

Review of Multivariate Calculus and Optimization by Gradient Descent

CMPUT 328

Nilanjan Ray

Computing Science, University of Alberta, Canada

Review of partial derivatives

- Need to review
 - Functions of several variables and vector valued functions
 - Partial derivatives and gradient
 - Hessian and Jacobian
- https://en.wikipedia.org/wiki/Partial_derivative
- Also look at basics of multi-variate calculus
<https://www.khanacademy.org/math/multivariable-calculus>

Gradient of a function

Example:

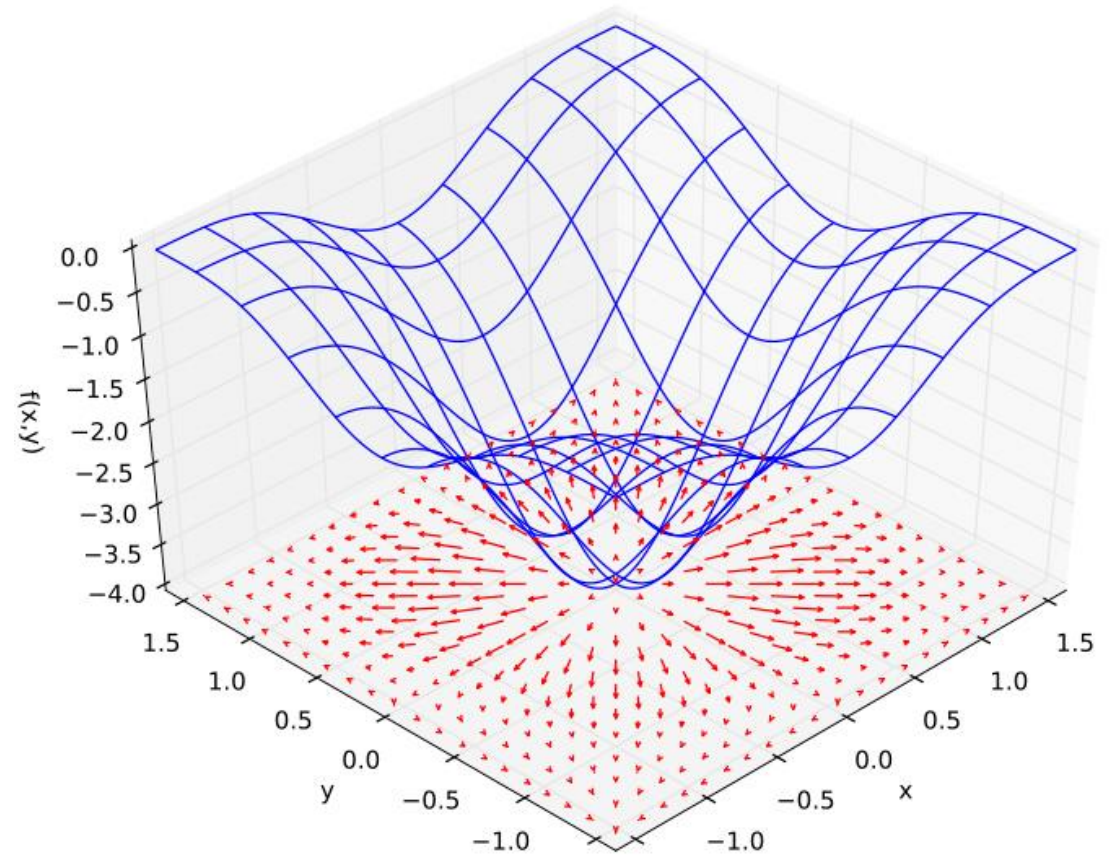
$$f(x, y) = -(\cos^2 x + \cos^2 y)^2$$

$$\nabla f(x, y) = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} 4(\cos^2(x) + \cos^2(y)) \cos(x) \sin(x) \\ 4(\cos^2(x) + \cos^2(y)) \cos(y) \sin(y) \end{bmatrix}$$

Note 1: f is a function of **two variables**,
so gradient of f is a **two dimensional vector**

