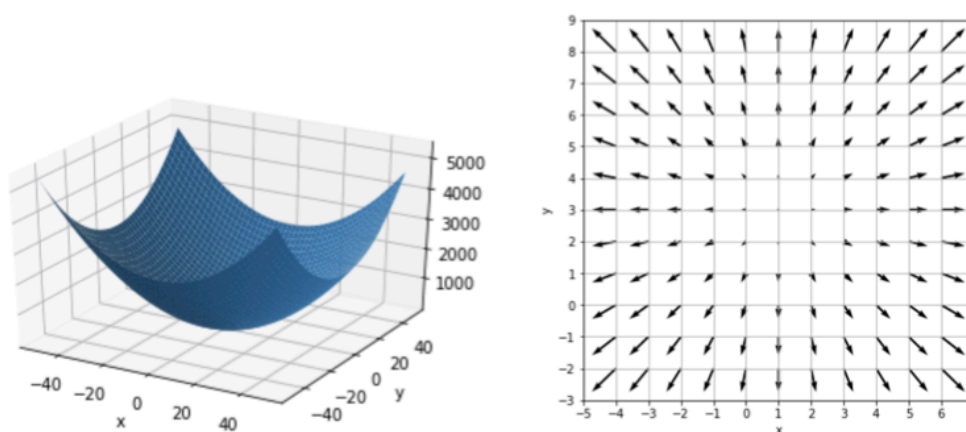# Assignment 4

# ECE 4710J

# Due 11/29/2021

## Visualizing Gradients

1. 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?

   Minimum value: $0$, occurs at $(1, 3)$

   (b) Calculate the gradient $\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \end{bmatrix}^T$.

   $$\nabla f = \begin{bmatrix} 2(x-1) \\ 2(y-3) \end{bmatrix}$$

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

   $\nabla f = \vec{0} \Rightarrow \begin{cases} 2(x-1) = 0 \\ 2(y-3) = 0 \end{cases} \Rightarrow \begin{cases} x = 1 \\ y = 3 \end{cases}$

# Gradient Descent Algorithm

2. Given the following loss function and $\vec{x} = [x_i]_{i=1}^n$, $\vec{y} = [y_i]_{i=1}^n$, and $\theta^t$, explicitly write out the update equation for $\theta^{t+1}$ in terms of $x_i$, $y_i$, $\theta^t$, and $\alpha$, where $\alpha$ is the constant learning rate.

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

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - \alpha \nabla_{\vec{\theta}} L(\vec{\theta}, \vec{x}, \vec{y})$$

$$= \vec{\theta}^{(t)} - \alpha \left( \frac{1}{n} \sum_{i=1}^n 2\theta x_i^2 \right)$$

$$= \vec{\theta}^{(t)} \left( 1 - \frac{2\alpha}{n} \sum_{i=1}^n x_i^2 \right)$$
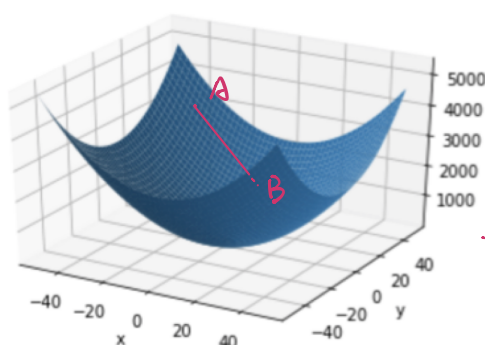
# Convexity

3. 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.

   a. Walking in a straight line between points on the function keeps you at or above the function. This works for any function.

   b. The tangent line at any point lies at or below the function, globally. To use this definition, the function must be differentiable.

   c. The second derivative is non-negative everywhere (in other words, the function is "concave up" everywhere). To use this definition, the function must be twice differentiable.

   Is the function described in Question 1 convex? Make an argument visually.

Yes, the function in Question 1 is convex.

According to definition c. $f(x,y) = (x-1)^2 + (y-3)^2$, $\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial^2 x} \\ \frac{\partial^2 f}{\partial^2 y} \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$.

which is non-negative.



The line A.B for $\forall$ A.B on the function keeps above the function.

# GPA Descent

4. Consider the following non-linear model with two parameters:

$$f_\theta(x) = \theta_0 \cdot 0.5 + \theta_0 \cdot \theta_1 \cdot x_1 + \sin(\theta_1) \cdot x_2$$

For some nonsensical reason, we decide to use the residuals of our model as the loss function. That is, the loss for a single observation is

$$L(\theta) = y_i - f_\theta(x_i)$$

We want to use gradient descent to determine the optimal model parameters, $\hat{\theta}_0$ and $\hat{\theta}_1$.

