

Discussion #9

Name:

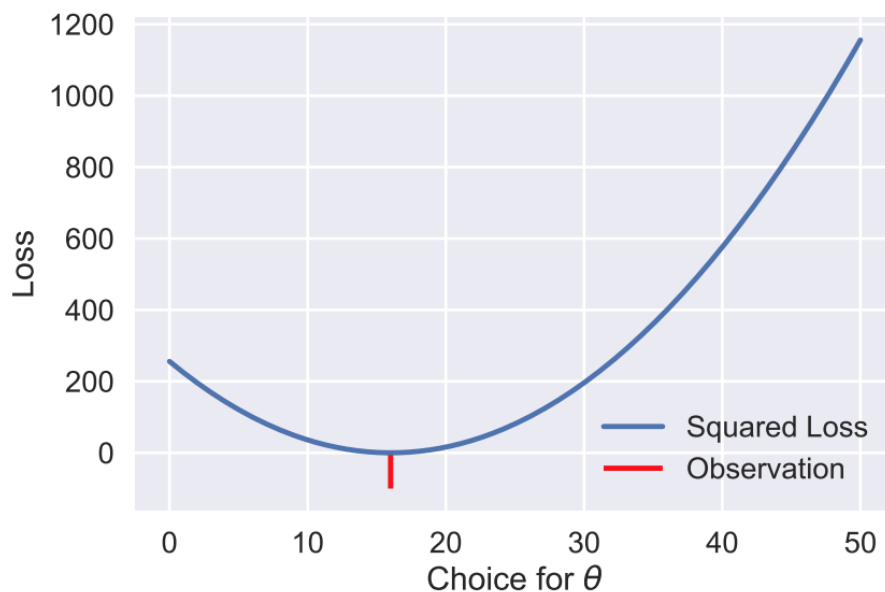
Loss Functions

1. Recall the loss functions discussed during lecture. θ represents our estimate of our parameter, and y represents a data point. Discuss the advantages and drawbacks of each of the following loss functions:

(a) Squared loss: $L(\theta, y) = (y - \theta)^2$

Solution:

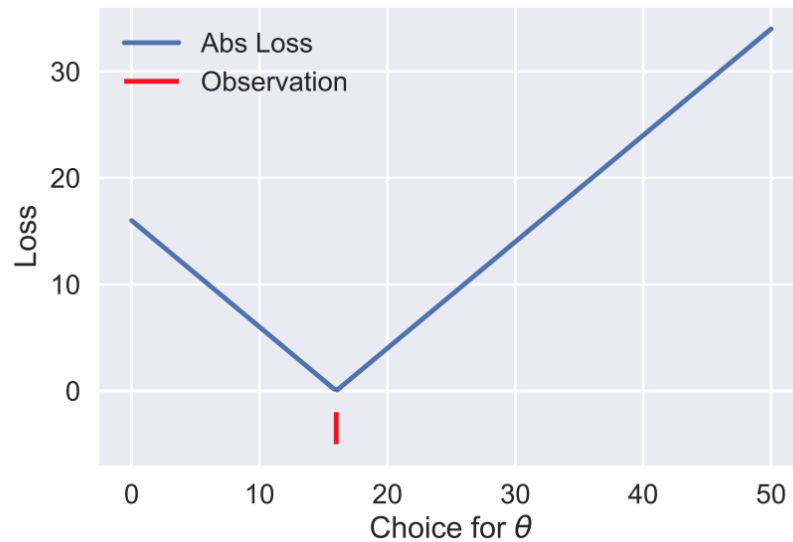
- Commonly used loss function.
- Sensitive to outliers.
- Differentiable, hence we can find solution analytically.



(b) Absolute Loss: $L(\theta, y) = |y - \theta|$

Solution:

