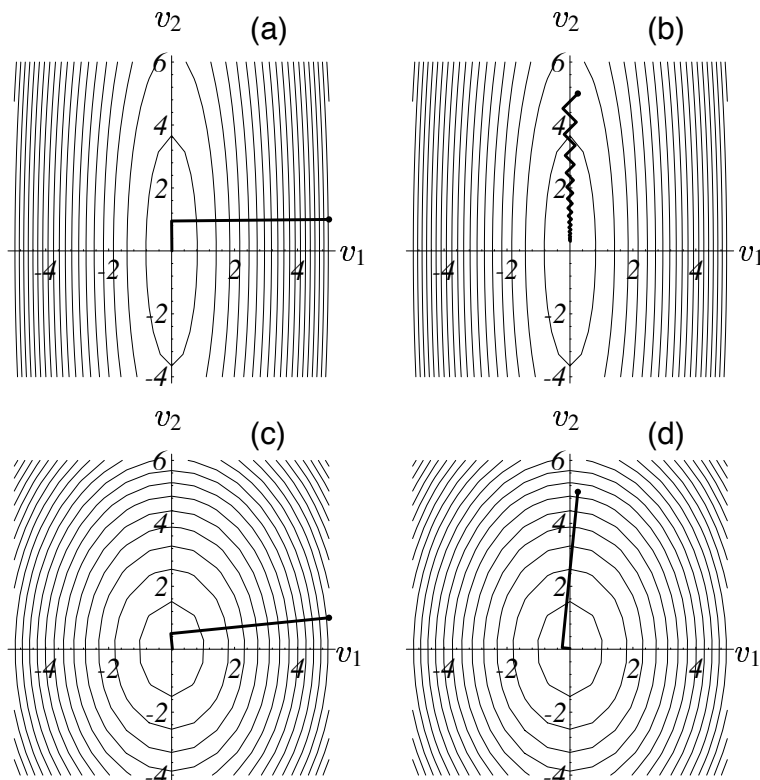


Convergence of gradient descent

The effectiveness of gradient descent depends critically on the “conditioning” of \mathbf{H} and the starting point. Consider the two examples below:



(from Shewchuk, “... without the agonizing pain”)

When the conditioning of \mathbf{H} is poor, which here corresponds to the case where the ellipses denoting the level sets of our objective function are more eccentric or “squished”, and we choose a bad starting point, convergence can take many iterations even in simple cases.

We can make this a bit more precise if we define mathematically what we really mean by the conditioning of \mathbf{H} . The **condition number** of a matrix \mathbf{H} , typically denoted $\kappa(\mathbf{H})$ is the ratio of the largest to smallest singular values of \mathbf{H} :

$$\kappa(\mathbf{H}) = \frac{\sigma_{\max}(\mathbf{H})}{\sigma_{\min}(\mathbf{H})}.$$

Note that by the $\sigma_{\max}(\mathbf{H})$ we mean the largest singular value and by $\sigma_{\min}(\mathbf{H})$ we mean the smallest *non-zero* singular value, i.e., σ_R where R is the rank of \mathbf{H} . For the case where \mathbf{H} is a square matrix (as it is in our context), we can also equivalently write

$$\kappa(\mathbf{H}) = \frac{\lambda_{\max}(\mathbf{H})}{\lambda_{\min}(\mathbf{H})},$$

where $\lambda_{\max}(\mathbf{H})$ and $\lambda_{\min}(\mathbf{H})$ denote the largest and smallest eigenvalues of \mathbf{H} , respectively. The condition number is a natural way of quantifying just how sensitive we are going to be to noise, but it also plays a key role in determining how computationally challenging it will be to solve the least squares problem using iterative methods.