(a) Suppose we have just one observation in our training data, $(x_1 = 1, x_2 = 2, y = 4)$. Assume that we set the learning rate $\alpha$ to 1. An incomplete version of the gradient descent update equation for $\theta$ is shown below. $\theta_0^{(t)}$ and $\theta_1^{(t)}$ denote the guesses for $\theta_0$ and $\theta_1$ at timestep $t$, respectively.

$$\begin{bmatrix} \theta_0^{(t+1)} \\ \theta_1^{(t+1)} \end{bmatrix} = \begin{bmatrix} \theta_0^{(t)} \\ \theta_1^{(t)} \end{bmatrix} - \begin{bmatrix} A \\ B \end{bmatrix}$$

Express both $A$ and $B$ in terms of $\theta_0^{(t)}$, $\theta_1^{(t)}$, and any necessary constants.

$$\begin{bmatrix} \theta_0^{(t+1)} \\ \theta_1^{(t+1)} \end{bmatrix} = \begin{bmatrix} \theta_0^{(t)} \\ \theta_1^{(t)} \end{bmatrix} - \alpha \nabla_{\vec{\theta}} L(\theta)$$

$\because x_1 = 1, \ x_2 = 2$

$$\Rightarrow \begin{bmatrix} A \\ B \end{bmatrix} = \alpha \nabla_{\vec{\theta}} L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial \theta_0} \\ \frac{\partial L}{\partial \theta_1} \end{bmatrix} = \begin{bmatrix} \frac{\partial(-f_\theta(x_i))}{\partial \theta_0} \\ \frac{\partial(-f_\theta(x_i))}{\partial \theta_1} \end{bmatrix} = \begin{bmatrix} -0.5 - \theta_1^{(t)} x_1 \\ -\theta_0^{(t)} x_1 - \cos(\theta_1^{(t)}) x_2 \end{bmatrix} = \begin{bmatrix} -0.5 - \theta_1^{(t)} \\ -\theta_0^{(t)} - 2\cos(\theta_1^{(t)}) \end{bmatrix}$$

$$\Rightarrow A = -0.5 - \theta_1^{(t)}, \quad B = -\theta_0^{(t)} - 2\cos(\theta_1^{(t)})$$

(b) Assume we initialize both $\theta_0^{(0)}$ and $\theta_1^{(0)}$ to 0. Determine $\theta_0^{(1)}$ and $\theta_1^{(1)}$ (i.e. the guesses for $\theta_0$ and $\theta_1$ after one iteration of gradient descent).

$$\begin{bmatrix} \theta_0^{(1)} \\ \theta_1^{(1)} \end{bmatrix} = \begin{bmatrix} \theta_0^{(0)} \\ \theta_1^{(0)} \end{bmatrix} - \begin{bmatrix} A \\ B \end{bmatrix}$$

$A = -0.5 - \theta_1^{(0)} = -0.5$    $B = -\theta_0^{(0)} - 2\cos(\theta_1^{(0)}) = -2$

$$= \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} -0.5 \\ -2 \end{bmatrix}$$

$$= \begin{bmatrix} 0.5 \\ 2 \end{bmatrix} \Rightarrow \theta_0^{(1)} = 0.5 \quad \theta_1^{(1)} = 2$$

(c) What happens to $\theta_0^{(t)}$ as $t \to \infty$ (i.e. as we run more and more iterations of gradient descent)?

$$\begin{bmatrix} \theta_0^{(t+1)} \\ \theta_1^{(t+1)} \end{bmatrix} = \begin{bmatrix} \theta_0^{(t)} \\ \theta_1^{(t)} \end{bmatrix} - \begin{bmatrix} -0.5 - \theta_1^{(t)} \\ -\theta_0^{(t)} - 2\cos(\theta_1^{(t)}) \end{bmatrix} = \begin{bmatrix} \theta_0^{(t)} + \theta_1^{(t)} + 0.5 \\ \theta_1^{(t)} + \theta_0^{(t)} + 2\cos(\theta_1^{(t)}) \end{bmatrix}$$

$$\begin{bmatrix} \theta_0^{(1)} \\ \theta_1^{(1)} \end{bmatrix} = \begin{bmatrix} 0.5 \\ 2 \end{bmatrix} \Rightarrow \begin{bmatrix} \theta_0^{(2)} \\ \theta_1^{(2)} \end{bmatrix} = \begin{bmatrix} 3 \\ 1.67 \end{bmatrix},$$ since $2\cos(\theta_1^{(t)}) \in [-2,2]$, $\theta_0^{(t)} - 2 > 0 \Rightarrow \theta_0^{(t)}$ and $\theta_1^{(t)}$ will

diverge.

meaning that $\theta_0^{(t)}$ will goes to infinity