## Gradient Descent

$\min f(x)$
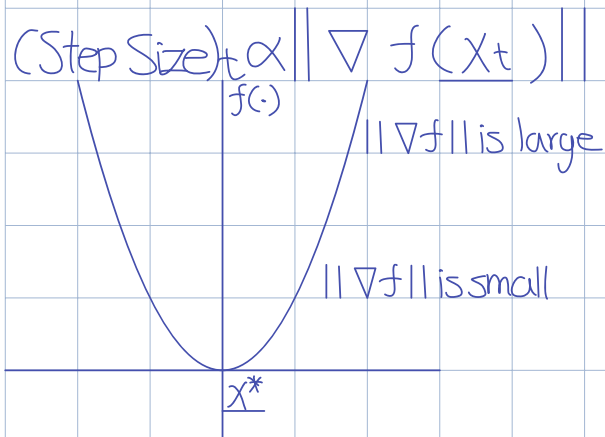
$x \in \mathbb{R}^N$

$f(x)$ is a convex function

## Gradient Update

$$X_{t+1} = X_t - (\text{step size})_t \, U_t$$

$$U_t = \frac{\nabla f(X_t)}{\|\nabla f(X_t)\|}$$

$$(\text{Step Size})_t \propto \|\nabla f(X_t)\|$$

$f(\cdot)$

$\|\nabla f\|$ is large

$\|\nabla f\|$ is small

$\underline{x^*}$

$$(\text{Step Size})_t = \varepsilon_t \|\nabla f(X_t)\|$$

$\downarrow$

Learning Rate

## Gradient Update

$$X_{t+1} = X_t - \varepsilon_t \nabla f(X_t)$$

$\downarrow$

Learning Rate

## Gradient Descent Algorithm

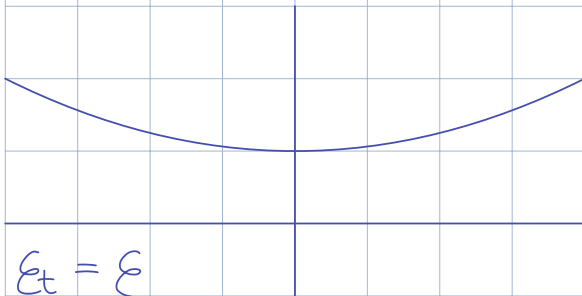Initialize $x_0$ in some arbitrary fashion

$t = 0, 1, 2, \ldots$

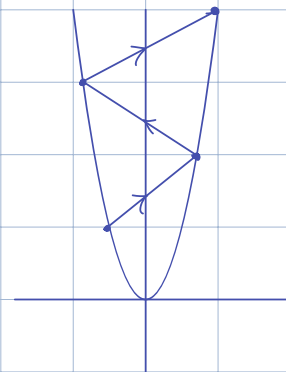Compute $g_t = \nabla f(x_t)$

Select Direction $u_t = -g_t$

Update $x_{t+1} = x_t + \varepsilon_t u_t$

Continue until stopping condition is reached for convec function

$\|\nabla f(x_t)\| \approx 0$



$\varepsilon_t = \varepsilon$
(Very Small)

$\varepsilon_t = \varepsilon$ (Very Large)

$\varepsilon_t \approx \dfrac{1}{t}$ (Proportional to $1/$iteration index)

## Linear Regression

Loss Function

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} e_n(w)$$

$$e_n(w) = (w^T x_n - y_n)^2$$

$$\underline{w_{LS}^*} = \underset{\underline{w} \in \mathbb{R}^{d+1}}{\arg\min} \; E_{in}(\underline{w})$$

$$= (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \underline{y}$$

Might Not Be Invertible

Use Gradient Descent To Minimize $E_{in}(\underline{w})$

Update Rule

$$\underline{w}_{t+1} = \underline{w}_t - \varepsilon_t \nabla E_{in}(\underline{w}_t)$$

$$\nabla E_{in}(\underline{w}_t) = \nabla \left( \frac{1}{N} \sum_{n=1}^{N} e_n(\underline{w}) \right)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \nabla e_n(\underline{w})$$

$$\nabla_{\underline{w}} e_n(\underline{w}) = \nabla_{\underline{w}} (w^T x_n - y_n)^2$$