Note 2: Gradient (vector) of f points toward the
steepest ascent for f

Note 3: At a (local) minimum of f its gradient
becomes a **zero vector**



Example source: Wikipedia

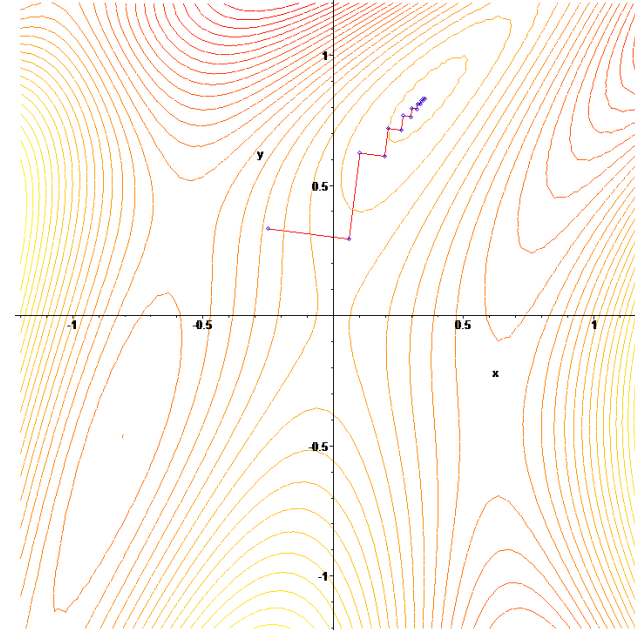
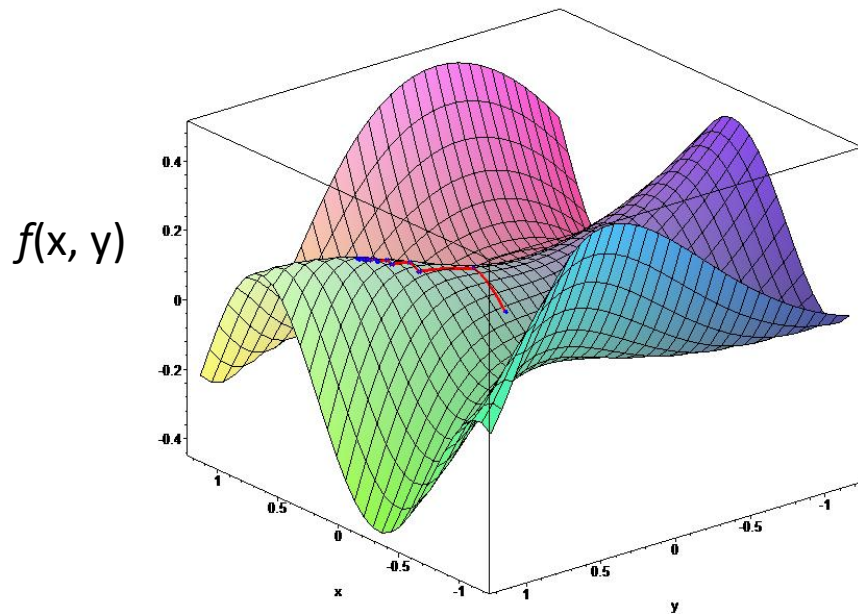
Gradient descent optimization

Start at an initial guess for the optimization variable: \mathbf{x}_0

Iterate until gradient magnitude becomes too small: $\mathbf{x}^{t+1} = \mathbf{x}^t - \alpha \nabla f(\mathbf{x}^t)$

} Gradient descent algorithm

α is called the step-length.



Gradient descent creates a zig-zag path leading to a local minimum of f

Example partial derivative computations

- Let's consider the following function:

$$f(x_1, x_2, x_3, x_4) = (x_1 + 10x_2)^2 + 5(x_3 - x_4)^2 + (x_2 + 2x_3)^4 + 10(x_1 - x_4)^4$$

- Let's compute derivative of this function at

$$[x_1, x_2, x_3, x_4] = [3, -1, 0, 1]$$

- Cross-verify PyTorch partial derivative computations with math formulas
- Gradient descent optimization

Chain rule of derivatives

- Let consider the same function as before:

$$f(x_1, x_2, x_3, x_4) = (x_1 + 10x_2)^2 + 5(x_3 - x_4)^2 + (x_2 + 2x_3)^4 + 10(x_1 - x_4)^4$$

