta_i x_i + \theta_0 = \theta \cdot x + \theta_0$ and $y^{(t)} \in \mathbb{R}$

Empirical Risk:

$$R_m(\theta) = \frac{1}{m} \sum_{t=1}^m \frac{(y^{(t)} - \theta x^{(t)})^2}{2}$$

\uparrow
m training examples

$\frac{1}{m}$? We want to for every single point in our training data and compute this extent of this deviation, compute some kind of loss.
 We can talk about how to define loss, sum them up, and then average

$\frac{(y^{(t)} - \theta x^{(t)})^2}{2}$? If the deviation is large, we're really penalized, and this is the behavior you are getting from the squared function, that the bigger difference would actually result in much higher loss.

Structural mistakes: Maybe the mapping between your training vectors and y 's is actually highly non-linear. You would be incurring a high mistake.

Estimation mistakes: Even if we know that the mapping itself is linear, but you have very limited training data, you cannot estimate it correctly.

Empirical Risk: Hinge loss vs Squared Error loss

$$(x^{(1)}, y^{(1)}) \quad x^{(1)} = [1, 0, 1]^T \quad y^{(1)} = 2$$

$$(x^{(2)}, y^{(2)}) \quad x^{(2)} = [1, 1, 1]^T \quad y^{(2)} = 2.7 \quad \text{and} \quad \theta = [0, 1, 2]^T$$

$$(x^{(3)}, y^{(3)}) \quad x^{(3)} = [1, 1, -1]^T \quad y^{(3)} = -0.7$$

$$(x^{(4)}, y^{(4)}) \quad x^{(4)} = [-1, 1, 1]^T \quad y^{(4)} = 2$$

$$R_m(\theta) = \frac{1}{n} \sum_{t=1}^n \underbrace{\text{loss}(y^{(t)} - \theta \cdot x^{(t)})}_z$$

empirical risk

$$\text{Hinge loss: } \text{loss}_h(z) = \begin{cases} 0 & \text{if } z \geq 1 \\ 1-z & \text{otherwise} \end{cases}$$

$$\text{Squared Error loss: } \text{loss}_s(z) = \frac{z^2}{2}$$

$R_m(\theta)$ with Hinge loss:

$$R_m(\theta) = \frac{1}{n} \sum_{t=1}^n \text{loss}_h(y^{(t)} - \theta \cdot x^{(t)}) \quad \text{with } \text{loss}(z) = \begin{cases} 0 & \text{if } z \geq 1 \\ 1-z & \text{otherwise} \end{cases}$$

$$z^{(1)}: y^{(1)} - \theta \cdot x^{(1)} = 2 - [0, 1, 2]^T \cdot [1, 0, 1]^T = 2 - (0 \cdot 1 + 1 \cdot 0 + 2 \cdot 1) = 0 < 1$$

$$z^{(2)}: y^{(2)} - \theta \cdot x^{(2)} = 2.7 - [0, 1, 2]^T \cdot [1, 1, 1]^T = 2.7 - (0 \cdot 1 + 1 \cdot 1 + 2 \cdot 1) = -0.3 < 1$$

$$z^{(3)}: y^{(3)} - \theta \cdot x^{(3)} = -0.7 - [0, 1, 2]^T \cdot [1, 1, -1]^T = -0.7 - (0 \cdot 1 + 1 \cdot 1 + 2 \cdot (-1)) = 0.3 < 1$$

$$z^{(4)}: y^{(4)} - \theta \cdot x^{(4)} = 2 - [0, 1, 2]^T \cdot [-1, 1, 1]^T = 2 - (0 \cdot (-1) + 1 \cdot 1 + 2 \cdot 1) = -1 < 1$$

$$R_m(\theta) = \frac{1}{4} ((1-0) + (1-(-0.3)) + (1-0.3) + (1-(-1))) = 1.25$$

$R_m(\theta)$ with Squared Error loss:

$$R_m(\theta) = \frac{1}{n} \sum_{t=1}^n \text{loss}_s(y^{(t)} - \theta \cdot x^{(t)}) \quad \text{with } \text{loss}(z) = \frac{z^2}{2}$$

$$R_m(\theta) = \frac{1}{4} \left(\frac{0^2}{2} + \frac{(-0.3)^2}{2} + \frac{(0.3)^2}{2} + \frac{(-1)^2}{2} \right) = 0.1675$$

Let's compute GBA for one example:

$$\nabla_{\theta} (y^{(t)} - \theta \cdot x^{(t)})^2 / 2 = \frac{2}{2} (y^{(t)} - \theta \cdot x^{(t)}) (-x^{(t)}) = -(y^{(t)} - \theta \cdot x^{(t)}) (x^{(t)})$$

Algorithm:

1. Initialize $\theta = 0$
2. Randomly pick $t = \{1, \dots, m\}$
3. update $\theta = \theta - \underbrace{\left[-(y^{(t)} - \theta \cdot x^{(t)}) \cdot (x^{(t)}) \right]}_{\uparrow \text{ minus because we are going against the direction of our gradient to minimize the expression.}} = \theta + \underbrace{\eta}_{\uparrow \text{ learning rate (how much to go in the opposite direction)}} (y^{(t)} - \theta \cdot x^{(t)}) x^{(t)}$

