

Gradient Descent (and Beyond)

Cornell CS 4/5780

Fall 2021

We want to minimize a convex, continuous and differentiable loss function $\ell(w)$. In this section we discuss two of the most popular "hill-climbing" algorithms, gradient descent and Newton's method.

Algorithm:

Initialize \mathbf{w}_0

Repeat until converge:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{s}$$

If $\|\mathbf{w}^{t+1} - \mathbf{w}^t\|_2 < \epsilon$, converged!

Trick: Taylor Expansion

How can you minimize a function ℓ if you don't know much about it? The trick is to assume it is much simpler than it really is. This can be done with Taylor's approximation. Provided that the norm $\|\mathbf{s}\|_2$ is small (*i.e.* $\mathbf{w} + \mathbf{s}$ is very close to \mathbf{w}), we can approximate the function $\ell(\mathbf{w} + \mathbf{s})$ by its first and second derivatives:

$$\ell(\mathbf{w} + \mathbf{s}) \approx \ell(\mathbf{w}) + g(\mathbf{w})^\top \mathbf{s}$$

$$\ell(\mathbf{w} + \mathbf{s}) \approx \ell(\mathbf{w}) + g(\mathbf{w})^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top H(\mathbf{w}) \mathbf{s}$$

Here, $g(\mathbf{w}) = \nabla \ell(\mathbf{w})$ is the gradient and $H(\mathbf{w}) = \nabla^2 \ell(\mathbf{w})$ is the Hessian of ℓ . Both approximations are valid if $\|\mathbf{s}\|_2$ is small, but the second one assumes that ℓ is twice differentiable and is more expensive to compute but also more accurate than only using gradient.

Gradient Descent: Use the first order approximation

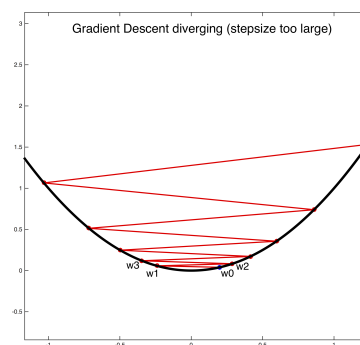
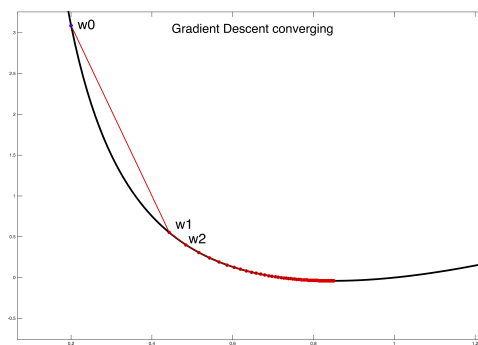
In gradient descent we only use the gradient (first order). In other words, we assume that the function ℓ around \mathbf{w} is linear and behaves like $\ell(\mathbf{w}) + g(\mathbf{w})^\top \mathbf{s}$. Our goal is to find a vector \mathbf{s} that minimizes this function. In steepest descent we simply set

$$\mathbf{s} = -\alpha g(\mathbf{w}),$$

for some small $\alpha > 0$. It is straight-forward to prove that in this case $\ell(\mathbf{w} + \mathbf{s}) < \ell(\mathbf{w})$.

$$\underbrace{\ell(\mathbf{w} + (-\alpha g(\mathbf{w})))}_{\text{after one update}} \approx \ell(\mathbf{w}) - \underbrace{\alpha g(\mathbf{w})^\top g(\mathbf{w})}_{>0} < \underbrace{\ell(\mathbf{w})}_{\text{before}}$$

Setting the **learning rate** $\alpha > 0$ is a dark art. Only if it is sufficiently small will gradient descent converge (see the first figure below). If it is too large the algorithm can easily *diverge* out of control (see the second figure below). A safe (but sometimes slow) choice is to set $\alpha = \frac{t_0}{t}$, which guarantees that it will eventually become small enough to converge (for any initial value $t_0 > 0$).



Adagrad

One option is to set the step-size adaptively for *every feature*. Adagrad keeps a running average of the squared gradient magnitude and sets a small learning rate for features that have large gradients, and a large learning rate for features with small gradients. Setting different learning rates for different features is particularly important if they are of different scale or vary in frequency. For

example, word counts can differ a lot across common words and rare words.