- But, this time x is a vector-valued function of variable z :

$$\begin{aligned}x_1 &= z_1 - z_2, \\x_2 &= z_1^2, \\x_3 &= z_2^2, \\x_4 &= z_1^2 + z_1 z_2\end{aligned}$$

- Let's compute gradient of f with respect to z using chain rule:
Jacobian vector product

$$\frac{\partial f}{\partial z_i} = \sum_j \frac{\partial x_j}{\partial z_i} \frac{\partial f}{\partial x_j}$$

A (toy) learning/optimization example

- You have a target probability vector: $y = [0. 0. 1. 0. 0. 0. 0. 0. 0. 0.]$
- Your model parameters are a 10-element vector z
- Your model: $p = \text{softmax}(z)$, where softmax function is defined as:

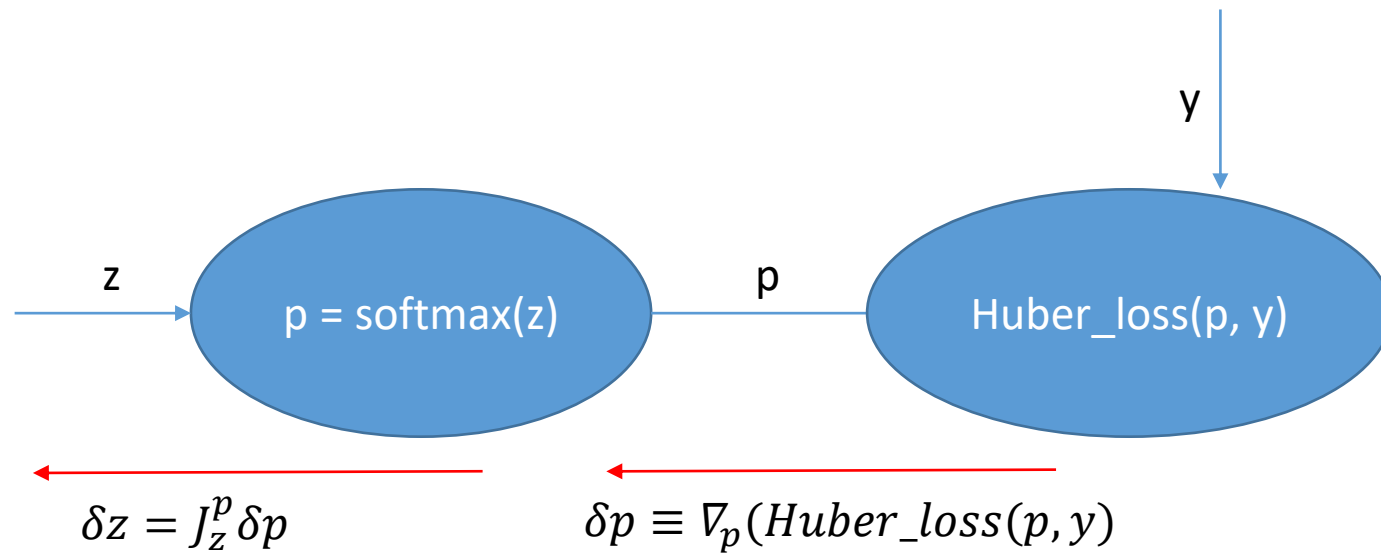
$$p_i = \frac{\exp(z_i)}{\sum_{j=1}^{10} \exp(z_j)}$$

Learning the model: set parameter vector z so that p is as close as possible to the target vector y

Loss function

- We need to measure the distance between output probability p and target y
- Let's work with an interesting loss function called Huber loss:
https://en.wikipedia.org/wiki/Huber_loss
- We need to find the distance, i.e., Huber loss between p and y
- Then, we will compute gradient of this loss with respect to model parameter vector z
- Next, we will adjust current value of z by taking a step in the opposite direction of the gradient
- We will repeat these steps to close the gap between p and y

Computational graph for the toy example



Your assignment #1 will ask you to compute Jacobian and chain rule. Here we will use autograd!

Gradient descent Algorithm

Repeat the following steps

1. (Forward pass) Compute p and Huber_loss
2. (Backward pass) Compute δp and δz
3. (Adjust parameter) $z = z - (\text{step_size})\delta z$

Jacobian-vector product
(chain rule)