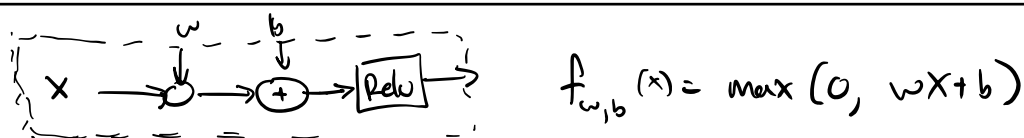


Basic principles review part 2



ReLU fn: $\geq \frac{-b}{w}$

Initialization:

- Current folk wisdom:

(1) use what worked on a related problem

↳ pretrained networks as an initializer

(2) random initializations using Gaussians

↳ how to set variance?

↳ Xavier initialization $\sim \mathcal{N}(0, ?)$ $? = \frac{1}{d}$ ← fan in of unit

↳ for ReLU: He initialization $? = \frac{2}{d}$ ← our assumption is half of ReLU outputs have 0 variance, so we double to get 1

↳ some ReLU elbows "dead"

because of early-dist learning

for out vals \rightarrow sets many to 0 \Rightarrow no ∇

Regularization

- focus on least squares linear regression

$$\text{Cost} = \|\vec{y} - X\vec{w}\|^2 + \lambda \|\vec{w}\|^2 \quad \vec{w}^* = (X^T X + \lambda I)^{-1} X^T \vec{y}$$

Data Augmentation

$$\begin{bmatrix} X \\ \sqrt{\lambda} I \end{bmatrix} \vec{w} = \begin{bmatrix} \vec{y} \\ 0 \end{bmatrix}$$

OVS \Rightarrow gets same as above
*Shown in hw.

Feature Augmentation

$$\begin{bmatrix} X & \sqrt{\lambda} I \end{bmatrix} \begin{bmatrix} \vec{w} \\ \phi \end{bmatrix} = \begin{bmatrix} \vec{s} \\ \vec{y} \end{bmatrix}$$

OBS: min norm solution w/ moore-penrose pseudo inverse

$$\begin{bmatrix} \vec{s} \\ \vec{y} \end{bmatrix} \Rightarrow \hat{\vec{w}} = X^T (X X^T + \lambda I)^{-1} \vec{y}$$

regularized ridge / dual view

$$\begin{bmatrix} X^T \\ \sqrt{\lambda} I \end{bmatrix} \left(\begin{bmatrix} X & \sqrt{\lambda} I \end{bmatrix} \begin{bmatrix} X^T \\ \sqrt{\lambda} I \end{bmatrix} \right)^{-1} \vec{y}$$

pseudo-inverse