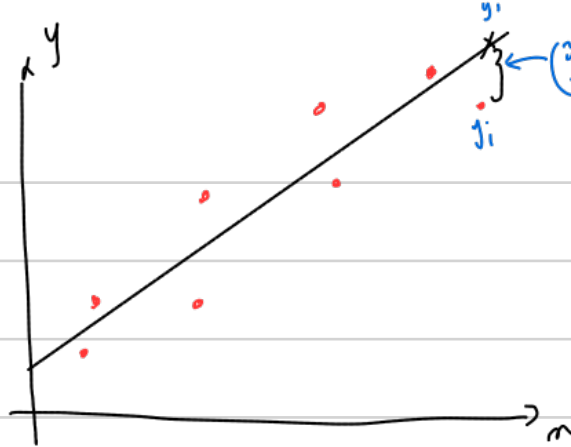


Regression



$$\text{mean square Error} \Rightarrow \text{loss} = w^2$$



Training loop

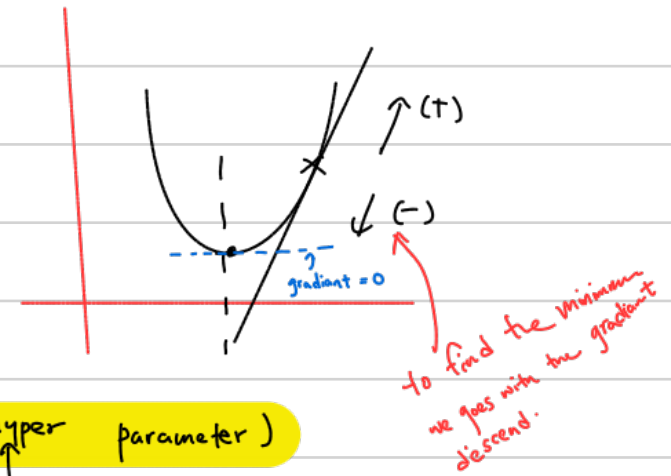
Tensors ← data containers
Autograds ← Numerical diffi...
Optimum ← Optimizer

Find the optimal

$$f(w) = w^2$$

$$f'(w) = 2w$$

moving downhill $\Rightarrow w_{\text{new}} = w - \eta f'(w)$
 $= w - \eta(2w)$

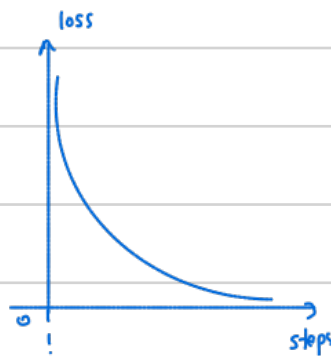


η ← learning rate (size of the step) (hyper parameter)

doesn't change the value in the learning process -

* use a threshold to find the minimal point.

why we not taking gradient = 0? (points may be overfitted/There can be a function which goes through all points)



```
# -----
w = -3.0          # starting point
lr = 0.3          # learning rate
steps = 25        # number of iterations
thresh = 0.001

for i in range(steps):

    gradient = grad_loss(w)
    w_next = w - lr * gradient

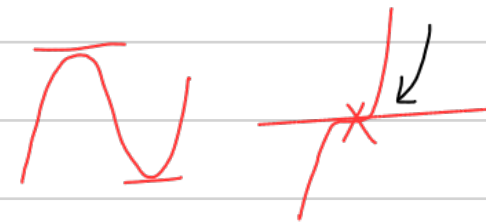
    # Plot the frame
    plot_step(w, w_next, gradient)

    # =====
    # 🐛 STUDENT TASK: Add your stopping criteria here
    # =====
    # OPTION CORRECT
    if (abs(w_next - w) < thresh) and (abs(gradient) < thresh):
        print(i)
        print("Stopping criteria met!")
        break

    # Update w
    w = w_next

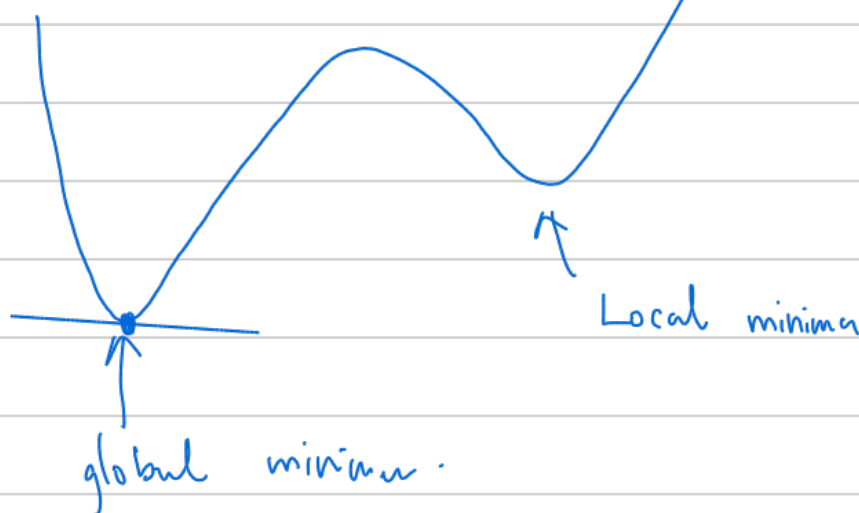
    time.sleep(0.4) # smooth pacing
```

Both we have to check



Importance of the learning Rate

"changing Learning rate we can find the global minima"



pytorch optimizers

SGD (Stochastic Gradient Descent)
Vanilla gradient Descent SGD with (batch size = full dataset)
Mini batch SGD
Momentum.
Adam.
RMSprop
Adagrad.

} automatical track all parameters.
update using gradients
apply learning rate scaling &
other things.