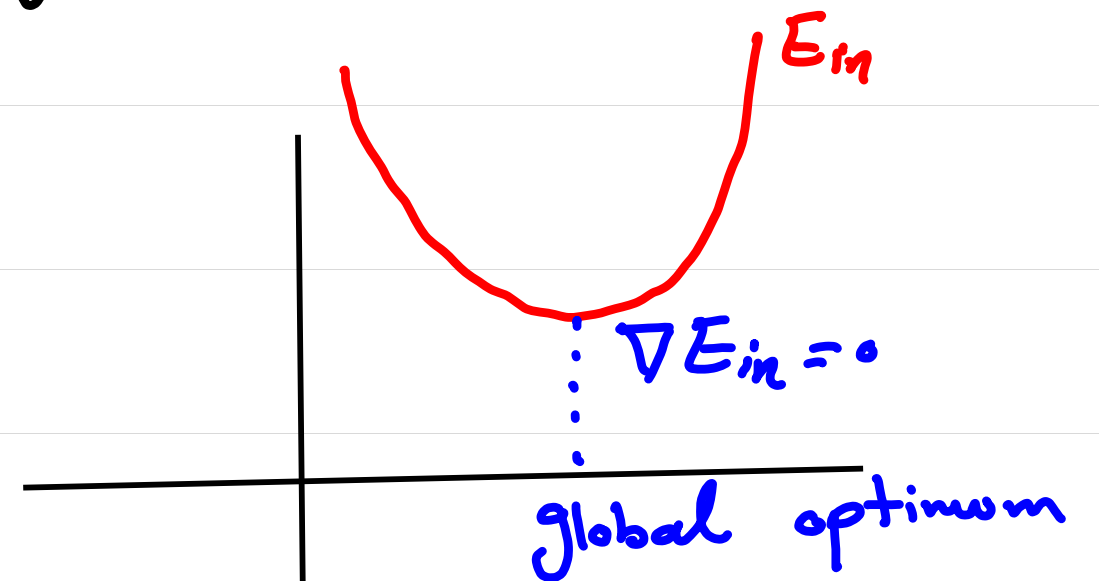


Week 04 - Part 03

So far: For Convex functions, GD Converges to the optimal point.

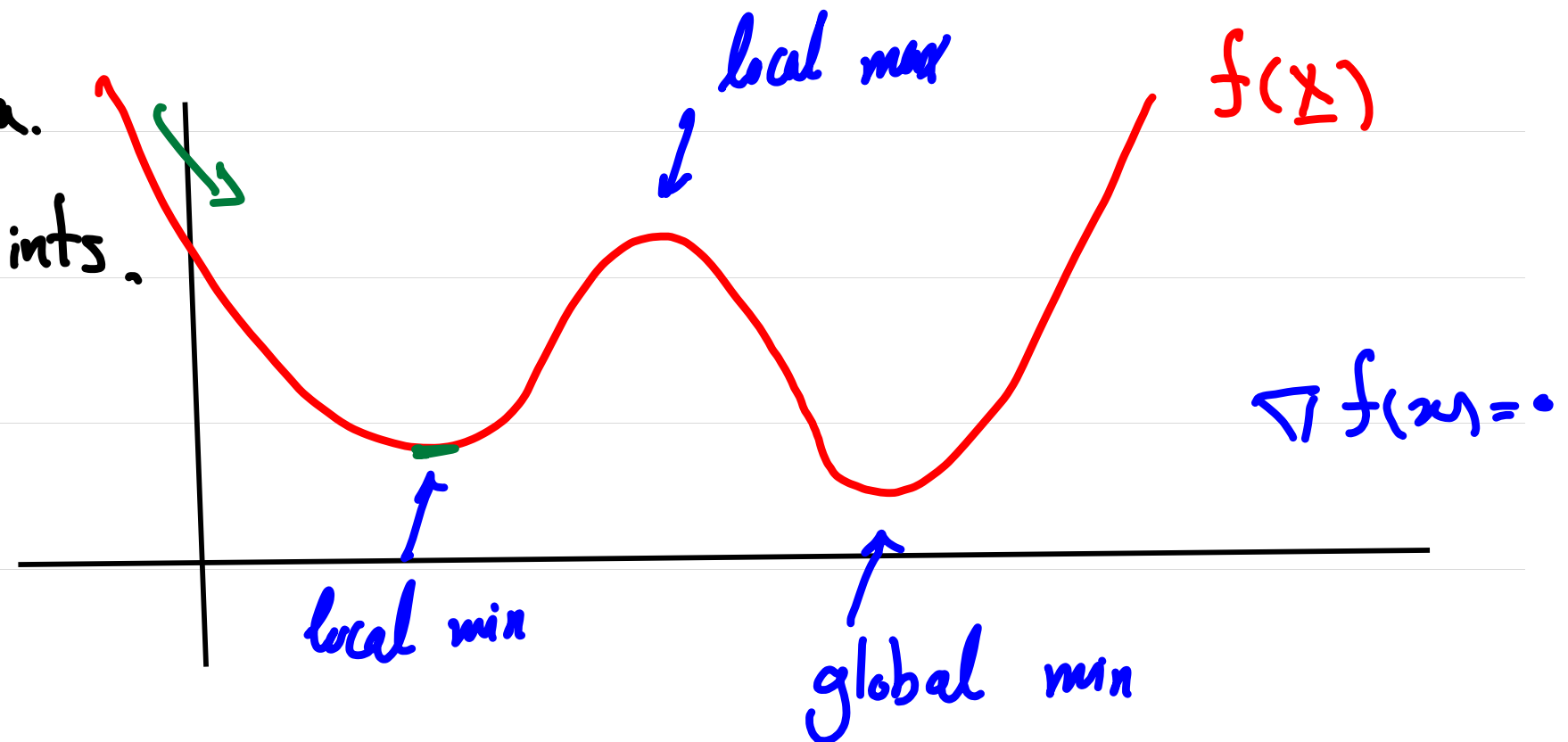
☞ E_{in} for linear / Logistic Regression is Convex.



