

# Deep Learning - Foundations and Concepts

## Chapter 7. Gradient Descent

nonlineark@github

February 24, 2025

# Outline

1 Error Surfaces

2 Gradient Descent Optimization

# Gradient and stationary points

## Theorem

Let the function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be differentiable at  $a \in \mathbb{R}^n$ :

- 1 Near  $a$  the function  $f$  increases fastest in the direction of  $\nabla f(a) \in \mathbb{R}^n$ .
- 2 The rate of increase in  $f$  is measured by the length of  $\nabla f(a)$ .
- 3 If  $f$  has a local extremum at  $a$  then  $\nabla f(a) = 0$ .

## Definition

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be differentiable at  $a \in \mathbb{R}^n$ . Then  $a$  is said to be a stationary point for  $f$  if  $Df(a) = 0$ , or, equivalently,  $\nabla f(a) = 0$ .

# Gradient and stationary points

It's easy to see the correctness of the theorem. Since the rate of increase of  $f$  at the point  $a$  in an arbitrary direction  $v \in \mathbb{R}^n$  is given by the directional derivative at  $v$ , we have:

$$|Df(a)v| = |v^T \nabla f(a)| \leq \|\nabla f(a)\| \|v\|$$

where we have used the Cauchy-Schwarz inequality. The rate of increase is maximal if  $v$  is a positive scalar multiple of  $\nabla f(a)$ . For the third claim, let's define  $g_j$  as:

$$g_j(t) = f(a_1, \dots, a_{j-1}, t, a_{j+1}, \dots, a_n)$$

then  $g'_j(a_j) = D_j f(a)$ . Since  $f$  has a local extremum at  $a$ ,  $g_j$  also has a local extremum at  $a_j$ . Thus  $g'_j(a_j) = 0$  for  $1 \leq j \leq n$ , and we have  $\nabla f(a) = 0$ .

# Gradient and stationary points

During training, we want to optimize the weights and biases  $w \in \mathbb{R}^W$  by using a chosen error function  $E(w)$ . From the previous theorem we see that, its smallest value will occur at a point in weight space such that:

$$\nabla E(w) = 0$$

But:

- Global minimum vs. local minimum.
- For any point  $w$  that is a local minimum, there will generally be other points in weight space that are equivalent minima (weight-space symmetries).

# Local quadratic approximation

## Theorem

Let  $U$  be a convex open subset of  $\mathbb{R}^n$  and let  $a \in U$  be a stationary point for  $f \in C^2(U)$ . Then we have the following assertions:

- ❶ If  $Hf(a)$  is positive definite, then  $f$  has a local strict minimum at  $a$ .
  - ❷ If  $Hf(a)$  is negative definite, then  $f$  has a local strict maximum at  $a$ .
- where  $Hf(a)$  is the Hessian of  $f$  at  $a$ .

# Local quadratic approximation

We only prove the first claim. Using Taylor expansion, we see that:

$$\begin{aligned} f(a+h) &= f(a) + h^T \nabla f(a) + \frac{1}{2} h^T H f(a) h + R_2(a, h) \\ &= f(a) + \frac{1}{2} h^T H f(a) h + R_2(a, h) \end{aligned}$$

where  $\lim_{h \rightarrow 0} \frac{R_2(a, h)}{\|h\|^2} = 0$ . Since  $f \in C^2(U)$ ,  $Hf(a)$  is a self-adjoint operator. Let  $\lambda$  be its smallest eigenvalue. Because  $Hf(a)$  is positive definite,  $\lambda > 0$ . Notice that:

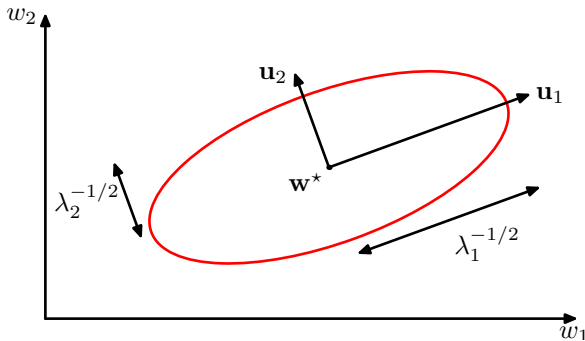
- $h^T H f(a) h \geq \lambda \|h\|^2$ .
- There is  $\delta > 0$ , such that  $\frac{|R_2(a, h)|}{\|h\|^2} < \frac{\lambda}{4}$  for  $\|h\| < \delta$ .

For  $\|h\| < \delta$ , we have:

$$f(a+h) - f(a) = \frac{1}{2} h^T H f(a) h + R_2(a, h) > \frac{\lambda}{2} \|h\|^2 - \frac{\lambda}{4} \|h\|^2 = \frac{\lambda}{4} \|h\|^2$$

From the previous theorem, we see that: A necessary and sufficient condition for  $w^*$  to be a local minimum of the error function  $E(w)$  is that the gradient of  $E(w)$  should vanish at  $w^*$  and the Hessian matrix evaluated at  $w^*$  should be positive definite.

