


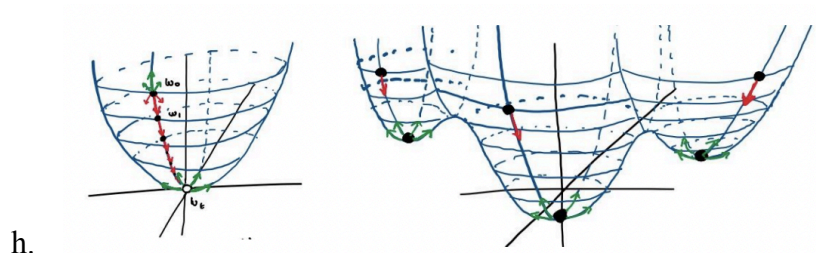
## Gradient Descent

### 1. Gradient Descent (Intuition)

- Optimization method when there is no closed form solution to finding the extrema of a function
- What to do if we don't have an optimization method?
- Example: Logistic Regression
- Goal: find a sequence of  $w_i$ 's and  $b$ 's that converge toward a minimum
- Consider a random weight  $w_0$ . What happens to  $\text{Loss}(w_0)$  as you nudge  $w_0$  slightly?

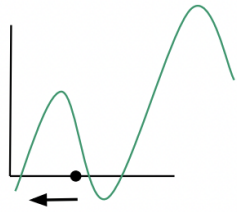


- As such we can define the following sequence:
  - $w_1 = \text{best nudge to } w_0$
  - $w_2 = \text{best nudge to } w_1$
  - ...
  - Until we reach  $w_t$  that looks like 
  - At this point we can stop updating  $w$



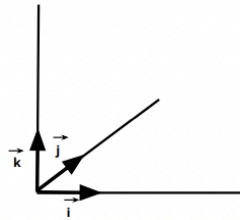
### 2. Gradients

- Intuitively the best nudge should be in the direction of the largest rate of change (steepness) of the function
- Rate of change  $\rightarrow$  think derivatives



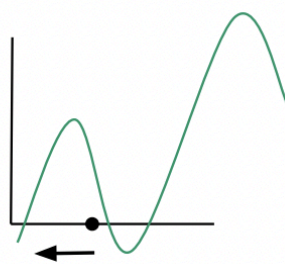
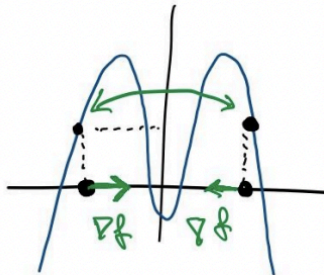
c.

- d. Intuitively, the rate of change of a multi-dimensional function should be a combination of the rate change in each dimension. For a 3-dimensional function, the rate of change would be



$$\nabla f(x, y, z) = \frac{\partial f}{\partial x} \vec{i} + \frac{\partial f}{\partial y} \vec{j} + \frac{\partial f}{\partial z} \vec{k}$$

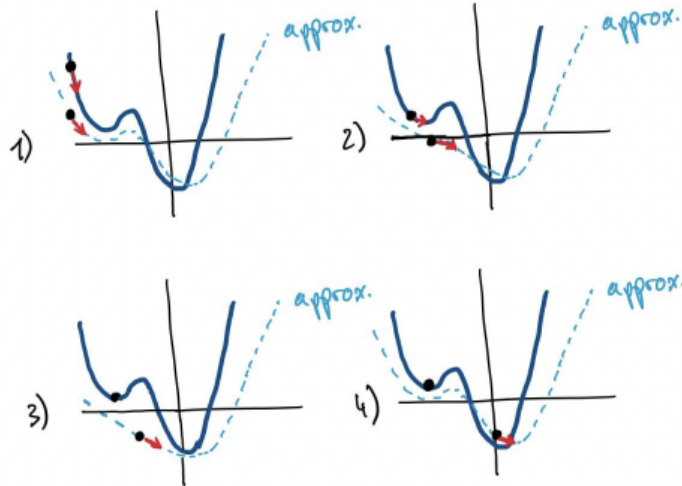
- e. However, the gradient expresses the instantaneous rate of change. At  $p$ , change in  $f_p$  is the steepest but the highest value of  $f$  will depend on how many units we step in that direction. If we step too many units away, the instantaneous change in  $f$  is no longer representative of what values  $f$  will take
- f. Example



- g. Given a smooth function  $f$  for which there exists no closed form solution for finding its maximum, we can find a local maximum through the following steps:
- Define a step size  $a$  (tuning parameter)
  - Initialize  $p$  to be random
  - $p_{\text{new}} = a * \text{change in } f_p + p$
  - $p = p_{\text{new}}$
  - Repeat 3 & 4 until  $p \sim p_{\text{new}}$
- h. To find a local minimum, just use  $-\text{change in } f_p$
- i. Notes about a:
- If  $a$  is too large, GD may overshoot the maximum, take a long time to or never be able to converge
  - If  $a$  is too small, GD may take too long to converge

j. Stochastic Gradient Descent

- i. Recall the cost is computed for the entire dataset. This has some limitations:
  1. Expensive to run
  2. Result we get depends only on the initial starting point
- ii. Goal: Approximate the gradient of the cost using a sample of the data (batch)



3. Note

- a. The magnitude of change in  $f_p$  depends on  $p$ . As  $p$  gets closer to the min/max, the size of change in  $f_p$  decreases
- b. This also means that points  $p$  that contain more information have larger gradients. So the order with which this process is exposed to examples matters