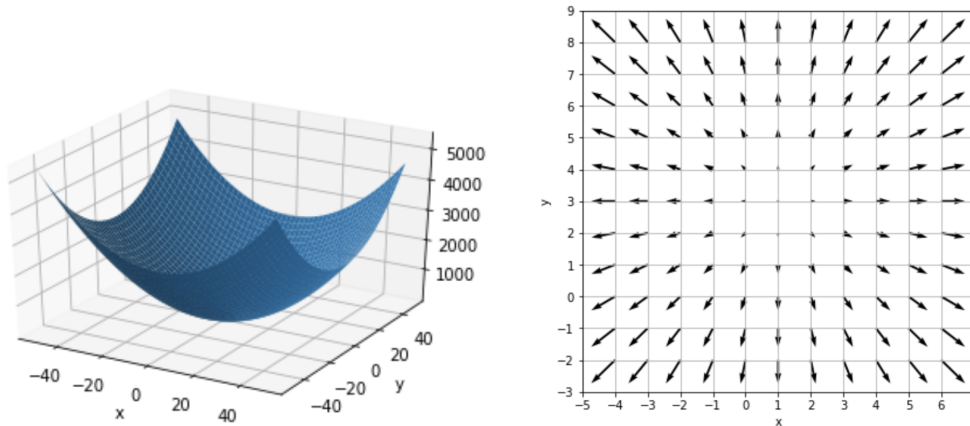## Discussion #11 Solutions

*Name:*

# 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 $\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \end{bmatrix}^T$.

> **Solution:**
> $$\begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \end{bmatrix}^T = \begin{bmatrix} 2(x - 1) & 2(y - 3) \end{bmatrix}^T.$$

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

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

1

> 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 $\mathbf{x} = (x_i)_{i=1}^n$, $\mathbf{y} = (y_i)_{i=1}^n$, $\beta^t$, explicitly write out the update equation for $\beta^{t+1}$ in terms of $x_i$, $y_i$, $\beta^t$, and $\alpha$, where $\alpha$ is the step size.

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

**Solution:**

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

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

# 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 above the function. This works for any function.

   b. 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.

   (a) Is the function described in question 1 convex? Make an argument visually.

> **Solution:** Yes, walking in a straight line between any two points on the graph will keep us above the graph.
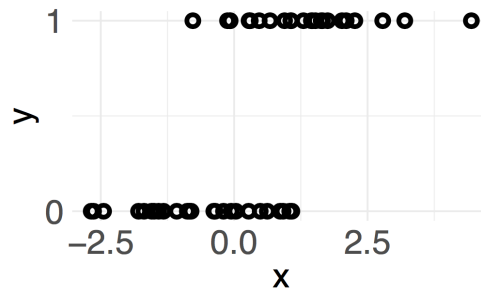
(b) Find a counterexample for the claim that the composition of two convex functions is also convex. $h = g(f(x))$

> **Solution:** Let $f(x) = x^2, g(x) = -x$. $g(f(x)) = -x^2$ which is not convex.
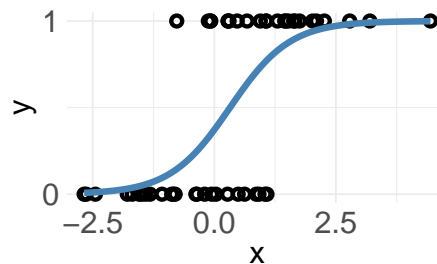
# Logistic Regression

The next two questions refer to a binary classification problem with a single feature $x$.

4. Based on the scatter plot of the data below, draw a reasonable approximation of the logistic regression probability estimates for $\mathbb{P}\left(Y = 1 \mid x\right)$.



**Solution:**



5. You have a classification data set consisting of two $(x, y)$ pairs $(1, 0)$ and $(-1, 1)$.

   The covariate vector $\mathbf{x}$ for each pair is a two-element column vector $\begin{bmatrix} 1 & x \end{bmatrix}^T$.

   You run an algorithm to fit a model for the probability of $Y = 1$ given $\mathbf{x}$:

   $$\mathbb{P}\left(Y = 1 \mid \mathbf{x}\right) = \sigma(\mathbf{x}^T \beta)$$

   where

   $$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

   Your algorithm returns $\hat{\beta} = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}^T$

(a) Calculate $\hat{\mathbb{P}}\left(Y = 1 \mid \mathbf{x} = [1 \quad 0]^T\right)$

