**Q(8)**

$$\frac{\partial L}{\partial w_i} = \delta \cdot x_i$$

where $x_i$ is input recvd by from a prev. layer that utilized activation fn.

$$\Rightarrow x_i > 0 \Rightarrow \frac{\partial L}{\partial w_i} \text{ depends on sign of } \delta.$$

$\Rightarrow$ Dirn of "step" $= \left(\dfrac{\partial L}{\partial w_1}, \dfrac{\partial L}{\partial w_2}, \cdots \rightarrow \dfrac{\partial L}{\partial w_n}\right)$

$$= (\delta \cdot x_1, \delta \cdot x_2, \longrightarrow \delta \cdot x_n)$$

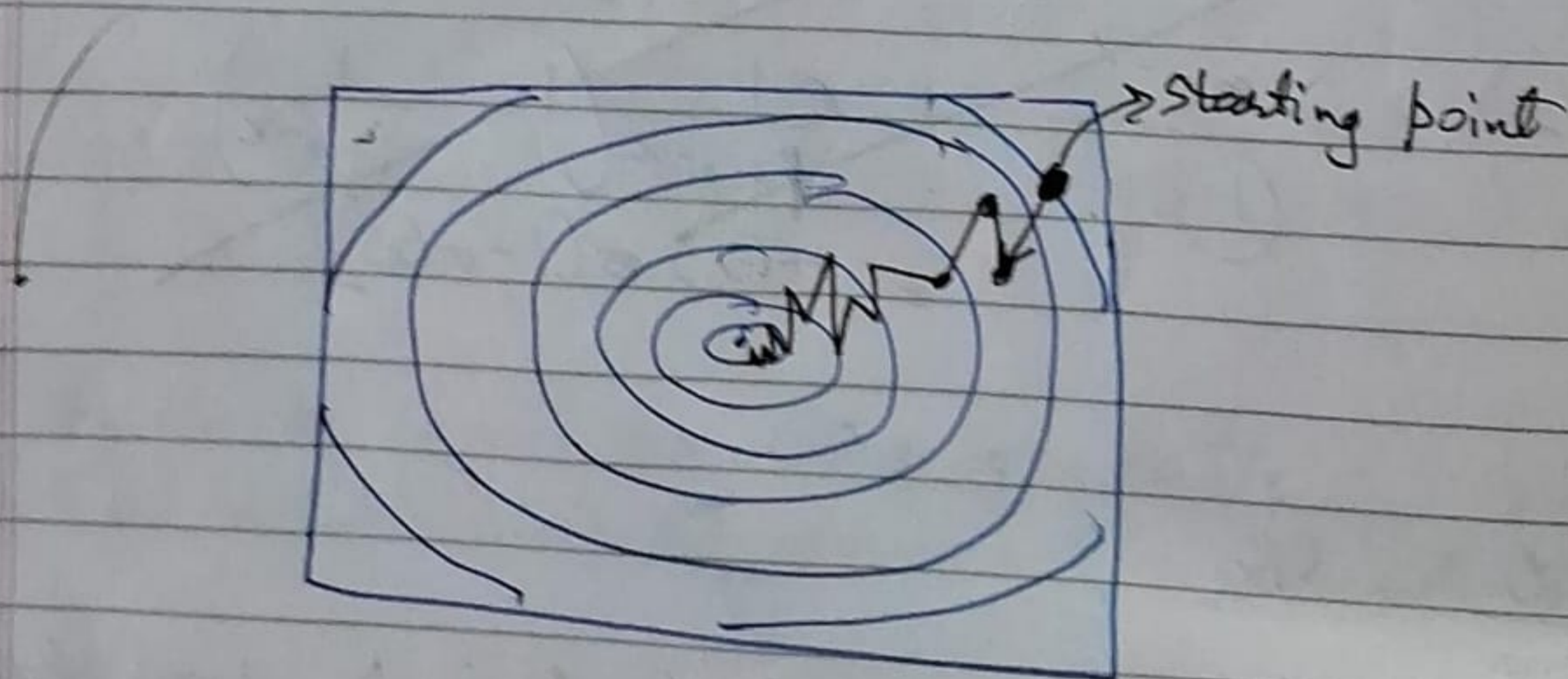$$= \delta (x_1, x_2, \longrightarrow x_n)$$

$$\Rightarrow \quad w_i \leftarrow w_i - \alpha \frac{\partial L}{\partial w_i}$$

$$\Rightarrow \quad w_i \leftarrow w_i - \alpha \delta ~~ x_i$$

This means all $w_i$'s ~~can~~ either increase or decrease simultaneously. The dirn of ~~motion~~ "step" is same and depends only on $\delta$.

## Problem with the above

Consider countour plot of ~~$L(w)$~~ $L(w, b)$:



Suppose ~~too~~ in two iterations we get $\delta$ with opposite signs, ~~which is~~ This is possible because ~~the~~ all inputs increasing/decreasing simultaneously does not guarantee ~~actual~~ dirn of steepest descent

In such a case, we first "overshoot" ($\uparrow$increase $w_i$'s in one dirn) and then retract ~~~~ (move in another dirn). This causes "zigzag" ~~~~ motion and chaotic convergence.