

# Explanatory Notes for 6.390

Shauntclair Ruiz (Current TA)

Fall 2022

## Application to Regression

One nice thing about **gradient descent** is that it is **easy** to switch the kind of problem you're applying it to: all you need is your **parameters**(s)  $\theta$ , and a function to optimize,  $J$ .

From there, you can just **compute** the gradient.

### Ordinary Least Squares

Our **loss** function is

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \left( (\theta^T x^{(i)} + \theta_0) - y^{(i)} \right)^2 \quad (1)$$

Or, in **matrix** terms,

Including the appended row of 1's from before.

$$J = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^T (\tilde{X}\theta - \tilde{Y})$$

Our gradient, according to **matrix derivative** rules, is

$$\nabla_{\theta} J(\theta) = \frac{2}{n} \tilde{X}^T (\tilde{X}\theta - \tilde{Y}) \quad (2)$$

Before, we set it equal to **zero**. But here, we can instead take **steps** towards the solution, using **gradient descent**.

We could use the **matrix** form, but sometimes it's easier to use a **sum**. Fortunately, derivatives are easy with a sum. If so, here's **another** way to write it:

$$\nabla_{\theta} J(\theta) = \frac{2}{n} \sum_{i=1}^n \left( \theta^T x^{(i)} - y^{(i)} \right) x^{(i)} \quad (3)$$

Either way, we use gradient descent **normally**:

Remember that  $\theta_{old}$  is an **input** to the gradient, not multiplied by it!

$$\theta_{new} = \theta_{old} - \eta \nabla_{\theta} J(\theta_{old})$$

Using  $\theta^{(t)}$  notation:

$$\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_{\theta} J(\theta^{(t-1)})$$

### Ridge Regression

Ridge regression is similar.

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \left( \underbrace{(\theta^T \mathbf{x}^{(i)} + \theta_0)}_{\text{guess}} - \underbrace{\mathbf{y}^{(i)}}_{\text{answer}} \right)^2 + \underbrace{\lambda \|\theta\|^2}_{\text{Regularizer}}$$

However, we have to treat  $\theta_0$  as **separate** from our other data points, because of **regularization**: remember that it **doesn't** apply to  $\theta_0$ .

For  $\theta$ :

$$\nabla_{\theta} J_{\text{ridge}}(\theta, \theta_0) = \frac{2}{n} \sum_{i=1}^n \left( (\theta^T \mathbf{x}^{(i)} + \theta_0) - \mathbf{y}^{(i)} \right) \mathbf{x}^{(i)} + 2\lambda \theta \quad (4)$$

For  $\theta_0$ :

$$\frac{\partial J_{\text{ridge}}(\theta, \theta_0)}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n \left( (\theta^T \mathbf{x}^{(i)} + \theta_0) - \mathbf{y}^{(i)} \right) \quad (5)$$

Notice that we used a **gradient** for our vector  $\theta$ , but since  $\theta_0$  is a single variable, we just used a **simple derivative**!

### Concept 1

The **gradient**  $\frac{dJ}{d\theta}$  must have the **same shape as  $\theta$** : this shape-matching is why we can easily **subtract** it during gradient descent.

$$\underbrace{\theta_{\text{new}}}_{(d \times 1)} = \underbrace{\theta_{\text{old}}}_{(d \times 1)} - \eta \underbrace{\nabla_{\theta} J(\theta_{\text{old}})}_{(d \times 1)}$$