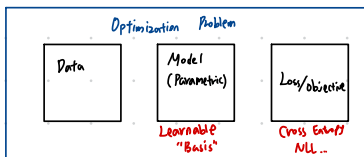


Gradient descent



Optimizer
(model parameters)
GD

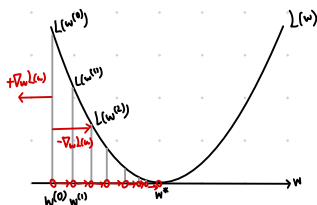
Gradient descent solution:

$$w^{(t+1)} = w^{(t)} - \epsilon \nabla_w L(w^{(t)})$$

where:

$$\nabla_w L(w^{(t)}) = \Phi \Phi^T w^{(t)} - \Phi y$$

$$L(w) = \|\Phi^T w - y\|^2$$



How should we select the learning rate?

1. Let's dive deep into the inner workings of gradient descent.

Let $g^{(t)} = \nabla_w L(w^{(t)})$ and let H be the Hessian matrix s.t.,

$$[H^{(t)}]_{ij} = \frac{\partial^2 L}{\partial w_i \partial w_j} (w^{(t)})$$

ex)
$$\begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

Recall the **Taylor Series** of a function (second order).

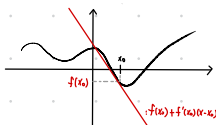
Taylor series of a function is an infinite sum of terms that are expressed in terms of the functions derivatives at a single point.

For most common functions, the function and the sum of its Taylor series are equal near this point.

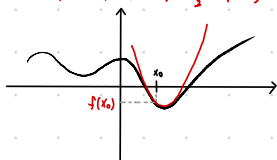
Definition: $f(x) \approx f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x_0)}{2!}(x-x_0)^2 + \dots + \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!}(x-x_0)^n$.

ex) $f: \mathbb{R} \rightarrow \mathbb{R}$

$f(x) \approx f(x_0) + f'(x_0)(x-x_0)$: First Order Taylor Series



$f(x) \approx f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x_0)}{2}(x-x_0)^2$: Second Order Taylor Series



Recap:

Examples of Hessian Matrices

① Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$H_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

② Let $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$$H_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}$$

let $f(x,y) = x^2 y^3$

$\frac{\partial f}{\partial x} = \frac{\partial}{\partial x}(x^2 y^3) = 2xy^3$

$\frac{\partial f}{\partial y} = \frac{\partial}{\partial y}(x^2 y^3) = 3x^2 y^2$

$\frac{\partial^2 f}{\partial x^2} = \frac{\partial}{\partial x}(\frac{\partial f}{\partial x}) = \frac{\partial}{\partial x}(2xy^3) = 2y^3$

$\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial}{\partial x}(\frac{\partial f}{\partial y}) = \frac{\partial}{\partial x}(3x^2 y^2) = 6xy^2$

③ Let $g(x,y) = x^2 + 2y^2 + 3xy^2$

$$H_g = \begin{bmatrix} \frac{\partial^2 g}{\partial x^2} & \frac{\partial^2 g}{\partial x \partial y} \\ \frac{\partial^2 g}{\partial y \partial x} & \frac{\partial^2 g}{\partial y^2} \end{bmatrix} = \begin{bmatrix} 2x & 6y \\ 6y & 4y + 6x \end{bmatrix}$$

$\frac{\partial}{\partial x}(\frac{\partial}{\partial x}(x^2 + 2y^2 + 3xy^2)) = \frac{\partial}{\partial x}(2x + 3y^2) = 2$

$\frac{\partial}{\partial x}(\frac{\partial}{\partial y}(x^2 + 2y^2 + 3xy^2)) = \frac{\partial}{\partial x}(4y + 6xy) = 6y$

$\frac{\partial}{\partial y}(\frac{\partial}{\partial x}(x^2 + 2y^2 + 3xy^2)) = \frac{\partial}{\partial y}(2x + 3y^2) = 6y$

$\frac{\partial}{\partial y}(\frac{\partial}{\partial y}(x^2 + 2y^2 + 3xy^2)) = \frac{\partial}{\partial y}(4y + 6xy) = 4 + 6x$

When $f: \mathbb{R}^d \rightarrow \mathbb{R}$,

$$f(w) \approx f(x_0) + \nabla f(x_0) \cdot (x - x_0) + \frac{1}{2} \underbrace{(x - x_0)^T}_{1 \times d} \underbrace{H(x_0)}_{d \times d} \underbrace{(x - x_0)}_{d \times 1}$$

★ What if we wanted Taylor approximation of degree 3?