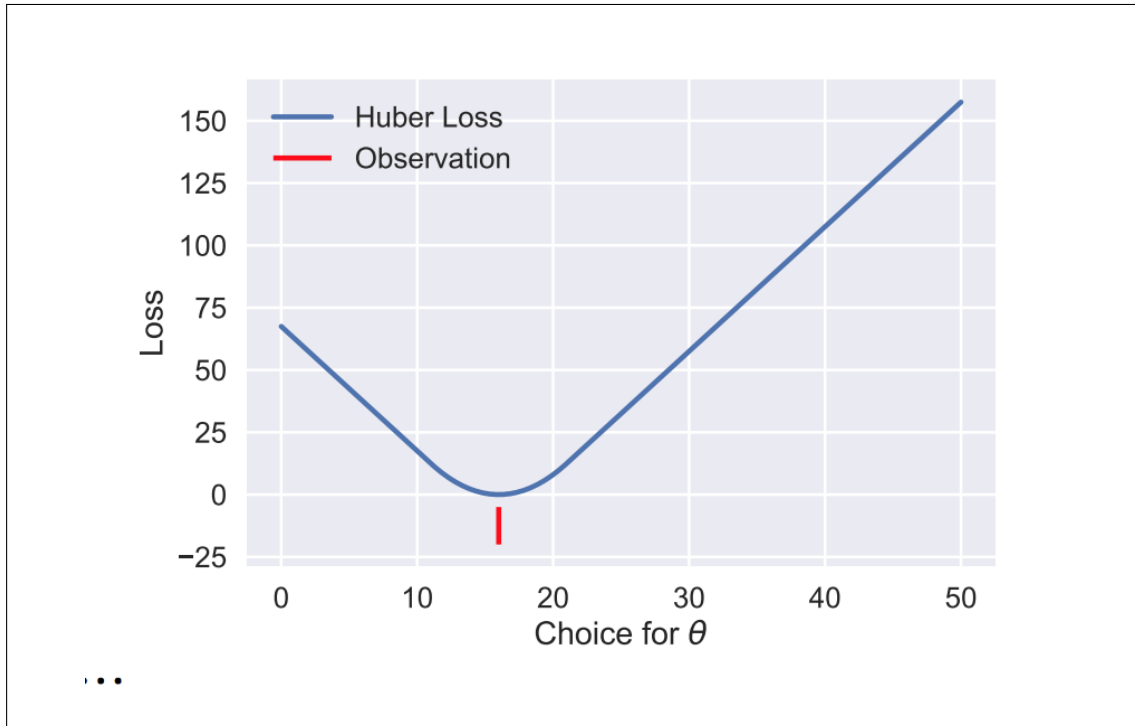
- Less sensitive to outliers.
- Not differentiable; hence, no analytic closed form solution.



(c) Huber Loss:
$$L_{\alpha}(\theta, y) = \begin{cases} (y - \theta)^2 & |y - \theta| < \alpha \\ \alpha(|y - \theta| - \alpha/2) & \text{otherwise} \end{cases}$$

Solution:

- Melding of both square and absolute loss. More sensitive to small deviations from θ than squared loss but less sensitive to outlier values.
- Differentiable; however, we cannot find an analytic solution because we cannot isolate the θ term
- Must choose where the hinges occur (another tuning parameter!)



Loss Minimization

2. Consider the following loss function:

$$L(\theta, x) = \begin{cases} 4(\theta - x) & \theta \geq x \\ x - \theta & \theta < x \end{cases}$$

Given a sample of x_1, \dots, x_n , find the optimal θ that minimizes the the average loss.

Solution:

$$\frac{\partial L(\theta, x)}{\partial \theta} = \begin{cases} 4 & \theta \geq x \\ -1 & \theta < x \end{cases}$$

$$\sum_{i=1}^n \frac{\partial L(\theta, x_i)}{\partial \theta} = \sum_{\theta < x_i} -1 + \sum_{\theta \geq x_i} 4 = 0$$

Which implies that

$$\#\{x_i > \theta\} = 4\#\{x_i \leq \theta\}$$

We also know that

$$\#\{x_i > \theta\} + \#\{x_i \leq \theta\} = n$$

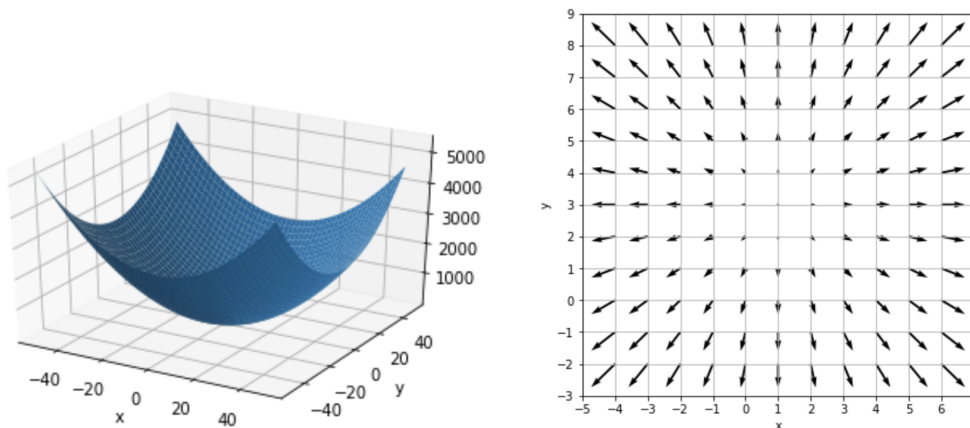
Substituting the first equation into the second, we get

$$n - \#\{x_i \leq \theta\} = 4\#\{x_i \leq \theta\} \implies 0.20n = \#\{x_i \leq \theta\}$$

Thus, θ is the 20th percentile of x_1, \dots, x_n .