**Solution:**

$$\hat{\mathbb{P}}\left(Y = 1 \mid \mathbf{X} = [1 \quad 0]^T\right) = \sigma\left([1 \quad 0]\begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{2} \end{bmatrix}\right)$$

$$= \sigma\left(1 \times -\frac{1}{2} + 0 \times -\frac{1}{2}\right)$$

$$= \sigma\left(-\frac{1}{2}\right)$$

$$= \frac{1}{1 + \exp(\frac{1}{2})}$$

$$\approx 0.38$$

(b) The empirical risk using log loss (a.k.a., cross-entropy loss) is given by:

$$R(\beta) = \frac{1}{n}\sum_{i=1}^{n} -\log \hat{\mathbb{P}}\left(Y = y_i \mid \mathbf{x_i}\right)$$

$$= -\frac{1}{n}\sum_{i=1}^{n} y_i \log \hat{\mathbb{P}}\left(Y = 1 \mid \mathbf{x_i}\right) + (1 - y_i)\log \hat{\mathbb{P}}\left(Y = 0 \mid \mathbf{x_i}\right)$$

And $\hat{\mathbb{P}}\left(Y = 1 \mid \mathbf{x_i}\right) = \frac{\exp(\mathbf{x_i}^T\beta)}{1+\exp(\mathbf{x_i}^T\beta)}$ while $\hat{\mathbb{P}}\left(Y = 0 \mid \mathbf{x_i}\right) = \frac{1}{1+\exp(\mathbf{x_i}^T\beta)}$. Therefore,

$$R(\beta) = -\frac{1}{n}\sum_{i=1}^{n} y_i \log \frac{\exp(\mathbf{x_i}^T\beta)}{1 + \exp(\mathbf{x_i}^T\beta)} + (1 - y_i)\log\frac{1}{1 + \exp(\mathbf{x_i}^T\beta)}$$

$$= -\frac{1}{n}\sum_{i=1}^{n} y_i\mathbf{x}_i^T\beta + \log(\sigma(-\mathbf{x}_i^T\beta))$$

Let $\beta = [\beta_0 \quad \beta_1]$. Explicitly write out the empirical risk for the data set $(1, 0)$ and $(-1, 1)$ as a function of $\beta_0$ and $\beta_1$.

**Solution:**

$$x_i^T\beta = [1 \quad x_i]\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \beta_0 + \beta_1 x_i$$

For the data point $(1, 0)$, $\mathbf{x}_i = [1 \quad 1]^T$ and $y_i = 0$, so:

$$y_i\mathbf{x}_i^T\beta = 0$$

$$-\mathbf{x}_i^T \beta = -(\beta_0 + \beta_1 \times 1) = -\beta_0 - \beta_1$$

For the data point $(-1, 1)$:

$$y_i x_i^T \beta = 1 \times (\beta_0 + \beta_1 \times -1) = \beta_0 - \beta_1$$

$$-x_i^T \beta = -(\beta_0 + \beta_1 \times -1) = -\beta_0 + \beta_1$$

We can then write the empirical risk as:

$$
\begin{aligned}
R(\beta) &= -\frac{1}{2} \left[ (0 + \log \sigma(-\beta_0 - \beta_1)) + (\beta_0 - \beta_1 + \log \sigma(-\beta_0 + \beta_1)) \right] \\
&= -\frac{1}{2} \left[ \beta_0 - \beta_1 + \log \sigma(-\beta_0 - \beta_1) + \log \sigma(-\beta_0 + \beta_1) \right] \\
&= -\frac{1}{2} \left[ \beta_0 - \beta_1 + \log \left( \frac{1}{1 + \exp(\beta_0 + \beta_1)} \right) + \log \left( \frac{1}{1 + \exp(\beta_0 - \beta_1)} \right) \right] \\
&= \frac{1}{2} \left[ \beta_1 - \beta_0 + \log \left( 1 + \exp(\beta_0 + \beta_1) \right) + \log \left( 1 + \exp(\beta_0 - \beta_1) \right) \right]
\end{aligned}
$$

(c) Calculate the empirical risk for $\hat{\beta} = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}^T$ and the two observations $(1, 0)$ and $(-1, 1)$.

**Solution:**

$$
\begin{aligned}
R(\hat{\beta}) &= \frac{1}{2} \left[ \beta_1 - \beta_0 + \log \left( 1 + \exp(\beta_0 + \beta_1) \right) + \log \left( 1 + \exp(\beta_0 - \beta_1) \right) \right] \\
&= \frac{1}{2} \left[ -\frac{1}{2} - \left( -\frac{1}{2} \right) + \log \left( 1 + \exp(-\frac{1}{2} + -\frac{1}{2}) \right) + \log \left( 1 + \exp(-\frac{1}{2} - -\frac{1}{2}) \right) \right] \\
&= \frac{1}{2} \left[ 0 + \log \left( 1 + \exp(-1) \right) + \log \left( 1 + \exp(0) \right) \right] \\
&= \frac{1}{2} \log(2 + 2e^{-1})
\end{aligned}
$$