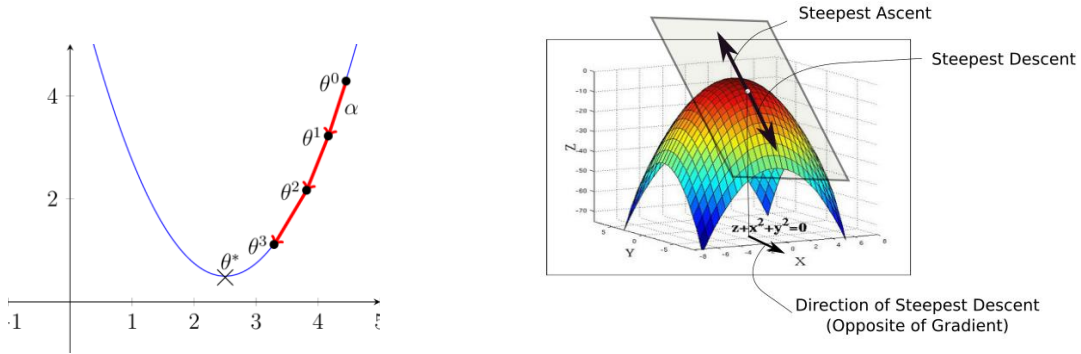


# Math foundations of ML – lesson 10

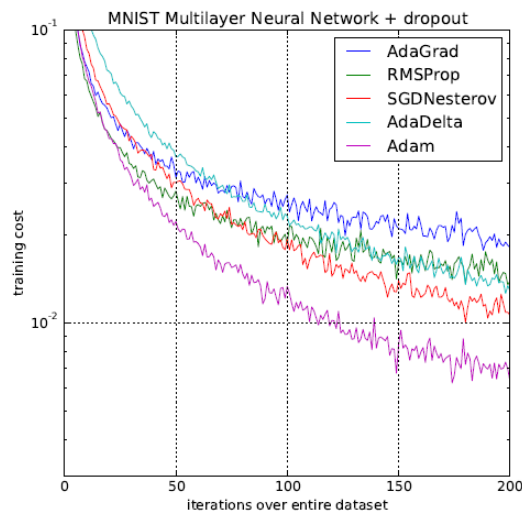
## Gradient descent



## Momentum SGD

$$\begin{aligned}
 v_t &= \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta^{(t)}) \\
 &= \gamma (\gamma v_{t-2} + \eta \nabla_{\theta} J(\theta^{(t-1)})) + \eta \nabla_{\theta} J(\theta^{(t)}) \\
 &= \gamma^2 v_{t-2} + \gamma \eta \nabla_{\theta} J(\theta^{(t-1)}) + \eta \nabla_{\theta} J(\theta^{(t)}) \\
 &\dots \\
 &= \gamma^n v_{t-n} + \eta (\gamma^{n-1} \nabla_{\theta} J(\theta^{(t-n+1)}) + \dots + \gamma \nabla_{\theta} J(\theta^{(t-1)}) + \nabla_{\theta} J(\theta^{(t)}))
 \end{aligned}$$

## Adam



## Loss functions for segmentation

Let  $p = (p_1, \dots, p_n)$  be the ground truth labels of the segmentation. This means that  $p_i \in \{0, 1\}$ . Let  $\tilde{p} = (\tilde{p}_1, \dots, \tilde{p}_n)$  be an approximation (e.g. as generated through a deep learning network) of the segmentation, with  $0 \leq \tilde{p}_i \leq 1$ . One can ensure that the network outputs ‘probability pixels’ by applying at the last layer at each pixel the logistic function

$$\sigma(t) := \frac{1}{1 + e^{-t}}.$$

### Negative log-likelihood as a loss for image segmentation

We can define a loss per image

$$-\sum_{i=1}^n p_i \log \tilde{p}_i + (1 - p_i) \log(1 - \tilde{p}_i)$$

### Simple differentiable approximations to the Jaccard index

It is clear that in unbalanced cases, where there are more ground truth background pixels than ground truth object pixels, the standard negative log-likelihood loss is potentially biased.

One would like to train a DL segmentation network and measure/maximize performance using Intersection Over Union (IoU) loss, which is also called the Jaccard loss

$$0 \leq \frac{|A \cap B|}{|A \cup B|} \leq 1.$$

The Jaccard index in our special is

$$J(p, \tilde{p}) = \frac{\#\{i : p_i = 1 \text{ and } \tilde{p}_i \geq 0.5\}}{\#\{i : p_i = 1 \text{ or } \tilde{p}_i \geq 0.5\}}.$$

Since the Jaccard index is not differentiable, we look for surrogates for which we can compute gradients and thus optimize. We use an approximation to the Jaccard index which is differentiable and based on the equality

$$\frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|},$$

Let

$$\langle p, \tilde{p} \rangle = \sum_{i=1}^n p_i \tilde{p}_i, \quad |p|_1 = \sum_{i=1}^n p_i, \quad |\tilde{p}|_1 = \sum_{i=1}^n \tilde{p}_i.$$

We define a differentiable ‘loss’ function that we wish to maximize

$$\tilde{J}(p, \tilde{p}) := \frac{\langle p, \tilde{p} \rangle}{|p|_1 + |\tilde{p}|_1 - \langle p, \tilde{p} \rangle}.$$

Note that in the special case where  $\tilde{p}_i \in \{0, 1\}$ ,  $i = 1, \dots, n$ , we get

$$\tilde{J}(p, \tilde{p}) = J(p, \tilde{p}).$$

The two losses are sometimes combined with a weight. There are other more complicated surrogates (e.g. Lovász).