Adagrad Algorithm:

Initialize \mathbf{w}_0 and \mathbf{z} : $\forall d: w_d^0 = 0$ and $z_d = 0$

Repeat until converge:

$\mathbf{g} = \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}}$ # Compute gradient

$\forall d: z_d \leftarrow z_d + g_d^2$

$\forall d: w_d^{t+1} \leftarrow w_d^t - \alpha \frac{g_d}{\sqrt{z_d + \epsilon}}$

If $\|\mathbf{w}^{t+1} - \mathbf{w}^t\|_2 < \delta$, **converged!** # for some small $\delta > 0$.

Newton's Method: Use 2nd order Approximation

Newton's method assumes that the loss ℓ is twice differentiable and uses the approximation with Hessian (2nd order Taylor approximation). The **Hessian Matrix** contains all second order partial derivatives and is defined as

$$H(\mathbf{w}) = \begin{pmatrix} \frac{\partial^2 \ell}{\partial w_1^2} & \frac{\partial^2 \ell}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 \ell}{\partial w_1 \partial w_n} \\ \vdots & \cdots & \cdots & \vdots \\ \frac{\partial^2 \ell}{\partial w_n \partial w_1} & \cdots & \cdots & \frac{\partial^2 \ell}{\partial w_n^2} \end{pmatrix},$$

and, because the convexity of ℓ , it is always a symmetric square matrix and positive semi-definite.

Note: A symmetric matrix \mathbf{M} is positive semi-definite if it has only non-negative eigenvalues or, equivalently, for any vector \mathbf{x} we must have $\mathbf{x}^\top \mathbf{M} \mathbf{x} \geq 0$.

It follows that the approximation

$$\ell(\mathbf{w} + \mathbf{s}) \approx \ell(\mathbf{w}) + g(\mathbf{w})^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top H(\mathbf{w}) \mathbf{s}$$

describes a convex parabola, and we can find its minimum by solving the following optimization problem:

$$\underset{\mathbf{s}}{\operatorname{argmin}} \ell(\mathbf{w}) + g(\mathbf{w})^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top H(\mathbf{w}) \mathbf{s}$$

To find the minimum of the objective, we take its first derivative with respect to \mathbf{s} , equate it with zero, and solve for \mathbf{s} :

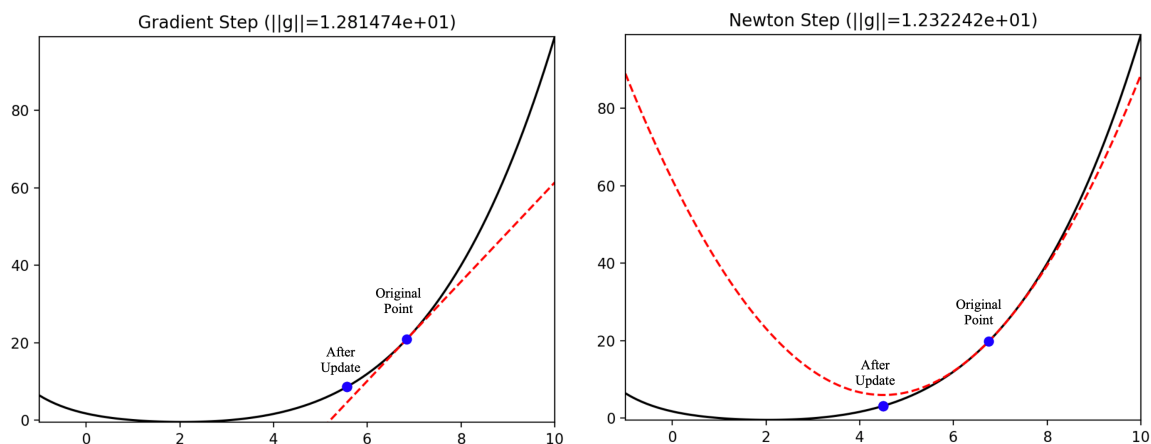
$$g(\mathbf{w}) + H(\mathbf{w})\mathbf{s} = 0$$

$$\Rightarrow \mathbf{s} = -[H(\mathbf{w})]^{-1}g(\mathbf{w}).$$

This choice of \mathbf{s} converges extremely fast if the approximation is sufficiently accurate and the resulting step sufficiently small. Otherwise it can diverge. Divergence often happens if the function is flat or almost flat with respect to some dimension. In that case the second derivatives are close to zero, and their inverse becomes very large - resulting in gigantic steps. Different from gradient descent, here there is no step-size that guarantees that steps are all small and local. As the Taylor approximation is only accurate locally, large steps can move the current estimates far from regions where the Taylor approximation is accurate.