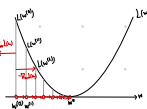
Taylor Series of $L(w)$ around $w^{(t)}$:

$$L(w) \approx L(w^{(t)}) + (w - w^{(t)})^T g^{(t)} + \frac{1}{2} (w - w^{(t)})^T H^{(t)} (w - w^{(t)})$$

① Substitute gradient descent update, $w^{(t+1)} = w^{(t)} - \epsilon g^{(t)}$, into the Taylor series: $L(w) \approx L(w^{(t+1)})$

Why?: To understand $L(w^{(t+1)}) - L(w^{(t)})$.

↳ We want this to be negative.
: We want $L(w^{(t+1)}) \leq L(w^{(t)})$



$$g^{(t)} = \nabla_w L(w^{(t)})$$

$$L(w^{(t+1)}) \approx L(w^{(t)}) + (w^{(t+1)} - w^{(t)})^T g^{(t)} + \frac{1}{2} (w^{(t+1)} - w^{(t)})^T H^{(t)} (w^{(t+1)} - w^{(t)})$$

Since $w^{(t+1)} = w^{(t)} - \epsilon g^{(t)}$, then $w^{(t+1)} - w^{(t)} = -\epsilon g^{(t)}$

$$\begin{aligned} \text{Then, } L(w^{(t+1)}) &\approx L(w^{(t)}) - \epsilon (g^{(t)})^T g^{(t)} - \frac{1}{2} \epsilon^2 (g^{(t)})^T H^{(t)} \epsilon g^{(t)} \\ &\approx L(w^{(t)}) - \epsilon (g^{(t)})^T g^{(t)} + \frac{1}{2} \epsilon^2 (g^{(t)})^T H^{(t)} g^{(t)} \end{aligned}$$

$$\text{And } L(w^{(t+1)}) = L(w^{(t)} - \epsilon g^{(t)})$$

• We want $L(w^{(t+1)}) \leq L(w^{(t)})$ for GD to converge

$$\text{Since } L(w^{(t+1)}) \approx L(w^{(t)}) - \epsilon (g^{(t)})^T g^{(t)} + \frac{1}{2} \epsilon^2 (g^{(t)})^T H^{(t)} g^{(t)}$$

$$L(w^{(t+1)}) - L(w^{(t)}) \approx -\epsilon (g^{(t)})^T g^{(t)} + \frac{1}{2} \epsilon^2 (g^{(t)})^T H^{(t)} g^{(t)}$$

Hence, since we want $L(w^{(t+1)}) \leq L(w^{(t)})$, which is equivalent to $L(w^{(t+1)}) - L(w^{(t)}) \leq 0$,

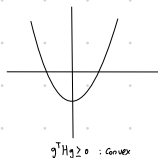
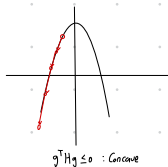
$$\text{we need to find } -\epsilon g^T g + \frac{\epsilon^2}{2} g^T H g \leq 0$$

① We know $g^T g \geq 0$ (always positive), hence $-\epsilon g^T g$ is always negative

② Then, what if $\frac{\epsilon^2}{2} g^T H g$ is negative?

: then $-\epsilon g^T g + \frac{\epsilon^2}{2} g^T H g \leq 0$ is satisfied

Hence, if $g^T H g \leq 0$ (as long as $\epsilon > 0$), then regardless of ϵ , $L(w^{(t+1)}) \leq L(w^{(t)})$



③ But, what if $g^T H g > 0$?

$$\text{: Since we want } -\epsilon g^T g + \frac{\epsilon^2}{2} g^T H g \leq 0$$

$$= \frac{\epsilon^2}{2} g^T H g \leq \epsilon g^T g$$

$$= \frac{\epsilon g^T g}{2} \leq g^T g$$

$$= \epsilon \leq \frac{2g^T g}{g^T H g}$$

Hence, when $\epsilon \leq \frac{2g^T g}{g^T H g}$, $L(w^{(t+1)}) \leq L(w^{(t)})$

In Summary:

- When $(g^{(t)})^T H^{(t)} g^{(t)} \leq 0$, we have no issue (i.e., if $L(w)$ is concave)
- But, when $(g^{(t)})^T H^{(t)} g^{(t)} > 0$, to satisfy $L(w^{(t+1)}) < L(w^{(t)})$, we need to pick ϵ carefully.

Convergence Rate: How fast is it going to converge?

- Assume that H has a spectral radius smaller than L (Lipschitz constant):

$$\text{s.t. } v^T H(w)v \leq L \|v\|^2$$

\uparrow
 Lipschitz constant

In other words, $\lambda_{\max}(H(w)) \leq L$ $\forall w$.