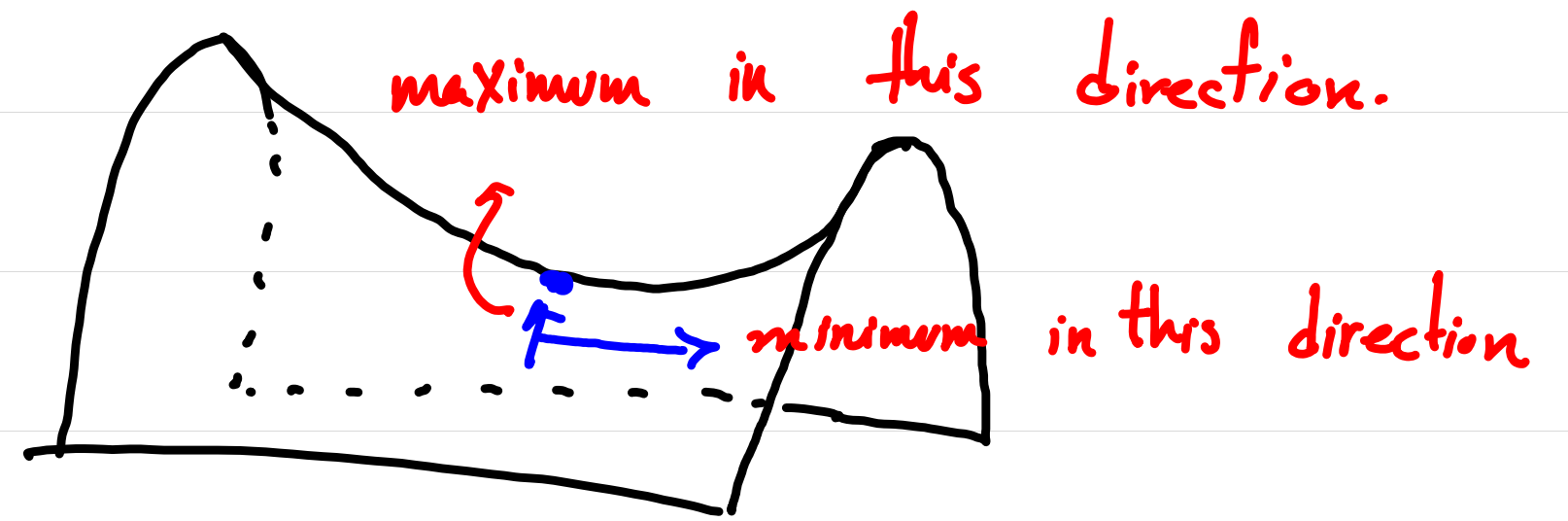
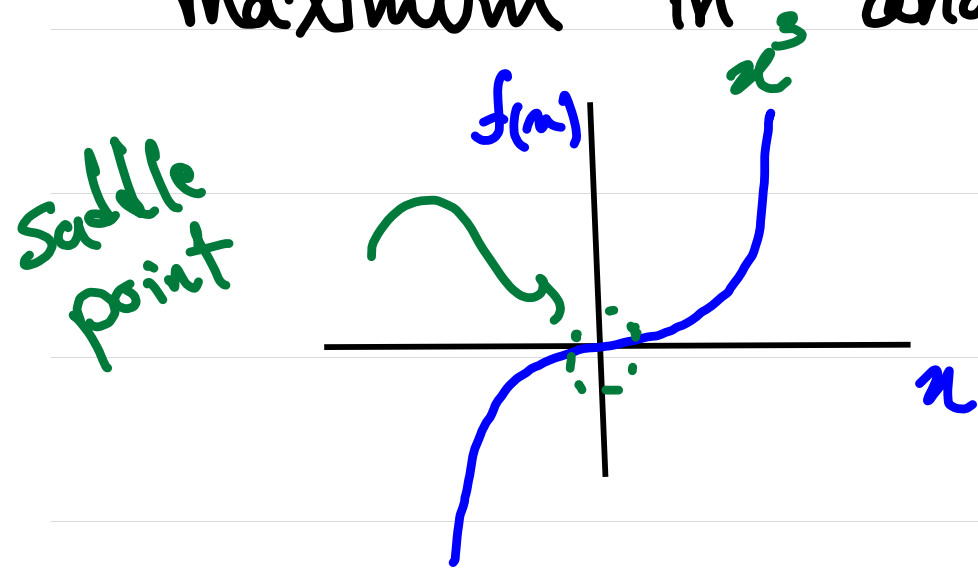
Today: Non-convex functions

☞ we often have multiple extrema.

☞ we may also have saddle points.



■ A saddle point is minimum in some direction, and maximum in another direction.



■ GD would have very slow progress near saddle points

- The preset # of iterations will run out

■ In high-dimensional space, it's highly likely to have saddle points.

■ Most often, we are dealing with high-dimensional spaces.

- How to use GD in this high-dimensional spaces?