Specifically, below we will provide a bound that shows how $f(\mathbf{x}_k)$ approaches $f(\mathbf{x}^*)$, where \mathbf{x}^* denotes the minimizer of f . Specifically, we will show that

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{1}{\kappa(\mathbf{H})}\right) (f(\mathbf{x}_k) - f(\mathbf{x}^*)). \quad (1)$$

Let's think a bit about what this says. Note that $1 - 1/\kappa(\mathbf{H})$ is always less than 1, so each iteration makes *some* progress. If $\kappa(\mathbf{H}) \leq 2$, then at each iteration we make *a lot* of progress – cutting the error in half with each iteration. However, if $\kappa(\mathbf{H})$ is very large, this constant becomes very close to 1, indicating only minor improvements.

Convergence analysis

Recall that we are trying to minimize

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{x}^T \mathbf{b}.$$

Our convergence analysis will rely on one very useful property of $f(\mathbf{x})$, namely that we can write¹

$$f(\mathbf{y}) = f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \mathbf{H}(\mathbf{y} - \mathbf{x}). \quad (2)$$

We can easily verify this by just plugging in $\nabla f(\mathbf{x}) = \mathbf{H} \mathbf{x} - \mathbf{b}$ and simplifying. Specifically, note that we can equivalently write (2) as

$$f(\mathbf{y}) - f(\mathbf{x}) = (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \mathbf{H}(\mathbf{y} - \mathbf{x}).$$

The right-hand side of this equation can be simplified as

$$\begin{aligned} & (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \mathbf{H}(\mathbf{y} - \mathbf{x}) \\ &= (\mathbf{y} - \mathbf{x})^T (\mathbf{H} \mathbf{x} - \mathbf{b}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \mathbf{H}(\mathbf{y} - \mathbf{x}) \\ &= \mathbf{y}^T \mathbf{H} \mathbf{x} - \mathbf{y}^T \mathbf{b} - \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{x}^T \mathbf{b} + \frac{1}{2}(\mathbf{y}^T \mathbf{H} \mathbf{y} + \mathbf{x}^T \mathbf{H} \mathbf{x} - 2\mathbf{x}^T \mathbf{H} \mathbf{y}) \\ &= \frac{1}{2}\mathbf{y}^T \mathbf{H} \mathbf{x} - \mathbf{y}^T \mathbf{b} - \frac{1}{2}\mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{x}^T \mathbf{b} \\ &= f(\mathbf{y}) - f(\mathbf{x}), \end{aligned}$$

as desired.

¹This is like taking a second-order Taylor approximation to f around the point \mathbf{x} , but since f is a quadratic function this is not an approximation but exact.

Equation (2) immediately tells us something about how much progress we make at each iteration. If we plug in $\mathbf{y} = \mathbf{x}_{k+1}$ and $\mathbf{x} = \mathbf{x}_k$ to (2), in which case $\mathbf{y} - \mathbf{x} = \mathbf{x}_{k+1} - \mathbf{x}_k = -\alpha_k \nabla f(\mathbf{x}_k)$, we obtain

$$f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k) - \alpha_k \|\nabla f(\mathbf{x}_k)\|_2^2 + \frac{\alpha_k^2}{2} (\nabla f(\mathbf{x}_k))^T \mathbf{H} \nabla f(\mathbf{x}_k). \quad (3)$$

Recall that we had set

$$\alpha_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{H} \mathbf{r}_k},$$

where $\mathbf{r}_k = -\nabla f(\mathbf{x}_k)$. Plugging this into (3) yields

$$\begin{aligned} f(\mathbf{x}_{k+1}) &= f(\mathbf{x}_k) - \frac{\|\mathbf{r}_k\|_2^4}{\mathbf{r}_k^T \mathbf{H} \mathbf{r}_k} + \frac{1}{2} \left(\frac{\|\mathbf{r}_k\|_2^2}{\mathbf{r}_k^T \mathbf{H} \mathbf{r}_k} \right)^2 \mathbf{r}_k^T \mathbf{H} \mathbf{r}_k \\ &= f(\mathbf{x}_k) - \frac{1}{2} \frac{\|\mathbf{r}_k\|_2^4}{\mathbf{r}_k^T \mathbf{H} \mathbf{r}_k}. \end{aligned}$$

