Input layer

Hidden layer

Output layer

$\vec{x}$

$f_w(\vec{x})$: error $\begin{pmatrix} \text{To be} \\ \text{minimized} \end{pmatrix}$

→ Goal: minimize $f_w(\vec{x}) = g(\vec{w})$ w.r.t. $\vec{w}$

saddle point

$\uparrow \nabla^2$ is +eV

Local min

$f_w(\vec{x})$

global min $g(w_{min}) = f_{w_{min}}(\vec{x})$

$$\nabla g(\vec{w}_{min}) = 0 \quad \text{at those three points}$$

To find $\vec{W}_{min}$; Use GD

$$\vec{W}_{t+1} = \vec{W}_t - \varepsilon \cdot \nabla g(\vec{w}_t)$$

Issue:
- cold start
- $\varepsilon$
- local min

Solution: adaptive learning rate

line search

→ Optimization methods:

- zero order:
  evaluate $g(\vec{w})$ at any $\vec{w}$

- 1st order:
  GD, KKT

- 2nd order:
  Newton's method
  (Quasi-Newton's method)

  $$\vec{W}_{t+1} = \vec{W}_t - \epsilon_t \cdot \nabla^2 g(\vec{w_t})^{-1} \nabla g(\vec{w_t})$$

→ Example: Linear regression

1-D



$\hat{y} = wx + b$

Training set $(X_1, Y_1) \cdots\cdots (X_n, Y_n)$

Higher dimension:

Least square error

$$\sum_{i=1}^{n} \left( \vec{y}_i - \vec{w}^T \vec{x}_i \right)$$

Note: LR actually has closed form solution

Using GD:

At time $t$:

$$\nabla g(\vec{w}) = -\sum \left(g_i - \vec{w}^\top \vec{x_i}\right) \cdot \vec{x_i}$$

Update $\vec{x}_{t+1}$

→ Error func

- It's usually decomposed into sum of data

$$g(\vec{w}) = \sum g_i(w, \vec{x}_{ij} \cdot y_i)$$

- It is costly to compute in high dimension

$\rightarrow$ SGD:

- Only compute part of $g(\vec{w})$

$$\vec{W}_{t+1} = \vec{W}_t - \epsilon_t \nabla g_i(\vec{W}_t), \text{ for some } i \in range(n)$$

- In the LR example, SGD can do n iterations in the time of <u>1</u> GD update

- Mini batch update

$\rightarrow$ Regularization

Again LR example:

With L2 norm regulation:

Error func $g(\vec{w}) = \sum (y_i - \vec{w}^T x_i)^2 + \frac{1}{2}\lambda \|\vec{w}\|_2^2$