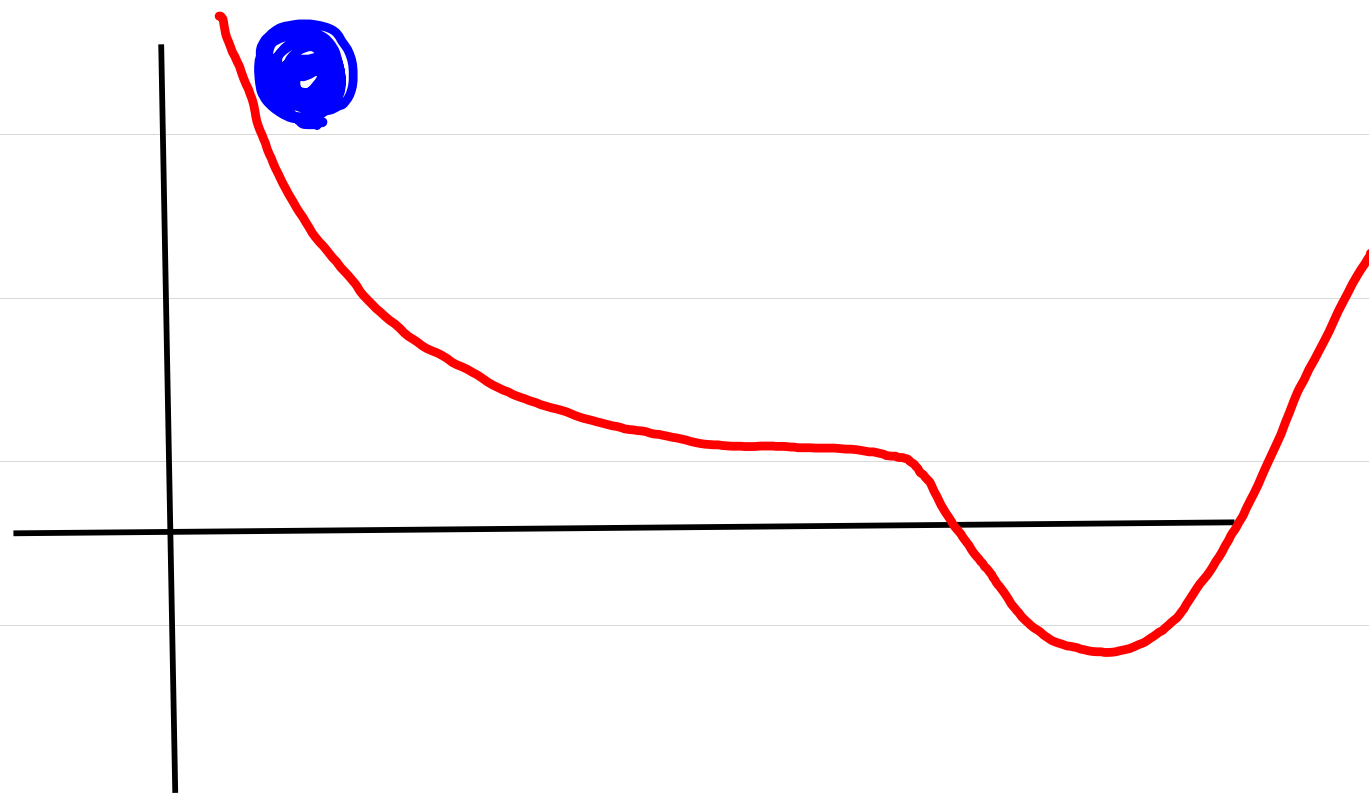
SGD with Momentum (Polyak, 1964)

Basic SGD: $\underline{g}_t = \nabla \ell_n(\underline{w}_t)$

$$\underline{v}_t = -\underline{g}_t$$

$$\underline{w}_{t+1} = \underline{w}_t + \epsilon_t \underline{v}_t$$

■ New idea: add a momentum (a push) so that SGD would not stop if $\nabla \ell_n(\underline{w}_t) \approx 0$.



□ When you let this heavy ball go, would it stop at the flat region? No

■ Why? Inertia/Momentum

SGD + Momentum:

this approach is
called "heavy ball"
momentum.

$$\underline{g}_t = \nabla e_n(\underline{w}_t)$$

$$\underline{v}_t = -\epsilon_t \underline{g}_t + \mu \underline{v}_{t-1}$$

$$\underline{w}_{t+1} = \underline{w}_t + \underline{v}_t$$

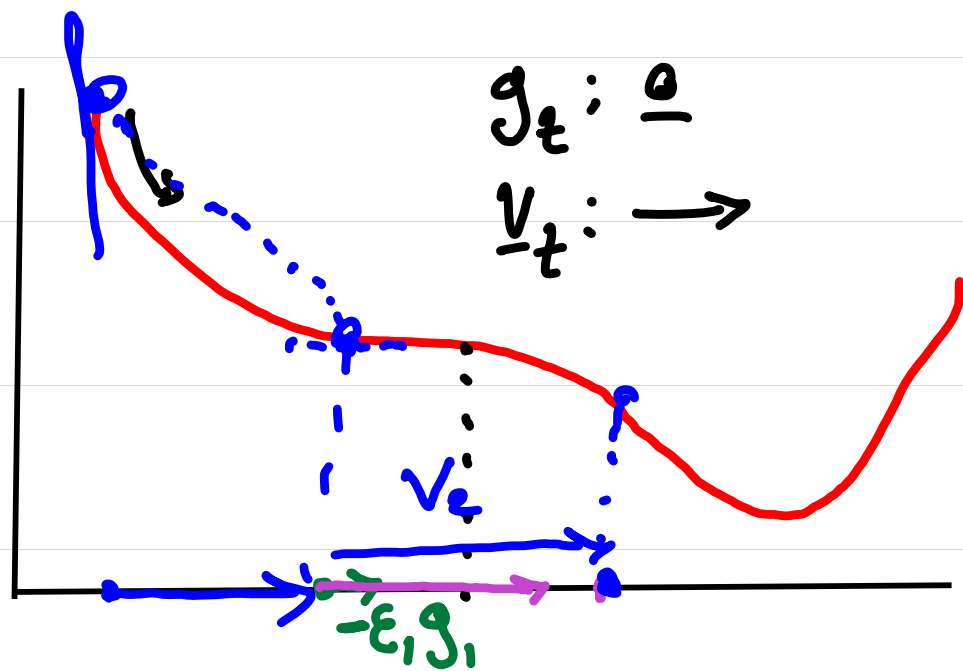
$$0 < \mu < 1$$

e.g. $\mu = 0.9$

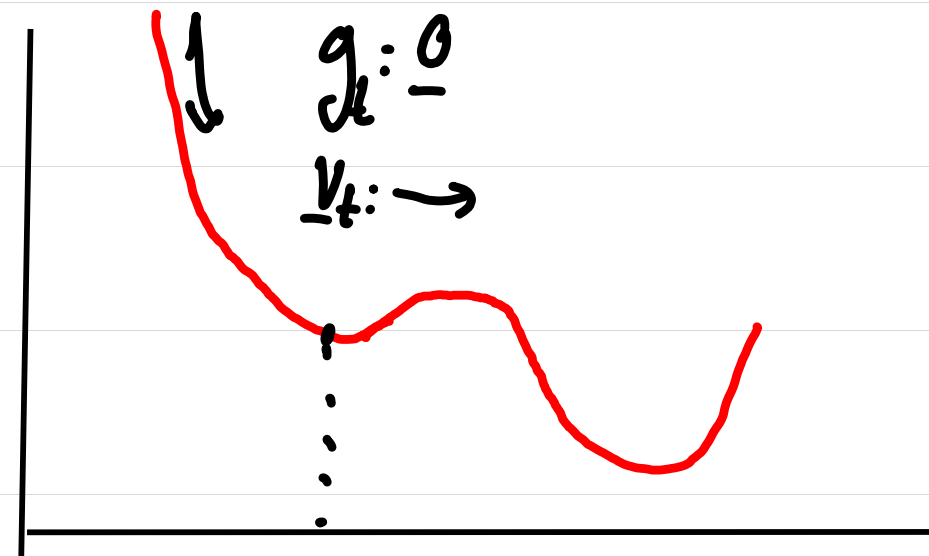
accumulation of your previous
movements.

