

Gradient Descent!

Given $f: \mathbb{R}^d \rightarrow \mathbb{R}$

The gradient of f at w is

$$\nabla f(w) = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \vdots \end{bmatrix}$$

or... $\nabla f(w)$ is the unique vector s.t. $\forall u \in \mathbb{R}^d$

$$\begin{aligned} u^T \nabla f(w) &= \lim_{\alpha \rightarrow 0} \frac{f(w + \alpha u) - f(w)}{\alpha} \\ &= \left. \frac{\partial}{\partial \alpha} f(w + \alpha u) \right|_{\alpha=0} \end{aligned}$$

e.g. if $f(w) = w^T A w$ for matrix A . Then:

$$\begin{aligned} \frac{\partial}{\partial \alpha} f(w + \alpha u) &= \frac{\partial}{\partial \alpha} (w + \alpha u)^T A (w + \alpha u) \\ &= \frac{\partial}{\partial \alpha} w^T A w + \alpha \cdot u^T A w + \alpha \cdot w^T A u + \alpha^2 u^T A u \\ &= u^T A w + w^T A u + 2\alpha u^T A u \end{aligned}$$

$$\textcircled{a} \alpha=0 \mid = u^T A w + w^T A u = u^T \underbrace{(A w + A^T w)}_{\nabla f(w)}$$

Hessian: $\nabla^2 f \approx H$

$$\nabla^2 f(w) = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 \partial w_2} & \dots & - \\ \frac{\partial^2 f}{\partial w_2 \partial w_1} & \dots & \dots & - \\ \dots & \dots & \dots & - \\ \dots & \dots & \dots & - \end{bmatrix} \in \mathbb{R}^{d \times d}$$

or: the unique symmetric matrix s.t.

$$u^T \nabla^2 f(w) u = \left. \frac{\partial^2}{\partial \alpha^2} f(w + \alpha u) \right|_{\alpha=0}$$

if $f: \mathbb{R}^d \rightarrow \mathbb{R}$, then $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ & $\nabla^2 f: \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$.

A general recipe: initialize $w_0 \in \mathbb{R}^d$

loop until converged:

- select some "step" $s \in \mathbb{R}^d$
- $w \leftarrow w + s$
- "converged" if $\|s\| \leq \delta$ (δ tolerance parameter)

Conclusion

$$l(w+s) \approx l(w) + s^T \nabla l(w)$$

$$l(w+s) \approx l(w) + s^T \nabla l(w) + \frac{1}{2} s^T \nabla^2 l(w) s$$

Gradient descent: idea look at $l(w+s)$ as a first order approximation

Direction of steepest descent is $-\nabla l(w)$

$$w \leftarrow w - \alpha \nabla l(w) \quad (\text{gradient descent})$$

step size
learning rate

$$l(w - \alpha \nabla l(w)) = l(w) - \alpha \nabla l(w)^T \nabla l(w) + \frac{1}{2} (\alpha \nabla l(w))^T \nabla^2 l(w) (\alpha \nabla l(w)) + \frac{1}{2} (\alpha \nabla l(w))^T \underbrace{\nabla^2 l(w + \beta \alpha \nabla l(w))}_{\text{bounded}} (\alpha \nabla l(w))$$

$$= l(w) - \alpha \|\nabla l(w)\|^2 + \underbrace{\frac{1}{2} \alpha^2}_{\text{bounded}} O(\alpha^2 \|\nabla l(w)\|^2)$$

$$= l(w) - \underbrace{(\alpha - O(\alpha^2))}_{\geq \frac{\alpha}{2}} \|\nabla l(w)\|^2$$

$$l(w - \alpha \nabla l(w)) \leq l(w) - \frac{\alpha}{2} \|\nabla l(w)\|^2 \quad \alpha - C\alpha^2$$

pick $\alpha \leq \frac{1}{2C}$

if we didn't converge, $\|s\| \geq \delta \Rightarrow \|\alpha \nabla l(w)\| \geq \delta$

$$l(w - \alpha \nabla l(w)) \leq l(w) - \frac{\delta^2}{2\alpha}$$

✓ GD converges

Newton's Method

$$l(\bar{w}) \approx l(w) + s^T \nabla l(w) + \frac{1}{2} s^T H(w) s$$

$$0 = \nabla l(w) + H(w)s \Rightarrow s = - (H(w))^{-1} \nabla l(w)$$

Newton's s

$$w \leftarrow w - (H(w))^{-1} \nabla l(w) \quad (\text{Newton's method})$$

e.g. $w_{t+1} = w_t - (H(w_t))^{-1} \nabla l(w_t)$