④ ① SGD:

$$X\hat{w} = t \implies X\hat{w} - t = 0 \quad \cdots ①$$

In SGD, all $x_i$ is contained in the span o X. And

the SGD update do steps don't ever leave the span of

X. Because, $\dfrac{d}{d\hat{w}_p}\left(x_i \hat{w}_p - t_i\right)^2 = 0$ will give update

$\hat{w}_p$ as so some combination of $x_i$, and $t_i$.

Thereby, we can assume the SGD solution is spanned

by X : $\hat{w} = X^T S$ , where S is a arbitrary matrix.

From ①, $X\hat{w} - t = 0$

$$\implies X X^T S - t = 0$$

$$\implies S = (X X^T)^{-1} t$$

So, $\hat{w} = X^T (X X^T)^{-1} t$

$$= w^*$$

[ Showed ]

④ ② Mini-batch SGD:

Yes, mini-batch SGD also obtains minimum norm solution on convergence.

Because the batch $B$ is taken from the rows of $X$. So, the solution $\hat{\omega}$ is spanned by the rows of $X$.

$$\hat{\omega} = B S = X\tilde{S}$$

$\downarrow$
Batch

So, $X\hat{\omega} - t = XX^T\tilde{S} - t = 0$

$$\Rightarrow \tilde{S} = (XX^T)^{-1} t$$

$$\therefore \hat{\omega} = X^T(XX^T)^{-1} t = \omega^*$$

④ ③ Adaptive Methods: Adagrad

$x_1 = [2, 1]$   $\omega_0 = [0, 0]$   $\ell = [2]$

Using minimum norm solution with GD,

we got   $\omega^* = \begin{bmatrix} 0.8 \\ 0.4 \end{bmatrix}$ and

$$\nabla_{\omega^*} \mathcal{L}(\omega) = -2x_1 \ell_1$$

Using Adagrad,

$$\omega_0 = [0, 0]$$

$$\omega_1 = \omega_0 - \frac{n}{\sqrt{G_{01}} + \epsilon} \nabla_{\hat{\omega}_0} \mathcal{L}(\omega)$$

$$G_1 = 0 + \left(\nabla_{\hat{\omega}_0} \mathcal{L}(\omega)\right)^2$$

let, assume, $\nabla_{\hat{\omega}_0} \mathcal{L}(\omega) = -2x_1 \ell_1$ [similar to the GD]

then, $\omega_1 = \omega_0 - \dfrac{n}{(-2x_1 \ell_1) + \epsilon} \cdot (-2x_1 \ell_1)$

as $\epsilon$ is small, $\omega_1$ looses $x_1$ term almost, ~~that means~~ the direct because numerator and

denominator both contains $(-2x_1 \ell_1)$, So,

$\omega_1$ has a little impact from $x_1$, which indicates

the direction of the gradient is no longer

along $x_1$ as much as the $0$ $\omega^*$ case.

(GD with minimum norm sol.)

Thereby, AdaGrad doesn't always obtain the

minimum norm solution.

This

~~Same~~ results. holds

true for other adaptive

~~models~~ methods (RMSProp,

Adam) in general.

gradient of the AdaGrad

gradient of the minimum norm sol.

Because the scaling part in the weight update

may divert solution gradient from

the span of X and it may get outside of

the span of X sometimes.