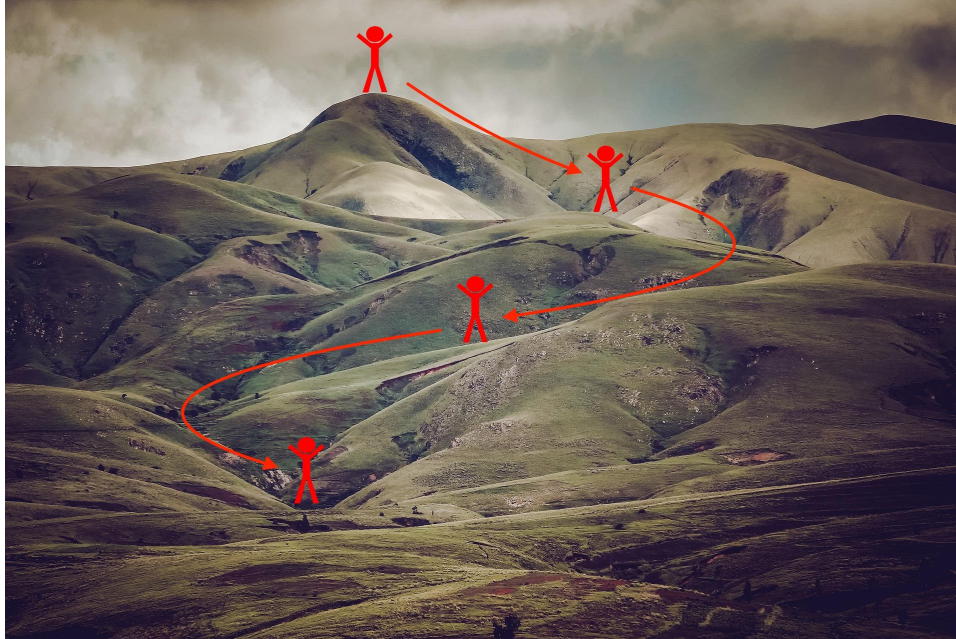


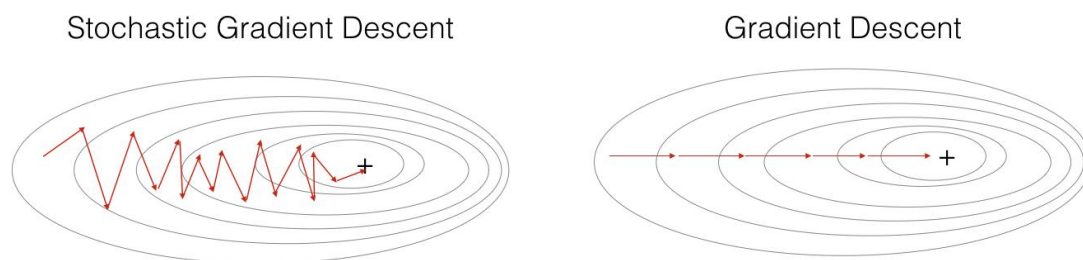
Optimization Methods

Supporting Figures

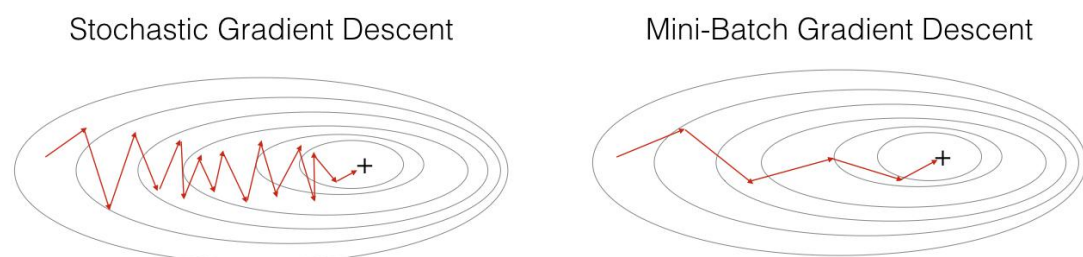
S0_Figure1_Assignment_optimization_methods