$$= 2(w^T x_n - y_n) \nabla_w (w^T x_n - y_n)$$

$$= 2(w^T x_n - y_n) x_n$$

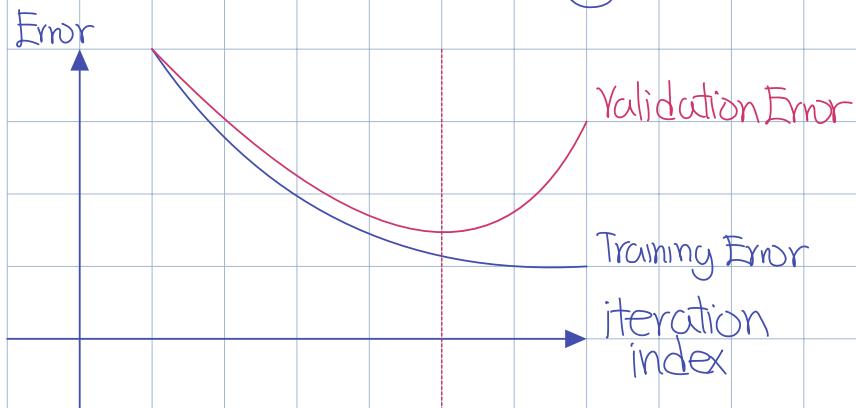$$\boxed{\underline{w}_{t+1} = \underline{w}_t - \varepsilon_t \frac{2}{N} \sum_{n=1}^{N} (w^T x_n - y_n) \underline{x_n}}$$

Batch Gradient Descent $O(Nd)$

Iterative Procedure

$\underline{w_t}$ will converge to $\underline{w_{LS}^*}$

Two Concerns

1) $(\mathcal{F}^T \mathcal{F})$ could be ill conditioned

2) $(\mathcal{F}^T \mathcal{F})^{-1}$ is computationally expensive

Error

Validation Error

Training Error

iteration index

Iterative Algorithm

Input Data

80% Training Data

20% Validation Data

Batch Gradient Descent

Complexity per update $O(Nd)$

Stochastic Gradient Descent

$$W_{t+1} = W_t - \mathcal{E}_t \nabla E_n(W_t)$$

$$n \in \text{Uniform} \{1, 2, \dots, N\}$$
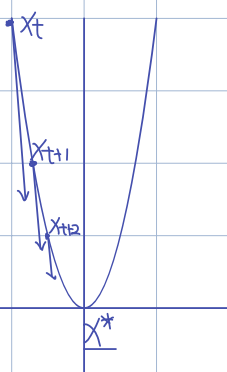
Random Variable

## Justification of SGD

$$\mathbb{E}_n\left[\nabla e_n(\underline{w_t})\right] = \sum_{n=1}^{N} Pr(n=1)\, \nabla e_n(\underline{w_t})$$
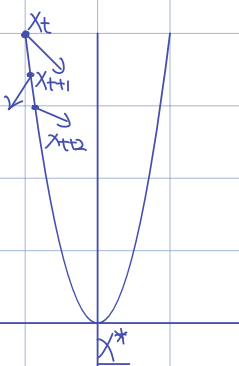
$$= \nabla E_{in}(\underline{w_t}) \qquad \text{Matches True Gradient In}$$

Expectation

Expectation



Batch Gradient Descent            SGD

## Mini Batch Gradient Descent

at each iteration draw $M$ examples from the dataset at random (without replacement)

$$S = \{ j_1, j_2, \ldots, j_M \}$$

Update Rule

$$\underline{w_{t+1}} = \underline{w_t} - \frac{\varepsilon_t}{M} \sum_{n=1}^{M} \nabla e_{(jn)}(\underline{w_t})$$

$M=N$, Batch Gradient Update

$M=1$, SGD

# GD For Logistic Regression

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} e_n(w)$$

$$e_n(w) = \log(1 + e^{-y_n w^T x_n})$$

## SGD Update

$$W_{k+1} = W_k - \epsilon_k \nabla e_n(W_k)$$

$$n \in \text{Uniform } \{1, 2, \dots, N\}$$

$$\nabla_{W_k} e_n(W_k) = \nabla_{W_k} \left[ \log(1 + e^{-y_n W_k^T x_n}) \right]$$

By Chain Rule and Calculus

$$= \frac{1}{1 + e^{-y_n W_k^T x_n}} \nabla_{W_k} (1 + e^{-y_n W_k^T x_n})$$

$$= \frac{e^{-y_n W_k^T x_n}}{1 + e^{-y_n W_k^T x_n}} \nabla_{W_k} (-y_n W_k^T x_n)$$

$$= \frac{1}{1 + e^{y_n W_k^T x_n}} (-y_n x_n)$$

## Update Rule

$$W_{k+1} = W_k + \epsilon_k \left\{ \frac{y_n x_n}{1 + e^{y_n W_k^T x_n}} \right\}$$