How Does Momentum Help?

① Momentum helps SGD escape flat regions, saddle points, and shallow local optimum.

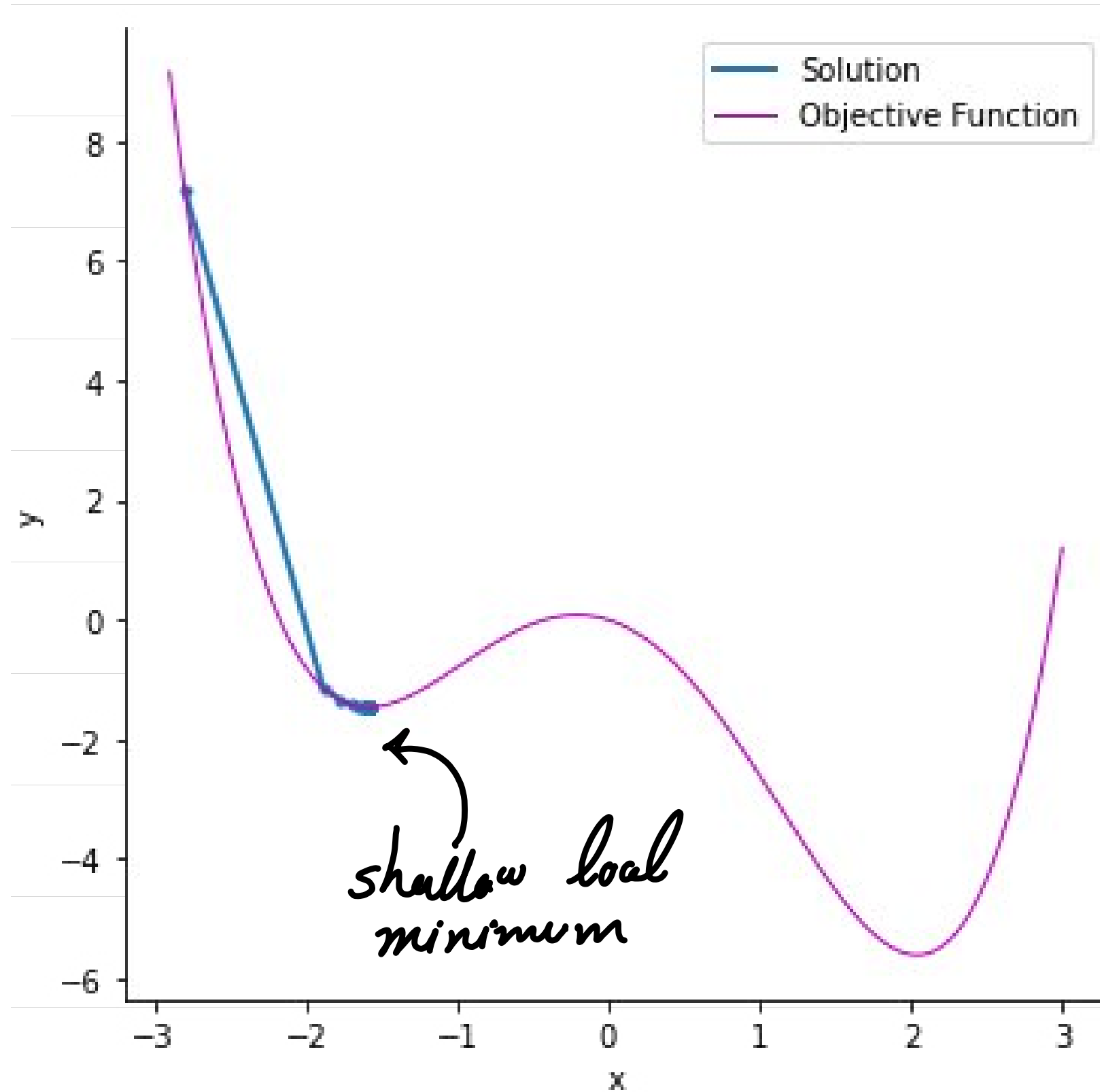


flat region/saddle point

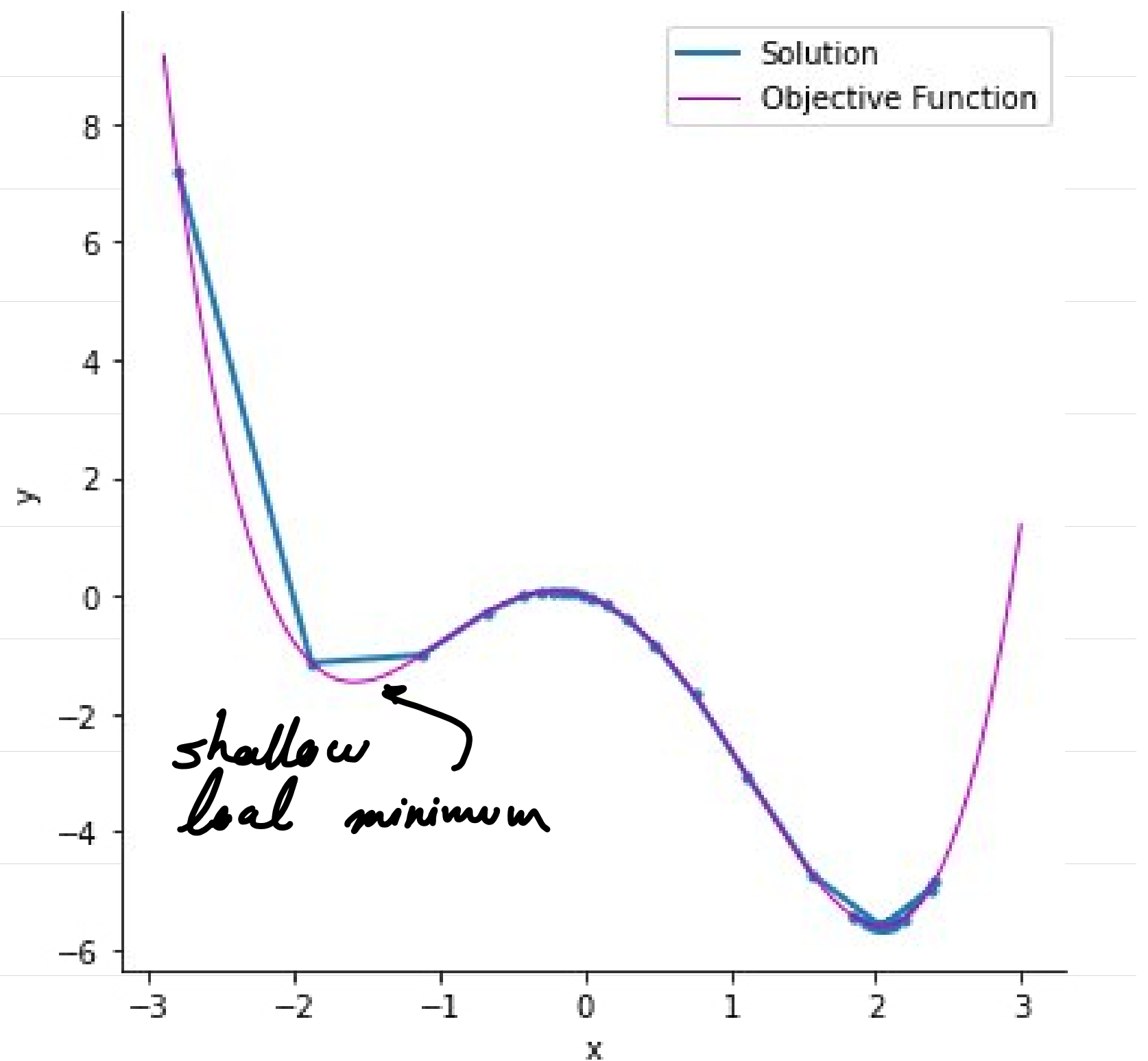


shallow local minimum

SGD



SGD + Momentum



How Can Momentum Help?

2-Momentum can lead to faster convergence, even for convex functions.

