

$$\min \underline{L} = -\frac{1}{N} \sum_{i=1}^N (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i))$$

$\hat{y}_i = \sigma(\underline{w}^T x_i + b)$ Convex, differentiable
 $a_i^{(i)} = \sigma(\underline{w}^T x_i + b)$

Algorithm: Gradient Descent $A \Rightarrow \hat{Y}$

Batch Gradient Descent

1) Compute gradients:

$$\frac{\partial L}{\partial \underline{w}} = \frac{1}{N} X (A - Y)^T, \quad \frac{\partial L}{\partial b} = \frac{1}{N} \sum_{i=1}^N (a_i^{(i)} - y_i)$$

2) Update Parameters

$$\underline{w}_{t+1} = \underline{w}_t - \eta \frac{\partial L}{\partial \underline{w}}$$

$$b_{t+1} = b_t - \eta \frac{\partial L}{\partial b}$$

→ No. of epochs (known domain)

→ Threshold

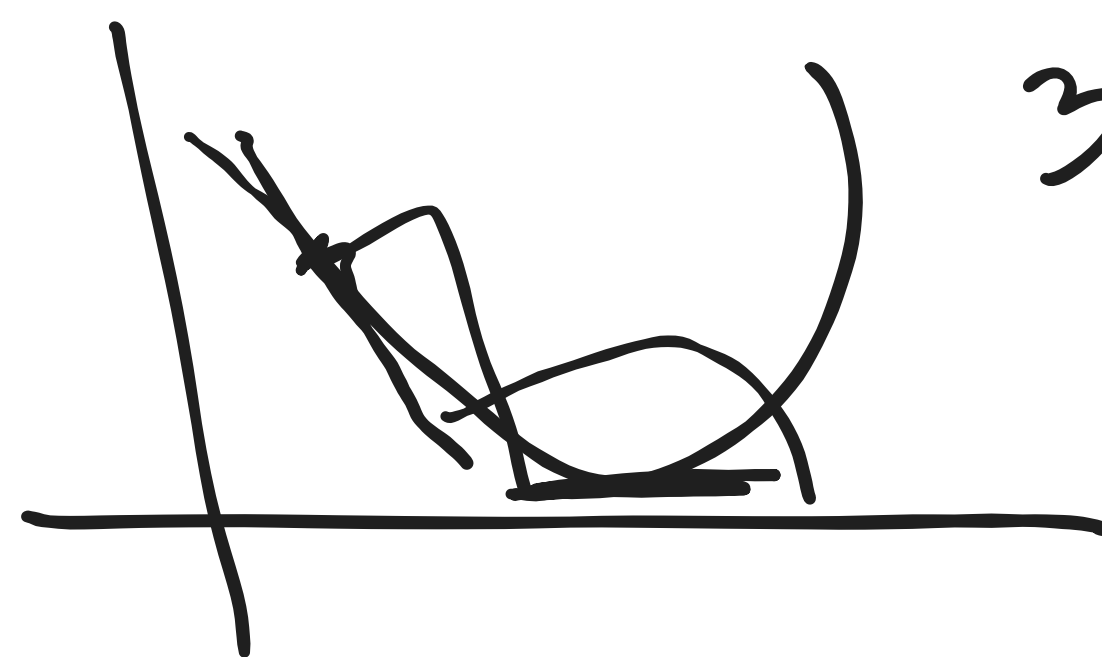
Vanilla Gradient Descent

$$\underline{w}_{t+1} = \underline{w}_t - \eta \frac{dL}{d\underline{w}}$$

1) Batch (entire Data)

2) Stochastic (One observation per update)

3) Mini-Batch $1 < B < N$



1) Momentum (SAD)

$$\underline{w}_{t+1} = \underline{w}_t - \boxed{V_t} \leftarrow$$

$$V_t = \rho V_{t-1} + \eta \nabla L(\underline{w})$$

$0.9 < \rho < 1$