\hookrightarrow Spectral radius of $H(w)$

$$\text{Then, } \frac{2g^T g}{L \cdot g^T g} \leq \frac{2g^T g}{g^T H g}, \text{ which implies } \epsilon < \frac{2}{L} \leq \frac{2g^T g}{g^T H g}$$

Hence, if we set $\epsilon < \frac{2}{L}$, we will guarantee convergence

I think it will be related to the concept of smoothness. But not sure yet.

Rewriting the Taylor expansion we have:

$$L(w^{(t+1)}) = L(w^{(t)}) - \epsilon (g^{(t)})^T g^{(t)} + \frac{1}{2} \epsilon^2 (g^{(t)})^T H^{(t)} g^{(t)} \leq L(w^{(t)}) - \epsilon \|g^{(t)}\|^2 + \frac{1}{2} \epsilon^2 \|g^{(t)}\|^2 \quad \text{since } g^T H g \leq L \|g\|^2$$

$$= L(w^{(t)}) + \left(\frac{1}{2} \epsilon^2 - \epsilon\right) \|g^{(t)}\|^2$$

Since $\epsilon < \frac{2}{L}$, let's assume $\epsilon = \frac{1}{L}$

Then,

$$L(w^{(t+1)}) \leq L(w^{(t)}) + \left(\frac{1}{2} \cdot \frac{1}{L^2} - \frac{1}{L}\right) \|g^{(t)}\|^2 = L(w^{(t)}) - \frac{\|g^{(t)}\|^2}{2L}$$

So, we have

$$L(w^{(t+1)}) \leq L(w^{(t)}) - \frac{\|g^{(t)}\|^2}{2L}$$

\uparrow
Lipschitz

$$\Rightarrow \frac{\|g^{(t)}\|^2}{2L} \leq L(w^{(t)}) - L(w^{(t+1)})$$

$$\Rightarrow \|g^{(t)}\|^2 \leq 2L(L(w^{(t)}) - L(w^{(t+1)}))$$

Let's now take sum over + steps of gradient descent

$$\sum_{k=0}^{t-1} \|g^{(k)}\|^2 \leq 2L \sum_{k=0}^{t-1} (L(w^{(k)}) - L(w^{(k+1)}))$$

Lower bound left-hand-side by:

$$t g_m \leq \sum_{k=1}^t \|g^{(k)}\|^2 \quad \text{where} \quad g_m = \min_{j \in \{1, \dots, t\}} \|g^{(j)}\|^2$$

The right-hand-side is equal to

$$\begin{aligned} \sum_{k=0}^{t-1} (\mathcal{L}(w^{(k)}) - \mathcal{L}(w^{(k+1)})) &= \mathcal{L}(w^{(0)}) - \mathcal{L}(w^{(1)}) + \mathcal{L}(w^{(1)}) - \mathcal{L}(w^{(2)}) + \dots + \mathcal{L}(w^{(t-1)}) - \mathcal{L}(w^{(t)}) \\ &= \mathcal{L}(w^{(0)}) - \mathcal{L}(w^{(t)}) \leq \mathcal{L}(w^{(0)}) - \mathcal{L}(w^*) \quad \text{because} \quad \mathcal{L}(w^*) \leq \mathcal{L}(w^{(t)}) \end{aligned}$$

Therefore, we can write:

$$t g_m \leq \sum_{k=0}^{t-1} \|g^{(k)}\|^2 \leq 2\mathcal{L} \sum_{k=0}^{t-1} (\mathcal{L}(w^{(k)}) - \mathcal{L}(w^{(k+1)})) \leq 2\mathcal{L}(\mathcal{L}(w^{(0)}) - \mathcal{L}(w^*))$$

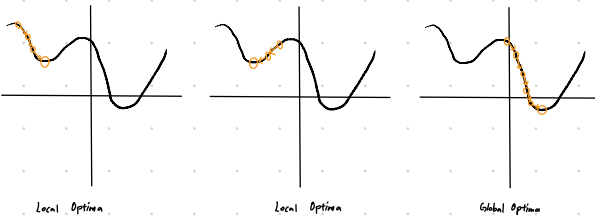
$$\Rightarrow g_m \leq \frac{2\mathcal{L}(\mathcal{L}(w^{(0)}) - \mathcal{L}(w^*))}{t}$$

Hence, $g_m = \nabla_w \mathcal{L}$ goes to zero at a rate $O(\frac{1}{t})$

How many iterations do we need to satisfy $\|\nabla_w \mathcal{L}\|^2 \leq \delta$? : $O(\frac{1}{\delta})$

Gradient Descent in Non-convex Problems.

• Sensitive to initialization and could get stuck in local optima



And Neural Network's Loss Function is Highly Nonconvex. Why does the optimization work?

Role of Overparameterization (Mikhail Belkin)