This tells us that we are guaranteed to make at least *some* progress at each iteration. Precisely how much depends on this rather strange looking function of \mathbf{r}_k , but we can actually get a much simpler expression by recalling that for any symmetric, positive semidefinite matrix \mathbf{H} we have that

$$\lambda_{\min}(\mathbf{H}) \leq \frac{\mathbf{x}^T \mathbf{H} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \leq \lambda_{\max}(\mathbf{H}) \quad (4)$$

for all \mathbf{x} , where $\lambda_{\max}(\mathbf{H})$ and $\lambda_{\min}(\mathbf{H})$ denote the largest and smallest eigenvalues of \mathbf{H} , respectively. This is a fact that we essentially proved in the discussion of least squares in noise, although we did not explicitly state this at the time. Using the upper half of (4) we can get a simpler bound on how much progress we make at each iteration:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{1}{2\lambda_{\max}(\mathbf{H})} \|\mathbf{r}_k\|_2^2. \quad (5)$$

This bound is nice, but it would be even better if we could say something concrete on how large $\|\mathbf{r}_k\|_2^2$ will be. In particular, our intuition should be that if we are far from the solution, the gradient (or \mathbf{r}_k) should be large. There is a clever way to prove exactly this. First, note that (4) applied to (2) also yields

$$f(\mathbf{y}) \geq f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^\top \nabla f(\mathbf{x}) + \frac{\lambda_{\min}(\mathbf{H})}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

We can obtain a simpler lower bound for $f(\mathbf{y})$ by determining the smallest value that the right-hand side of this could ever take over all possible choices of \mathbf{y} . To do this, we simply minimize this lower bound by taking the gradient with respect to \mathbf{y} and setting it equal to zero:

$$\nabla f(\mathbf{x}) + \lambda_{\min}(\mathbf{H})(\mathbf{y} - \mathbf{x}) = 0,$$

From this we obtain that the lower bound will be minimized by

$$\mathbf{y} - \mathbf{x} = -\frac{1}{\lambda_{\min}(\mathbf{H})} \nabla f(\mathbf{x}).$$

Plugging this in yields

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}) - \frac{1}{\lambda_{\min}(\mathbf{H})} \|\nabla f(\mathbf{x})\|_2^2 + \frac{1}{2\lambda_{\min}(\mathbf{H})} \|\nabla f(\mathbf{x})\|_2^2 \\ &= f(\mathbf{x}) - \frac{1}{2\lambda_{\min}(\mathbf{H})} \|\nabla f(\mathbf{x})\|_2^2. \end{aligned}$$

In particular, this applies when $\mathbf{y} = \mathbf{x}^*$ (where \mathbf{x}^* denotes the minimizer of $f(\mathbf{x})$), which after some rearranging yields

$$\|\nabla f(\mathbf{x})\|_2^2 \geq 2\lambda_{\min}(\mathbf{H}) (f(\mathbf{x}) - f(\mathbf{x}^*)). \quad (\text{PL})$$

This is a famous and useful result, often referred to as the **Polyak-Łojasiewicz inequality**.

From our previous bound in (5) we have that

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq f(\mathbf{x}_k) - f(\mathbf{x}^*) - \frac{1}{2\lambda_{\max}(\mathbf{H})} \|\mathbf{r}_k\|_2^2.$$

Combining this with the PL inequality we obtain

$$\begin{aligned} f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) &\leq f(\mathbf{x}_k) - f(\mathbf{x}^*) - \frac{\lambda_{\min}(\mathbf{H})}{\lambda_{\max}(\mathbf{H})} (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \\ &= \left(1 - \frac{\lambda_{\min}(\mathbf{H})}{\lambda_{\max}(\mathbf{H})}\right) (f(\mathbf{x}_k) - f(\mathbf{x}^*)). \end{aligned}$$

That is, the gap between the current value of the objective function and the optimal value is cut down by a factor of $1 - 1/\kappa(\mathbf{H}) < 1$ at each iteration.