■ Consider the stretched ellipsoid $f(x_1, x_2) = 0.1x_1^2 + 2x_2^2$.

□ its minimum is at $(0,0)$

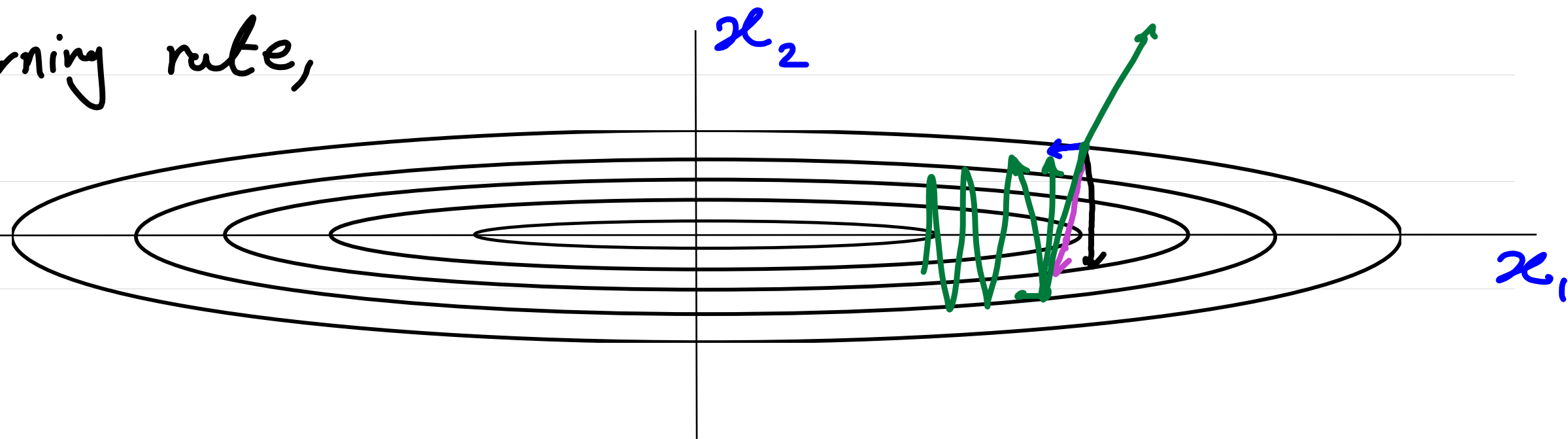
□ This function is very flat in the direction of x_1 ,

□ The gradient in x_2 direction is much higher than x_1 .

□ with a small learning rate, SGD won't diverge in x_2 direction but is very slow in x_1 direction.

