

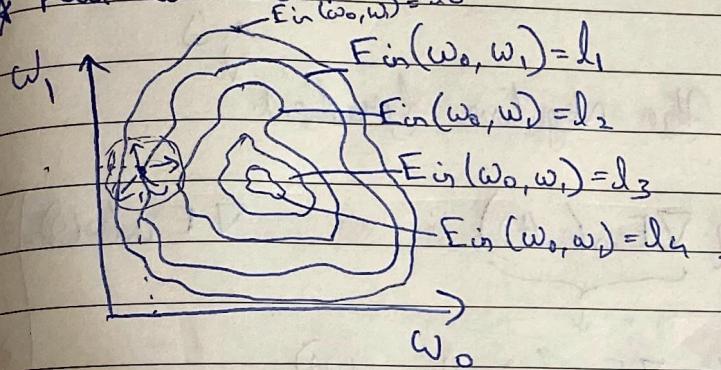
Lecture - 10. FT2qJM

$$\text{Minimize } E_{\text{in}}(\omega) = \frac{1}{N} \sum_{n=1}^N \ln (1 + e^{y_n \omega^T x_n})$$

- \* Method of Maximum Likelihood: Minimize Cross Entropy Error.
- No Analytic Solution, cannot  $(\nabla \omega) = 0$ .

To maximise the probability i.e to minimise the cross entropy error, take derivative.

\* Hill descent:



Contour map.

$$l_0 > l_1 > l_2 > l_3 > l_4$$

To find the required point, take the gradient i.e take the derivative till it becomes 0.

difference  $\Delta E_{\text{in}}(\omega) = E_{\text{in}}(\omega_{\text{new}}) - E_{\text{in}}(\omega)$  (Not a formal defn).

When choosing the direction where our new  $(\omega_0, \omega_1)$  will be based on the step size (radius)

Here the  $\Delta E_{\text{in}}(\omega)$  derivative at point  $h_1 \neq 0$   
at point  $h_3 = 0$ .

takes a

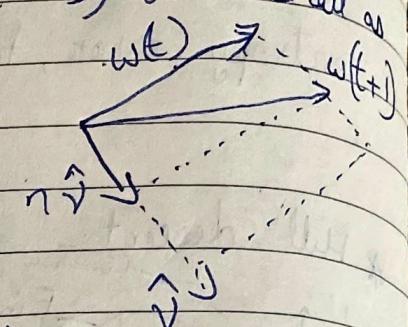
## \* How to Roll down the hill?

- At any step time step  $t$ , suppose we are at weights  $w(t)$ 
  - take a step of size  $\eta$  in direction  $\hat{v}$

$$w(t+1) = w(t) + \eta \hat{v}$$

Greedy: Pick  $\hat{v}$  so  $E_{in}(w(t+1))$  is as small as possible.

$\hat{v}$  = Greedy direction.



Rolling down = Iterating the Negative Gradient.

$$w(t+1) = w(t) - \eta \nabla E_{in}(w(t))$$

$d+1 \times 1$   
vector

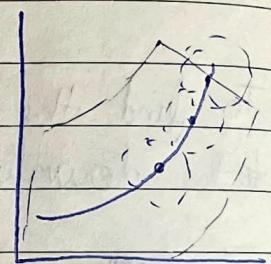
Scalar  
value

$\nabla E_{in}(w_0, w_1, \dots, w_d)$

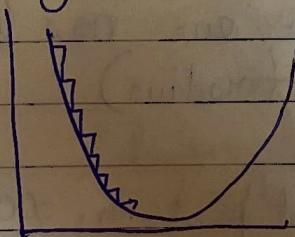
$d+1 \times 1$   
vector.

$$\nabla E_{in}(w(t)) = \begin{bmatrix} \frac{\partial E_{in}}{\partial w_0} \\ \vdots \\ \frac{\partial E_{in}}{\partial w_d} \end{bmatrix}$$

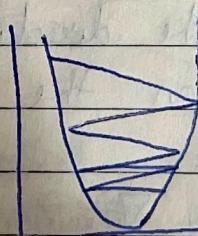
$\hat{v}$  choose direction to reduce  $E_{in}$ .  
Problem is with choosing the right value  $\eta$ .