Figure: Geometry of the error surface in the neighbourhood of a minimum  $w^*$





# Use of gradient information

There is no way we can find an analytical solution to the equation  $\nabla E(w) = 0$ , so we resort to iterative numerical procedures to minimize the error function:

$$w^{(\tau)} = w^{(\tau-1)} + \Delta w^{(\tau-1)}$$

# Use of gradient information

Let's estimate the computational effort needed. The Taylor expansion of  $E(w)$  around some point  $\hat{w}$  in weight space is given by:

$$E(w) \approx E(\hat{w}) + (w - \hat{w})^T \nabla E(\hat{w}) + \frac{1}{2} (w - \hat{w})^T H E(\hat{w}) (w - \hat{w})$$

We see that in order to locate the minimum, we need at least  $\mathcal{O}(W^2)$  pieces of information, because:

- $\nabla E(\hat{w})$  contains  $W$  independent parameters.
- $H E(\hat{w})$  contains  $\frac{W(W+1)}{2}$  independent parameters.

# Use of gradient information

- If we do not make use of gradient information, we would expect to have to perform  $\mathcal{O}(W^2)$  function evaluations, each of which would require  $\mathcal{O}(W)$  steps. Thus the total cost would be  $\mathcal{O}(W^3)$ .
- If we make use of the gradient information, each evaluation of  $\nabla E$  brings  $W$  pieces of information, so we would expect to find the minimum of the function in  $\mathcal{O}(W)$  gradient evaluations. In the next chapter we shall see, by using error backpropagation, each such evaluation takes only  $\mathcal{O}(W)$  steps. Thus the total cost would be  $\mathcal{O}(W^2)$ .

# Batch gradient descent

The simplest approach to using gradient information is to take a small step in the direction of the negative gradient:

$$w^{(\tau)} = w^{(\tau-1)} - \eta \nabla E(w^{(\tau-1)})$$

where the parameter  $\eta > 0$  is known as the learning rate. At each step:

- The weight vector is moved in the direction of the greatest rate of decrease of the error function ("gradient descent" in the name).
- The entire training set will be processed ("batch" in the name).

# Stochastic gradient descent

Note that error functions based on maximum likelihood for a set of independent observations comprise a sum of terms:  $E(w) = \sum_{n=1}^N E_n(w)$ . We can update the weight vector based on one data point at a time:

$$w^{(\tau)} = w^{(\tau-1)} - \eta \nabla E_n(w^{(\tau-1)})$$

Advantages of stochastic gradient descent:

- It is much more efficient on very large data sets.
- It works when the data arises from a continuous stream of new data points (a.k.a., online algorithm).
- It handles redundancy in the data much more efficiently.
- The possibility of escaping from local minima.

# Stochastic gradient descent

---

**Algorithm 1:** Stochastic gradient descent

---

 $n \leftarrow 1;$ **repeat** $\quad w \leftarrow w - \eta \nabla E_n(w);$   
 $\quad n \leftarrow n + 1(\text{mod} N);$ **until** *convergence*;**return**  $w;$ 

---

# Mini-batches

---

**Algorithm 2:** Mini-batch stochastic gradient descent

---

```
 $n \leftarrow 1;$   
repeat  
   $w \leftarrow w - \eta \nabla E_{n:n+B-1}(w);$   
   $n \leftarrow n + B;$   
  if  $n > N$  then  
    shuffle data;  
     $n \leftarrow 1;$   
  end  
until convergence;  
return  $w;$ 
```

---

# Parameter initialization

The initial setting for the parameters being learned can have a significant effect on:

- How long it takes to reach a solution.
- The generalization performance of the resulting trained network.

Although there is relatively little theory to guide the initialization strategy, one key consideration is symmetry breaking, by initializing parameters randomly from some distribution:

- Uniform distribution in the range  $[-\epsilon, \epsilon]$ .
- Zero-mean Gaussian of the form  $\mathcal{N}(0, \epsilon^2)$ .



# Parameter initialization

The approach of He initialization for choosing the value of  $\epsilon$ . Consider a network in which layer  $l$  evaluates the following transformations:

$$a_i^{(l)} = \sum_{m=1}^M w_{im} z_m^{(l-1)}$$
$$z_i^{(l)} = \text{ReLU}(a_i^{(l)})$$

Suppose we initialize the weights using a Gaussian  $\mathcal{N}(0, \epsilon^2)$ , and suppose that the pre-activations  $a_m^{(l-1)}$  in layer  $l - 1$  have variance  $\lambda^2$ .

# Parameter initialization

We see that:

$$\begin{aligned}
 E(a_i^{(l)}) &= \sum_{m=1}^M E(w_{im})E(z_m^{(l-1)}) = 0 \\
 \text{var}(a_i^{(l)}) &= E((a_i^{(l)})^2) - (E(a_i^{(l)}))^2 = E((a_i^{(l)})^2) \\
 &= E\left(\sum_{1 \leq m, n \leq M} w_{im}w_{in}z_m^{(l-1)}z_n^{(l-1)}\right) = E\left(\sum_{m=1}^M w_{im}^2(z_m^{(l-1)})^2\right) \\
 &= \sum_{m=1}^M E(w_{im}^2)E((z_m^{(l-1)})^2) = \sum_{m=1}^M \epsilon^2 \frac{1}{2} E((a_m^{(l-1)})^2) \\
 &= \sum_{m=1}^M \epsilon^2 \frac{1}{2} \lambda^2 = \frac{M}{2} \epsilon^2 \lambda^2
 \end{aligned}$$

# Parameter initialization

If we require that the pre-activations at layer  $l$  also have variance  $\lambda^2$  then we arrive at the following choice for  $\epsilon$  to initialize the weights that feed into a unit with  $M$  inputs:

$$\epsilon = \sqrt{\frac{2}{M}}$$

The bias parameters are typically set to small positive values to ensure that most pre-activations are initially active during learning.