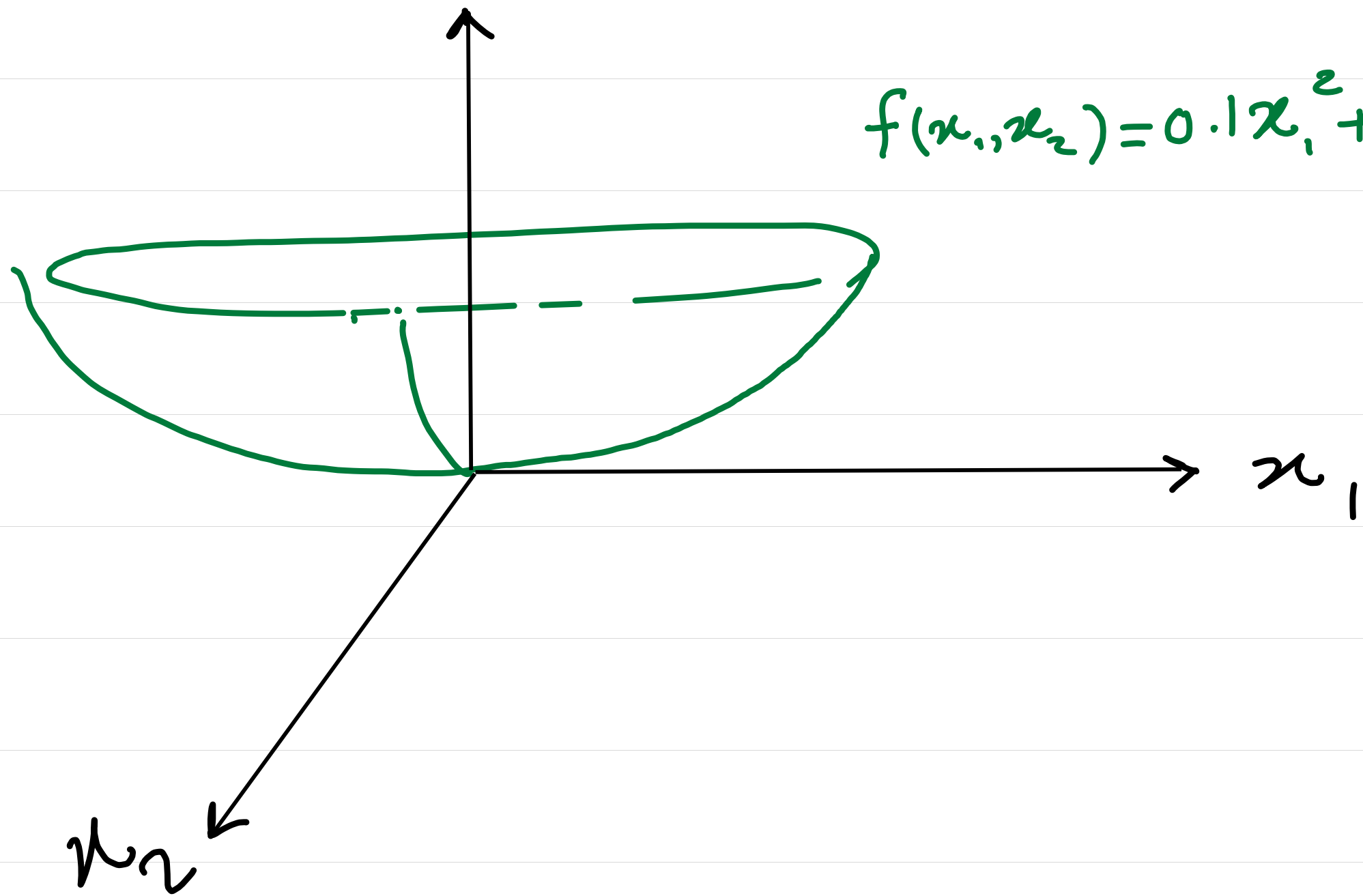
□ with large learning rate,

SGD progresses
more rapidly in x_1 ,
but diverges in x_2 .