Gradients

3. On the left is a 3D plot of $f(x, y) = (x - 1)^2 + (y - 3)^2$. On the right is a plot of its gradient field. Note that the arrows show the relative magnitudes of the gradient vector.



- (a) From the visualization, what do you think is the minimal value of this function and where does it occur?

Solution: Since $(x - 1)^2$ and $(y - 3)^2$ are both nonnegative, the minimum function value of $f(x, y)$ is attained when both are equal to zero. This occurs at $(1, 3)$ where the gradient field shows the smallest (in magnitude) vectors.

- (b) Calculate the gradient $\nabla f = \left[\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \right]^T$.

Solution:

$$\left[\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \right]^T = [2(x - 1) \quad 2(y - 3)]^T.$$

(c) When $\nabla f = \mathbf{0}$, what are the values of x and y ?

Solution:

$$\nabla f = \mathbf{0} \implies 2(x - 1) = 2(y - 3) = 0 \implies x = 1, y = 3.$$

If the gradient is equal to zero, then the function must be at a local minima. The only minima in this case is the global minima, meaning it must be at $(1, 3)$, due to part (e).

4. In this question, we will explore some basic properties of the gradient.

Note: In this class, we use the following conventions:

- x represents a scalar
- X represents a random variable
- \mathbf{x} represents a vector
- \mathbf{X} represents a matrix or a random vector (context will tell)

(a) Determine the derivative of $f(x) = a_0 + a_1x$ and gradient of $g(x_1, x_2) = a_0 + a_1x_1 + a_2x_2$.

Solution:

$$\frac{df}{dx} = a_1$$
$$\nabla g = \left[\frac{\partial g}{\partial x_1} \quad \frac{\partial g}{\partial x_2} \right]^T = [a_1 \quad a_2]^T$$

(b) Suppose $\mathbf{x} = [x_1 \quad x_2 \quad \dots \quad x_n]^T$, and $h(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$, where $\mathbf{a}, \mathbf{x} \in \mathbb{R}^n$. Determine ∇h .

Solution: Note that $h(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$ is just a concise way of writing

$$h(\mathbf{x}) = \sum_{i=1}^n a_i x_i = a_1 x_1 + a_2 x_2 + \dots + a_n x_n$$

So as in (a), we have

$$\nabla h = \left[\frac{\partial h}{\partial x_1} \quad \frac{\partial h}{\partial x_2} \quad \dots \quad \frac{\partial h}{\partial x_n} \right]^T = [a_1 \quad a_2 \quad \dots \quad a_n]^T = \mathbf{a}$$

(c) Determine the gradient of $f(\mathbf{x}) = \mathbf{x}^T \mathbf{x}$. (Hint: f is a scalar-valued function. How can you write $\mathbf{x}^T \mathbf{x}$ as a sum of scalars?)

Solution: $f(\mathbf{x})$ can also be expanded as $\sum_{i=1}^n x_i^2 = x_1^2 + x_2^2 + \dots + x_n^2$

$$\nabla f = [2x_1 \quad 2x_2 \quad \dots \quad 2x_n]^T = 2\mathbf{x}$$

Gradient Descent Algorithm

5. Given the following loss function and $\mathbf{x} = (x_i)_{i=1}^n$, $\mathbf{y} = (y_i)_{i=1}^n$, θ^t , explicitly write out the update equation for θ^{t+1} in terms of x_i , y_i , θ^t , and α , where α is the step size.

$$L(\theta, \mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (\theta^2 x_i^2 - \log(y_i))$$

Solution:

$$\theta^{t+1} \leftarrow \theta^t - \alpha \left. \frac{\partial L}{\partial \theta} \right|_{\theta=\theta^t}$$

$$\frac{\partial L}{\partial \theta} = \frac{1}{n} \sum_{i=1}^n 2\theta x_i^2$$

6. (a) In your own words, describe how to use the update equation in the gradient descent algorithm.
- (b) Say that x and y are your model parameters and f as defined in question 1 is your loss function. Describe in your own words what happens “visually” as the gradient descent algorithm runs.

Convexity

Convexity allows optimization problems to be solved more efficiently and for global optimums to be realized. Mainly, it gives us a nice way to minimize loss (i.e. gradient descent). There are three ways to informally define convexity.

- Walking in a straight line between points on the function keeps you above the function. This works for any function.
- The tangent line at any point lies below the function (globally). The function must be differentiable.

- c. The second derivative is non-negative everywhere (aka "concave up" everywhere). The function must be twice differentiable.
7. Find a counterexample for the claim that the composition of two convex functions is also convex. $h = g(f(x))$

Solution: Let $f(x) = x^2, g(x) = -x$. $g(f(x)) = -x^2$ which is not convex.