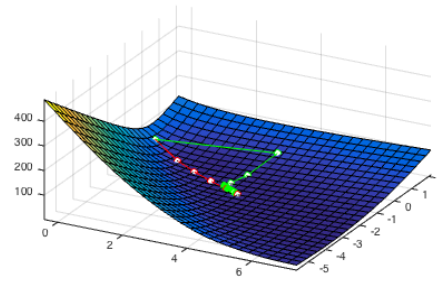
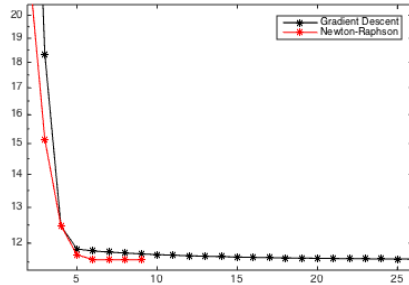
Best practices

1. The matrix $H(\mathbf{w})$ scales $d \times d$ and is expensive to compute. A good approximation can be to only compute its diagonal entries and multiply the update with a small step-size. Essentially you are then doing a hybrid between Newton's method and gradient descent, where you weigh the step-size for each dimension by the inverse Hessian.
2. To avoid divergence of Newton's method, a good approach is to start with gradient descent (or even stochastic gradient descent) and then finish the optimization Newton's method. Typically, the second order approximation, used by Newton's Method, is more likely to be appropriate near the optimum.

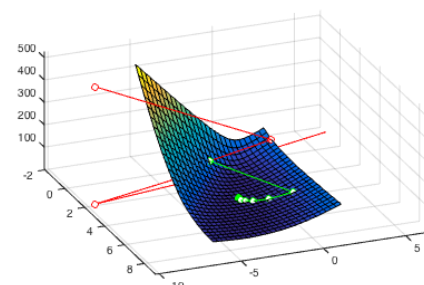
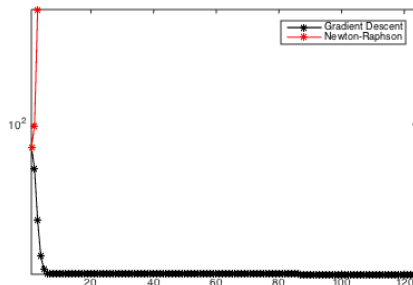


A gradient descent step (left) and a Newton step (right) on the same function. The loss function is depicted in black, the approximation as a dotted red line. The gradient step moves the point downwards along the linear approximation

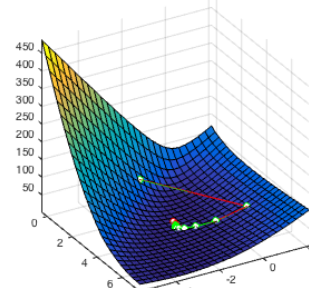
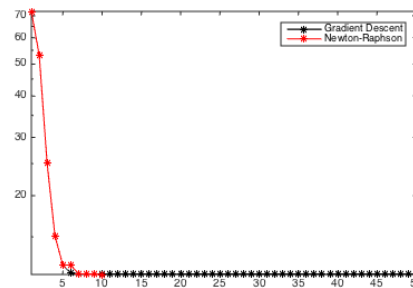
of the function. The Newton step moves the point to the minimum of the parabola, which is used to approximate the function.



(a) A starting point where Newton's Method converges in 8 iterations.



(b) A starting point where Newton's Method diverges.



(c) same starting point as in Figure 2, however Newton's method is only used after 6 gradient steps and converges in a few steps.

The three plots show a comparison of Newton's Method and Gradient Descent. Gradient Descent always converges after over 100 iterations from all initial starting points. If it converges (Figure 1), Newton's Method is much faster (convergence after 8 iterations) but it can diverge (Figure 2). Figure 3 shows the hybrid approach of taking 6 gradient descent steps and then switching to Newton's Method. It still converges in only 10 updates.