$\eta_k = \frac{1}{1+k}$: with k number of iteration

Closed form solution:

$$R_m(\theta) = \frac{1}{m} \sum_{t=1}^m (y^{(t)} - \theta \cdot x^{(t)})^2 / 2$$

$$\nabla_{\theta} R_m(\theta) \Big|_{\theta=\hat{\theta}} = \frac{1}{m} \sum_{t=1}^m \left[\nabla_{\theta} (y^{(t)} - \theta \cdot x^{(t)})^2 / 2 \right] \Big|_{\theta=\hat{\theta}}$$

$\hat{\theta}$ is where this expression is = 0

$$= -\frac{1}{m} \sum_{t=1}^m (y^{(t)} - \theta \cdot x^{(t)}) x^{(t)}$$

$$= -\frac{1}{m} \underbrace{\sum_{t=1}^m y^{(t)} x^{(t)}}_{\text{vector } b} + \frac{1}{m} \sum_{t=1}^m \hat{\theta} \cdot \overset{\text{dot product}}{x^{(t)} x^{(t)}}$$

$$= -b + \frac{1}{m} \sum_{t=1}^m x^{(t)} \hat{\theta} \cdot x^{(t)}$$

$$\hat{\theta} \cdot x^{(t)} = (x^{(t)})^T \hat{\theta}$$

$$= -b + \frac{1}{m} \underbrace{\sum_{t=1}^m x^{(t)} (x^{(t)})^T}_{A} \hat{\theta}$$

$$= -b + A \hat{\theta} \Rightarrow \text{we are trying to solve } A \hat{\theta} = b \quad (-b + A \hat{\theta} = 0)$$

$$\Rightarrow \hat{\theta} = A^{-1} b \quad A \text{ is not always reversible, } A \text{ is reversible when } m \gg d$$

Cost of reversibility of A has a huge cost $\mathcal{O}(d^3)$

num of training examples \uparrow dimension \uparrow

Regularization:

We completed the discussion of 2 ways of doing the linear algorithm for the linear Regression

What happens if we don't have enough training examples?

What happens if our training examples contain some noise?

We will solve these questions with a mechanism called Regularization.

What Regularization will do?

It will push you away from trying to perfectly fit your training examples.

Reach Regression: How well we are fitting the training examples like the Empirical Risk

$$J_{\lambda, m}(\theta) = \underbrace{\frac{\lambda}{2} \|\theta\|^2}_{\text{keep all the } \theta \text{ very close to zero, because it's norm}} + \underbrace{R_m(\theta)}_{\text{Empirical Risk: Controls how much loss we're incurring in our training data set.}}$$

keep all the θ very close to zero, because it's norm

Empirical Risk: Controls how much loss we're incurring in our training data set.

We want to keep θ 's grounded in some area and only push them when we have enough evidence.

$$J_{\lambda, m}(\theta) = \frac{\lambda}{2} \|\theta\|^2 + \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \theta \cdot x^{(i)})^2 / 2$$

$$\begin{aligned} \nabla_{\theta} J_{\lambda, m}(\theta) &= \nabla_{\theta} \left(\frac{\lambda}{2} \|\theta\|^2 + (y^{(i)} - \theta x^{(i)})^2 / 2 \right) \\ &= \lambda \theta - (y^{(i)} - \theta x^{(i)}) x^{(i)} \end{aligned}$$

1. Initialize: $\theta = 0$

2. Randomly pick $t = \{1, \dots, m\}$

3. $\theta = \theta - \eta (\lambda \theta - (y^{(t)} - \theta x^{(t)}) x^{(t)})$

$$\theta = \theta - \eta (\lambda \theta - (y^{(t)} - \theta x^{(t)}) x^{(t)})$$

$$\theta = (1 - \eta \lambda) \theta + \eta (y^{(t)} - \theta x^{(t)}) x^{(t)}$$



Goals of this expression:

1. It's self-correcting it pushes the parameters in the right direction to minimize the loss.

2. It also pushes our thetas down

Regularization: Extreme case 1

$$J_{m,\lambda}(\theta, \theta_0) = \frac{1}{n} \sum_{t=1}^n \frac{(y^{(t)} - \theta \cdot x^{(t)} - \theta_0)^2}{2} + \frac{\lambda}{2} \|\theta\|^2$$

* If we increase λ to ∞ , minimizing J is equivalent to minimizing $\|\theta\|$.

Thus θ will have to be a zero vector $\Rightarrow f(x) = \theta \cdot x + \theta_0$ becomes $f(x) = \theta_0$, a horizontal line

* If we decrease λ to zero, minimizing J is equivalent to minimizing $\frac{1}{n} \sum_{t=1}^n \frac{(y^{(t)} - \theta \cdot x^{(t)} - \theta_0)^2}{2}$

which is the "fit" \Rightarrow fitting the data

$\lambda \nearrow \Leftrightarrow$ more mistakes, we are pushing $\theta \rightarrow 0$

$\lambda \searrow \Leftrightarrow$ Fitting training data