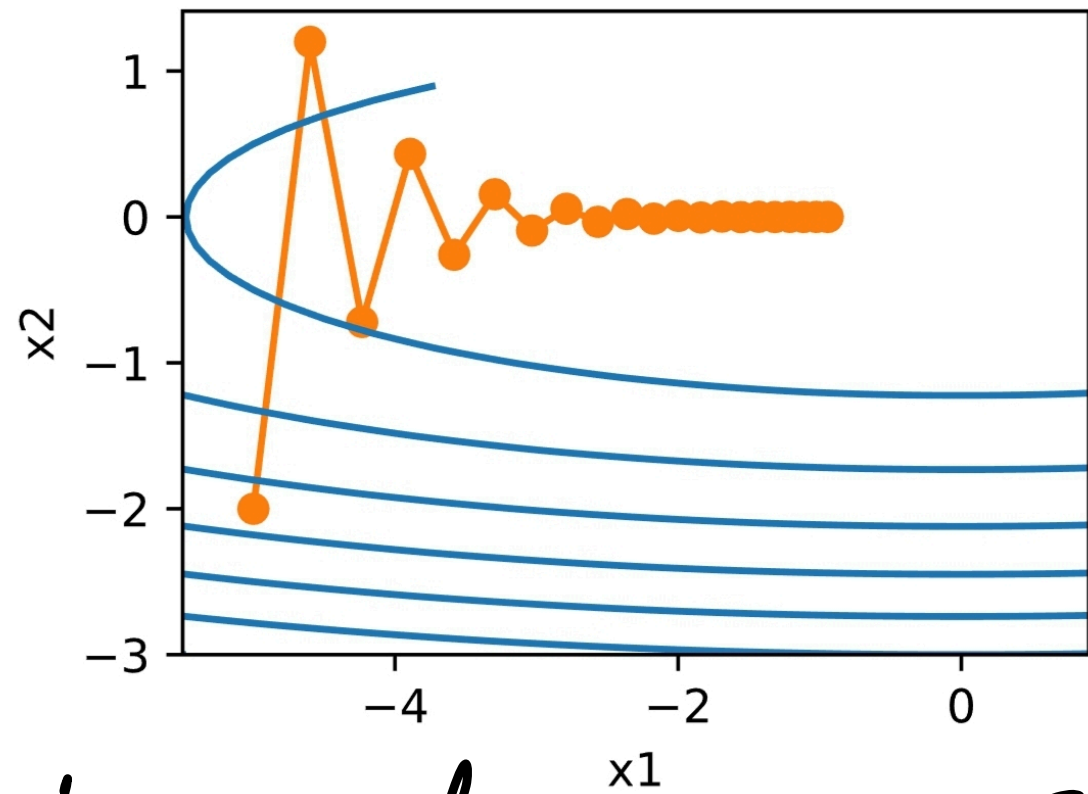


$$f(x_1, x_2)$$

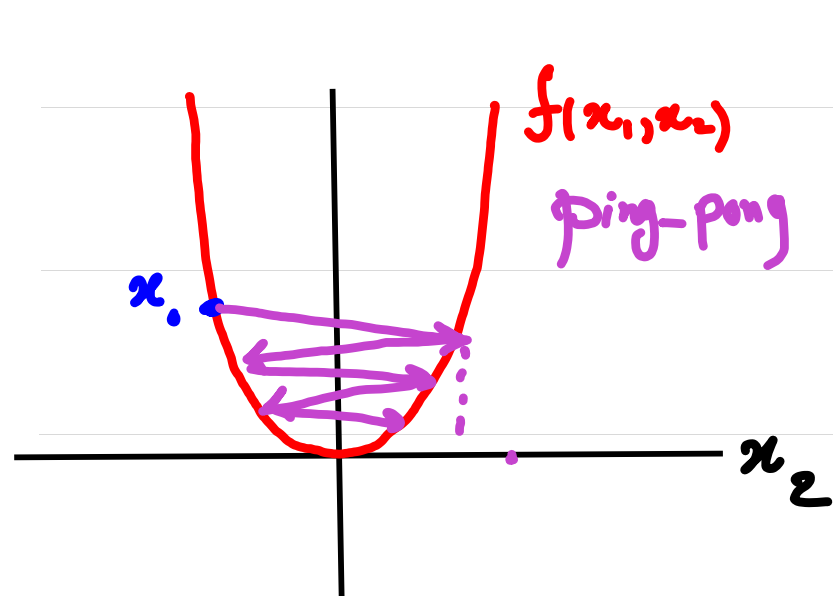
$$f(x_1, x_2) = 0.1x_1^2 + 2x_2^2$$



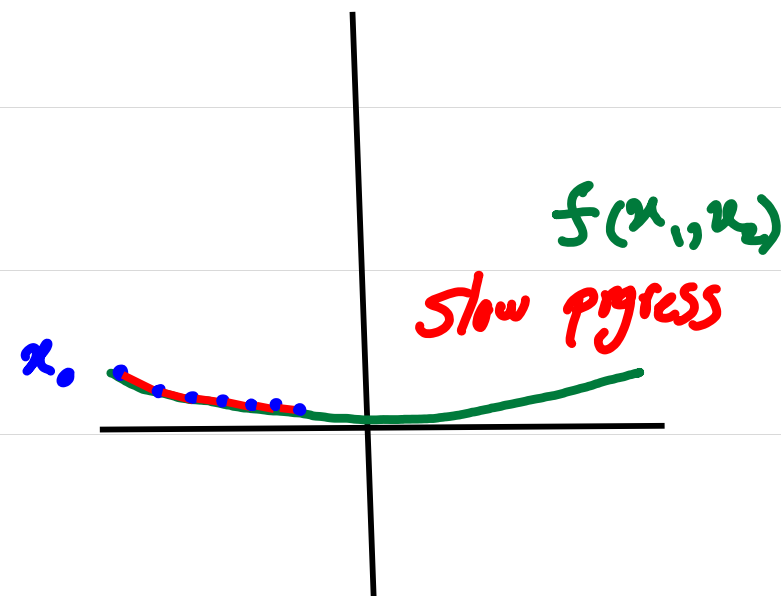
SGD with Learning rate = 0.4



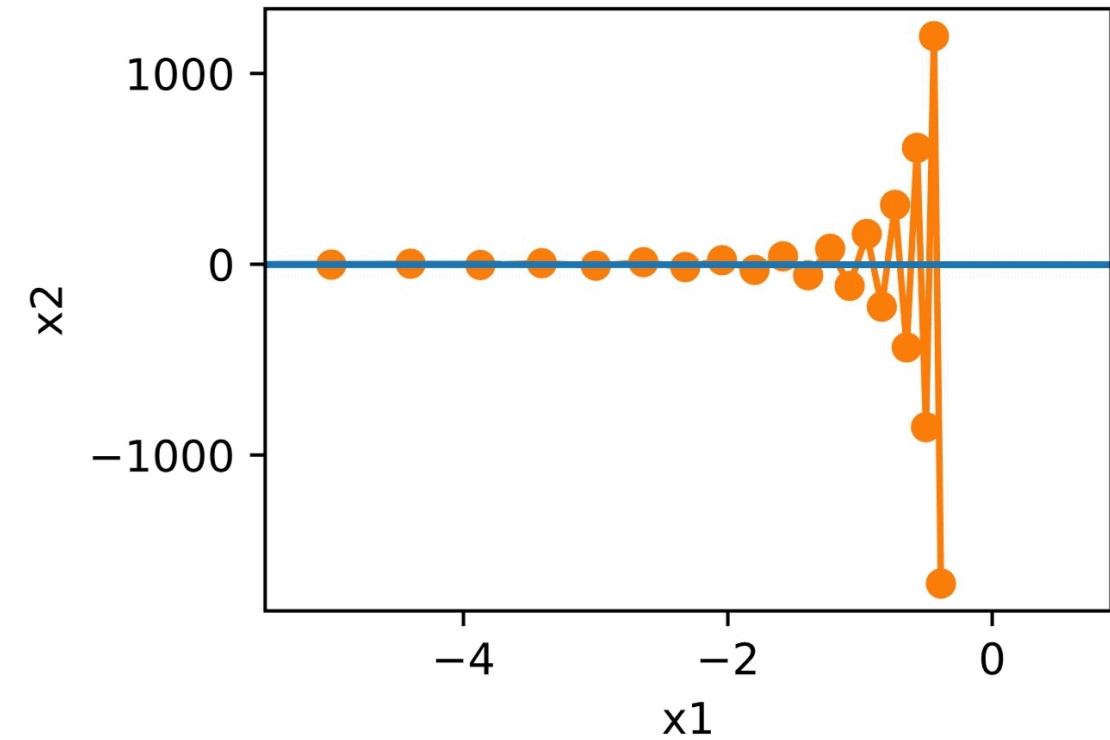
■ Slicing along x_2 direction



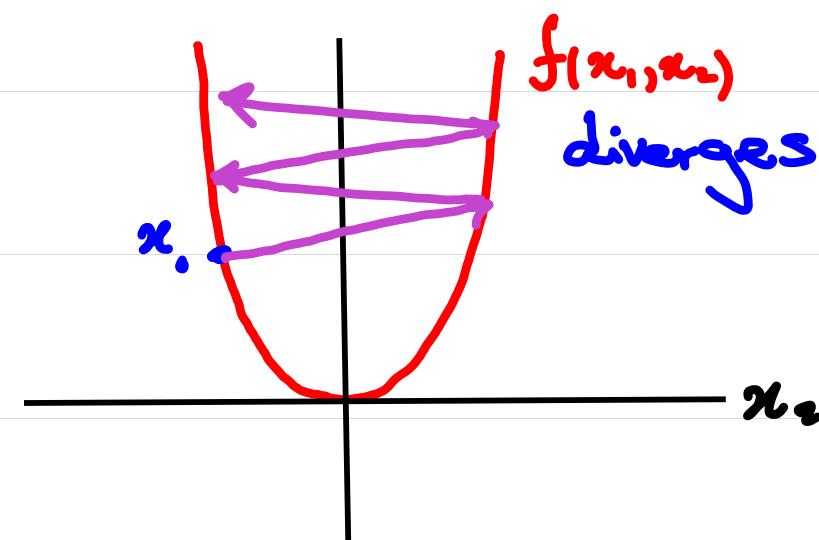
■ Slicing along x_1 direction



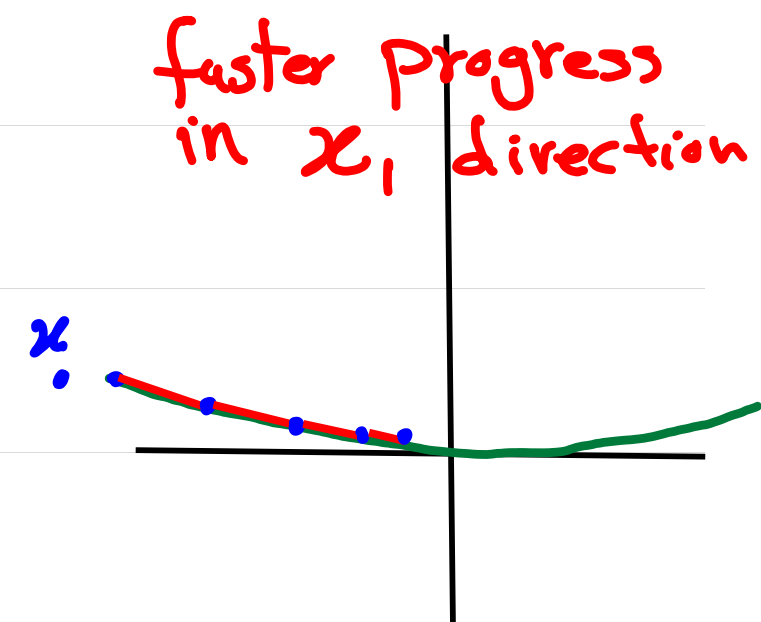
SGD with Learning rate = 0.6



■ Slicing along x_2 direction

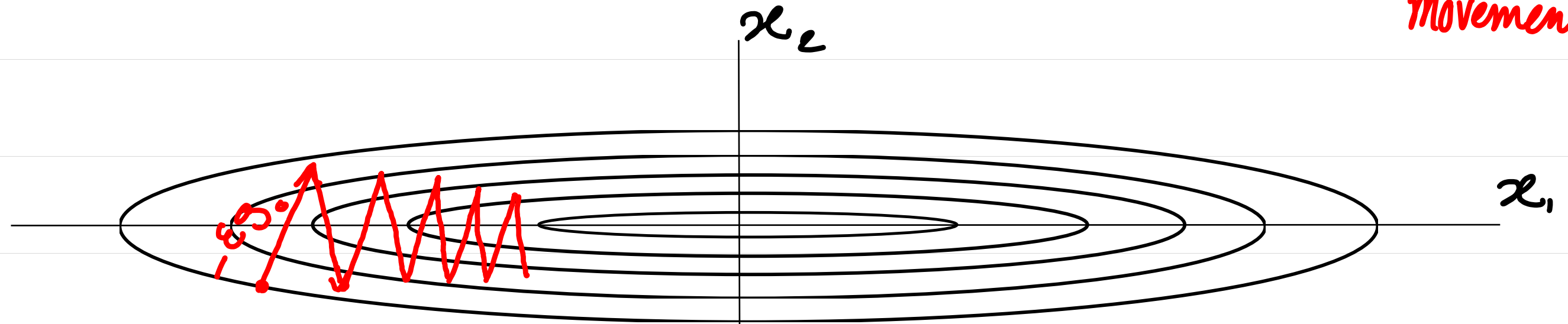


■ Slicing along x_1 direction



SGD+Momentum: $V_t = -\epsilon_t g_t + \mu V_{t-1}$ accumulation of your previous movements.

→ $-\epsilon_t g_t$

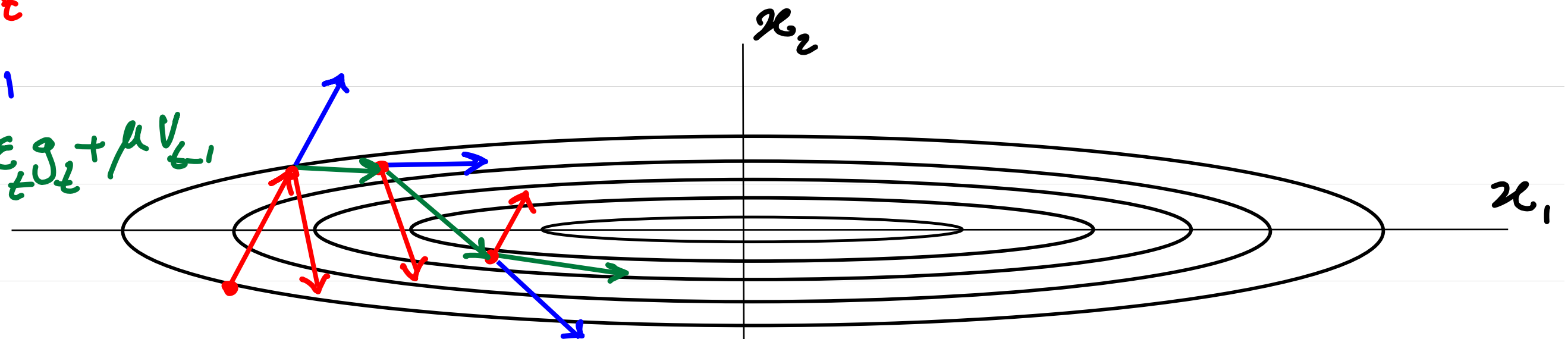


SGD without Momentum

→ $-\epsilon_t g_t$

→ μV_{t-1}

→ $V_t = -\epsilon_t g_t + \mu V_{t-1}$



SGD with heavy ball momentum

- Heavy-ball momentum method works well in practice.
- But we don't have any theoretical proof for it.
- Nesterov, in 1983, modified the momentum and could prove nice theoretical guarantees.

Nestrov Momentum

■ You update your location with your velocity first, and then take the gradient

intermediate point.

$$\underline{V}_t = -\epsilon_t \nabla \ell_n(\underline{W}_t + \mu \underline{V}_{t-1}) + \mu \underline{V}_{t-1}$$

$$\underline{W}_{t+1} = \underline{W}_t + \underline{V}_t$$

■ Provably better convergence for convex functions.

□ Full GD: $|f(\underline{w}_t) - f(\underline{w}^*)| = O(\frac{1}{t})$

□ With Nestrov: $|f(\underline{w}_t) - f(\underline{w}^*)| = O(\frac{1}{t^2})$