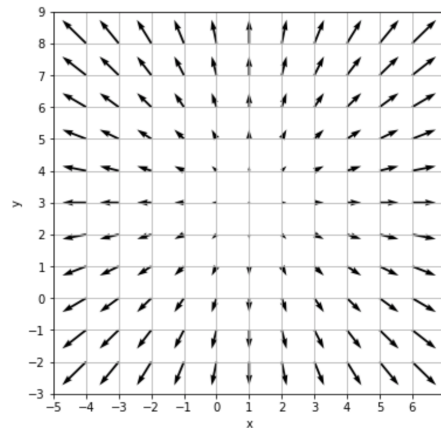
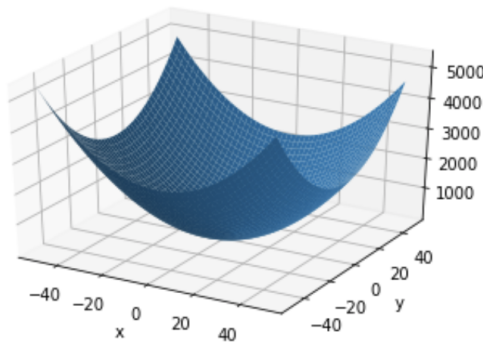


Discussion #7 Solutions

Name:

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?

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).

Gradient Descent Algorithm

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

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

Solution:

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

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

3. (a) The learning rate α can *potentially* affect which of the following? Select all that apply. Assume nothing about the function being minimized other than that its gradient exists. You may assume the learning rate is positive.

- ☐ A. The speed at which we converge to a minimum.
- ☐ B. Whether gradient descent converges.
- ☐ C. The direction in which the step is taken.
- ☐ D. Whether gradient descent converges to a local minimum or a global minimum.

- (b) Suppose we run gradient descent with a fixed learning rate of $\alpha = 0.1$ to minimize the 2D function $f(x, y) = 5 + x^2 + y^2 + 5xy$.

The gradient of this function is

$$\nabla_{x,y} f(x, y) = \begin{bmatrix} 2x + 5y \\ 2y + 5x \end{bmatrix}$$

If our starting guess is $x^{(0)} = 1$, $y^{(0)} = 2$, what will be our next guess $x^{(1)}$, $y^{(1)}$?

$x^{(1)} =$ $y^{(1)} =$

Solution: The gradient is $= [2*1 + 5*2, 2*2 + 5*1] = [12, 9]$ so next guess is $[1, 2] - 0.1 * [12, 9] = -0.2, 1.1$

- (c) Suppose we are performing gradient descent to minimize the empirical risk of a linear regression model $y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2$ on a dataset with 100 observations. Let \mathcal{D} be the number of components in the gradient, e.g. $\mathcal{D} = 2$ for the equation in part b. What is \mathcal{D} for the gradient used to optimize this linear regression model?

☐ A. 2 ☐ B. 3 ☒ C. 4 ☐ D. 8 ☐ E. 100 ☐ F. 200 ☐ G. 300
☐ H. 400 ☐ I. 800