S1_Figure1_SGD vs GD



S2_Figure 2_SGD vs Mini-Batch GD



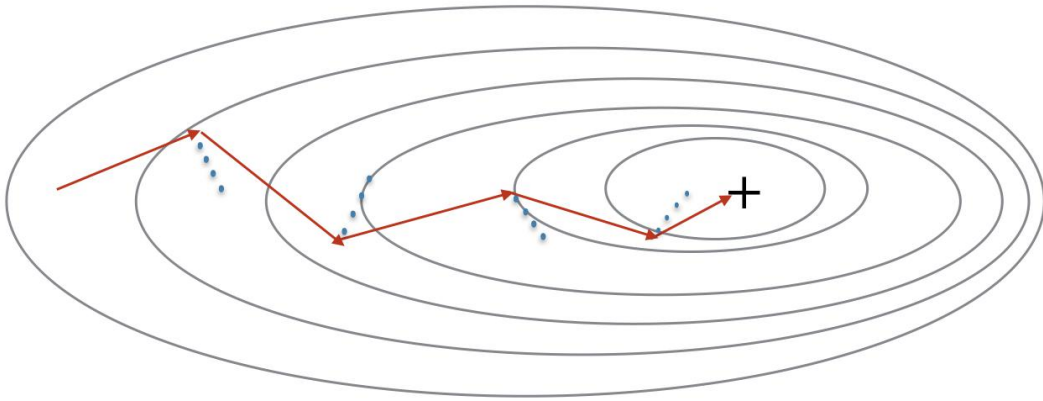
S3_Shuffle

$$\begin{aligned}
 X &= \begin{pmatrix} x_0^{(1)} & x_0^{(2)} & \dots & x_0^{(m-1)} & x_0^{(m)} \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m-1)} & x_1^{(m)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{12286}^{(1)} & x_{12286}^{(2)} & \dots & x_{12286}^{(m-1)} & x_{12286}^{(m)} \\ x_{12287}^{(1)} & x_{12287}^{(2)} & \dots & x_{12287}^{(m-1)} & x_{12287}^{(m)} \end{pmatrix} \\
 Y &= \begin{pmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m-1)} & y^{(m)} \end{pmatrix} \\
 X &= \begin{pmatrix} x_0^{(5)} & x_0^{(16)} & \dots & x_0^{(2)} & x_0^{(m-1)} \\ x_1^{(5)} & x_1^{(16)} & \dots & x_1^{(2)} & x_1^{(m-1)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{12286}^{(5)} & x_{12286}^{(16)} & \dots & x_{12286}^{(2)} & x_{12286}^{(m-1)} \\ x_{12287}^{(5)} & x_{12287}^{(16)} & \dots & x_{12287}^{(2)} & x_{12287}^{(m-1)} \end{pmatrix} \\
 Y &= \begin{pmatrix} y^{(5)} & y^{(16)} & \dots & y^{(2)} & y^{(m-1)} \end{pmatrix}
 \end{aligned}$$

S4_Partion

$$\begin{aligned}
 X &= \begin{array}{|c|c|c|c|c|c|c|c|} \hline 64 \text{ training} & 64 \text{ training} & 64 \text{ training} & \dots & \dots & \dots & 64 \text{ training} & <64 \\ \text{examples} & \text{examples} & \text{examples} & & & & \text{examples} & \text{training} \\ & & & & & & & \text{examples} \\ \hline \end{array} \\
 Y &= \begin{array}{|c|c|c|c|c|c|c|c|} \hline 64 \text{ training} & 64 \text{ training} & 64 \text{ training} & \dots & \dots & \dots & 64 \text{ training} & <64 \\ \text{examples} & \text{examples} & \text{examples} & & & & \text{examples} & \text{training} \\ & & & & & & & \text{examples} \\ \hline \end{array} \\
 &\quad \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \quad \dots \quad \dots \quad \dots \quad \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \\
 &\quad \text{mini_batch} \quad \text{mini_batch} \quad \text{mini_batch} \quad \dots \quad \dots \quad \dots \quad \text{mini_batch} \quad \text{mini_batch} \\
 &\quad \quad \quad 1 \quad \quad \quad 2 \quad \quad \quad 3 \quad \quad \quad \quad \quad \quad \lfloor m/64 \rfloor \quad \lfloor m/64 \rfloor + 1
 \end{aligned}$$

S5_ Figure 3: The red arrows show the direction taken by one step of mini-batch gradient descent with momentum. The blue points show the direction of the gradient (with respect to the current mini-batch) on each step. Rather than just following the gradient, the gradient is allowed to influence v and then take a step in the direction of v .



S6_How does Adam work

The update rule is, for $l = 1, \dots, L$:

$$\begin{cases} v_{dW^{[l]}} = \beta_1 v_{dW^{[l]}} + (1 - \beta_1) \frac{\partial \mathcal{J}}{\partial W^{[l]}} \\ v_{dW^{[l]}}^{corrected} = \frac{v_{dW^{[l]}}}{1 - (\beta_1)^t} \\ s_{dW^{[l]}} = \beta_2 s_{dW^{[l]}} + (1 - \beta_2) \left(\frac{\partial \mathcal{J}}{\partial W^{[l]}} \right)^2 \\ s_{dW^{[l]}}^{corrected} = \frac{s_{dW^{[l]}}}{1 - (\beta_2)^t} \\ W^{[l]} = W^{[l]} - \alpha \frac{v_{dW^{[l]}}^{corrected}}{\sqrt{s_{dW^{[l]}}^{corrected} + \epsilon}} \end{cases}$$

where:

- t counts the number of steps taken of Adam
- L is the number of layers
- β_1 and β_2 are hyperparameters that control the two exponentially weighted averages.
- α is the learning rate
- ϵ is a very small number to avoid dividing by zero

S7_Excercise6_ update_parameters_with_adam

Exercise 6 - update_parameters_with_adam

Now, implement the parameters update with Adam. Recall the general update rule is, for $l = 1, \dots, L$:

$$\begin{cases} v_{dW^{[l]}} = \beta_1 v_{dW^{[l]}} + (1 - \beta_1) \frac{\partial \mathcal{J}}{\partial W^{[l]}} \\ v_{dW^{[l]}}^{corrected} = \frac{v_{dW^{[l]}}}{1 - (\beta_1)^t} \\ s_{dW^{[l]}} = \beta_2 s_{dW^{[l]}} + (1 - \beta_2) \left(\frac{\partial \mathcal{J}}{\partial W^{[l]}} \right)^2 \\ s_{dW^{[l]}}^{corrected} = \frac{s_{dW^{[l]}}}{1 - (\beta_2)^t} \\ W^{[l]} = W^{[l]} - \alpha \frac{v_{dW^{[l]}}^{corrected}}{\sqrt{s_{dW^{[l]}}^{corrected} + \epsilon}} \end{cases}$$

Note that the iterator `l` starts at 1 in the `for` loop as the first parameters are $W^{[1]}$ and $b^{[1]}$.