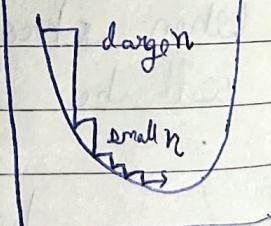
\* Picking the "Learning Rate"  $\eta$  (Convex Func).



$\eta$  too small



$\eta$  too large



Variable  $\eta$

how to drop to small  $\eta$  when the ball jumps to

\* Rolling down.

$$E_{in}(w(t+1)) = E_{in}(w(t) + \eta \hat{v})$$

$$[\text{Taylor approximation}] \approx E_{in}(w(t)) + \eta \hat{v}^T \nabla E_{in}(w(t)) + \frac{\eta^2}{2!} \hat{v}^T \nabla^2 E_{in}(w(t)) \dots$$

$$\approx E_{in}(w(t)) + \eta \hat{v}^T \nabla E_{in}(w(t)) + O(\eta^2)$$

$$\Delta E_{in} = E_{in}(w(t+1)) - E_{in}(w(t))$$

$$\approx \eta \hat{v}^T \nabla E_{in}(w(t)) + O(\eta^2)$$

$\leq 0$ . (We want to reduce the in sample error)

$$\text{Do } E_{in}(w(t+1)) < E_{in}(w(t))$$

Fastest Way to Roll down. (Is the gradient)

$$-1 \leq \cos(\alpha) = \frac{\hat{v}^T \nabla E_{in}(w(t))}{\|\hat{v}\| \|\nabla E_{in}(w(t))\|} \leq 1$$

(substitute  $\|\hat{v}\| \|\nabla E_{in}(w(t))\|$  as c)

$$-c \leq c \cdot \cos(\alpha) = \hat{v}^T \nabla E_{in}(w(t)) \leq c.$$

for  $\cos 180^\circ = -1$ .  $\rightarrow E_{in}(w(t))$

Do we should choose  $\alpha = 180^\circ$ ?  $\xrightarrow{\text{d.}}$

## \* Gradient Descent Algo :-

1. Initialize at step  $t=0$  to  $w(0)$ .
2. For  $t = 0, 1, 2, \dots$  do.
3. Compute the gradient  $g_t = \nabla E_{in}(w(t)) \rightarrow$  large error
4. Move in direction  $v_t = -g_t$ .
5. Update the weights  $w(t+1) = w(t) + \eta v_t$ .
6. Iterate until stopping condition is met.
7. End for.
8. Return final weights.

$$\text{Minimize } E_{in}(w) = \frac{1}{N} \sum_{n=1}^N \ln \left( 1 + e^{-y_n w^T x_n} \right)$$

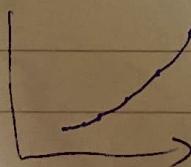
By moving small steps  $-\eta \nabla E_{in}(w)$ .

## \* Stochastic Gradient Descent (SGD) :-

→ Pick a data point  $(x_*, y_*)$  from  $D$  at random.

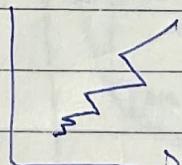
→ Run one iteration of GD on error  $(w, y_*, y_*)$   
 $w(t+1) = w(t) - \eta \nabla \text{error}(w, x_*, y_*)$

$$w(t+1) = w(t) + y_* x_* - \frac{\eta}{1 + e^{y_* w^T x_*}} \quad \boxed{\text{SGD.}}$$



GD

10 steps



SGD

30 steps.

Not necessary to have  
higher lower  $E_{in}$ .  
that may increase low  
by overfitting.