## "Highly Misclassified Point"

$$y_n W_k^T X_n \ll 0 \quad (\text{Large And Negative})$$

$$W_{k+1} = W_k + \epsilon_k y_n X_n$$

## Well Classified Point

$$y_n W_k^T X_n \gg 0 \quad (\text{Large And Possitive})$$

$$W_{k+1} \approx W_k$$

## Multiple Classes Logistic Regression

$$y \in \{1, 2, \ldots, c\}, \quad c > 2$$

**Input Data**

$$D = \{(X_1, y_1), \ldots, (X_N, y_N)\}$$

$$X_i \in \mathbb{R}^{d+1}, \quad y_i \in \{1, 2, \ldots, c\}$$

We will consider the case: $c = 3$

Let $w(1), w(2)$ and $w(3)$ be the weight be the weight vectors for class 1, 2 and 3

$$\Omega = \{w(1), w(2), w(3)\}$$

**Output:**

$$[\hat{P}_\Omega(1 \mid x), \ \hat{P}_\Omega(2 \mid x), \ \hat{P}_\Omega(3 \mid x)]$$

$$\hat{P}_{\Omega}(i \mid X) = \frac{e^{w^T(i)X}}{e^{w^T(1)X} + e^{w^T(2)X} + e^{w^T(3)X}}$$

$i \in \{1, 2, 3\}$

## Loss Function

Log Loss Function

Given $(X_n, y_n)$

$y_n \in \{1, 2, 3\}$

$$E_n(\Omega) = -\log \hat{P}_{\Omega}(y_n \mid X_n)$$

$$\Omega^* = \underset{\Omega = \{w(1), w(2), w(3)\}}{\arg\min} E_n(\Omega) = \{w(1)^*, w(2)^*, w(3)^*\}$$

$C = 2$ (Binary Classification)

$\Omega = \{w(1), w(2)\}$

$$\hat{P}_{\Omega}(1 \mid X) = \frac{e^{w(1)^T X}}{e^{w(1)^T X} + e^{w(2)^T X}} = \frac{e^{(w(1) - w(2))^T X}}{1 + e^{(w(1) - w(2))^T X}}$$

$$\hat{P}_{\Omega}(2 \mid X) = \frac{e^{w(2)^T X}}{e^{w(1)^T X} + e^{w(2)^T X}} = \frac{1}{1 + e^{(w(1) - w(2))^T X}}$$

$$\Omega = \{w\}$$

$$\hat{P}_{\Omega}(1 \mid x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

$$\hat{P}_{\Omega}(2 \mid x) = \frac{1}{1 + e^{w^T x}}$$

Thus if $w(1) - w(2) = w$

Then the models will output same probabilities

SGD Update Rule for

Mutti Class Logistic Regression

Iteration #k

Model Parameter

$$\Omega_k = \{ w_k(1), w_k(2), w_k(3) \}$$

$$w_{k+1}(1) = w_k(1) - \varepsilon_k \nabla_{w_k(1)} e_n(\Omega_k)$$

$$w_{k+1}(2) = w_k(2) - \varepsilon_k \nabla_{w_k(2)} e_n(\Omega_k)$$

$$w_{k+1}(3) = w_k(3) - \varepsilon_k \nabla_{w_k(3)} e_n(\Omega_k)$$

$$
\begin{bmatrix} W_{k+1}(1) \\ W_{k+1}(2) \\ W_{k+1}(3) \end{bmatrix} = \begin{bmatrix} W_k(1) \\ W_k(2) \\ W_k(3) \end{bmatrix} - \varepsilon_k \nabla_{\Omega_k} E_n(\Omega_k)
$$

$$
\Omega_{k+1} \qquad \Omega_k
$$

Compute

$$
\nabla_{W_k(\ell)} \left[ -\log \hat{P}_\Omega(y_n | x_n) \right] \qquad \ell \in \{1,2,3\}
$$

$$
= \nabla_{W_k(\ell)} \left[ -\log \frac{e^{W_k^T(y_n) x_n}}{e^{W^T(1) x} + e^{W^T(2) x} + e^{W^T(3) x}} \right]
$$

$$
= \nabla_{W_k(\ell)} \left[ - W_k^T(y_n) x_n + \log \left( e^{W^T(1) x} + e^{W^T(2) x} + e^{W^T(3) x} \right) \right]
$$

If $\ell = y_n$

$$
= -x_n + \nabla_{W_k(y_n)} \left[ \log \left( \sum_{c=1}^{3} e^{W_k^T(c) x_n} \right) \right]
$$

$$
= -x_n + \frac{1}{\sum_{c=1}^{3} e^{W_k^T(c) x_n}} e^{W_k^T(\ell) x} x_n
$$

If $c \neq y_n$

$$= \frac{1}{\sum_{i=1}^{3} e^{w_k^T c(i) x_n}} e^{w_k^T c(1) x} x_n$$