

# Convex Functions

- Any local minimizer  $\mathbf{x}^*$  is a global minimizer of  $f$  if  $f$  is convex.
- Further if  $f$  is differentiable, then any stationary point  $\mathbf{x}^*$  is a global minimizer of  $f$

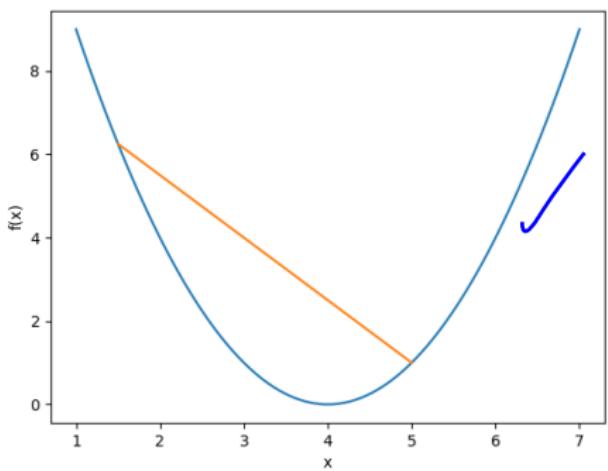


Figure:  $f(x) = (x - 4)^2$  is convex on  $\mathbb{R}$  and  $x^* = 4$  is both global and local minimizer of  $f(x)$

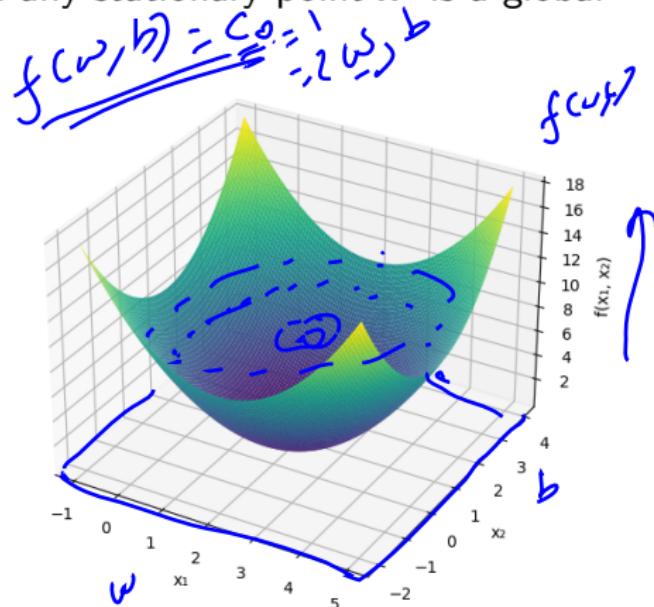


Figure:  $f(x_1, x_2) = (x_1 - 2)^2 + (x_2 - 1)^2$  is convex on  $\mathbb{R}^2$  and  $x^* = (2, 1)$  is both global and local minimizer of  $f(x_1, x_2)$

# Unconstrained Minimization

## Unconstrained Minimization Algorithm:

- ① Initialize  $\mathbf{x}_0$  and  $i = 0$
- ② If stopping condition is not satisfied, then continue, else stop
  - ① Calculate  $\mathbf{x}_{i+1}$  so that  $f(\mathbf{x}_{i+1}) < f(\mathbf{x}_i)$
  - ② Update  $i$  with  $i + 1$
- ③ Final  $\mathbf{x}_i$  is the local minimum  $\mathbf{x}^*$  of  $f(\mathbf{x})$

## Following questions need to be answered:

- What stopping condition can be used?
- What is the speed of convergence?
- How to calculate  $\mathbf{x}_{i+1}$ ?

# Unconstrained Minimization

## Stopping Conditions:

- $\|\nabla f(\mathbf{x}_i)\| \leq \epsilon$  ✓
- $\frac{f(\mathbf{x}_i) - f(\mathbf{x}_{i+1})}{|f(\mathbf{x}_i)|} \leq \epsilon$

## Speed of Convergence:

The sequence  $\{\mathbf{x}_i\}$  converges to  $\mathbf{x}^*$  with order  $p$  and convergence rate  $\beta$  if

$$\lim_{i \rightarrow \infty} \frac{\|\mathbf{x}_{i+1} - \mathbf{x}^*\|}{\|\mathbf{x}_i - \mathbf{x}^*\|^p} = \beta, \quad \beta \in \mathbb{R}$$

- Convergence is faster for higher  $p$
- Linear Convergence:  $p = 1, 0 < \beta < 1$
- Quadratic Convergence:  $p = 2, \beta > 0$

$$\|\mathbf{x}_{i+1} - \mathbf{x}^*\| \leq \beta \|\mathbf{x}_i - \mathbf{x}^*\|^p$$

# Unconstrained Minimization

Calculating  $\mathbf{x}_{i+1}$  for the Unconstrained Minimization Algorithm:

- First Order Methods
  - Gradient Descent
  - Stochastic Gradient Descent
  - Mini-Batched Gradient Descent
  - Stochastic Average Gradient
  - Optimizers - AdaGrad, RMSProp & Adam
- Second Order Methods
  - Newton's Method
  - Quasi-Newton Method

# Nonlinear Least Squares Problem

The problem is to fit  $g(x) = \frac{1}{1+e^{-(wx+b)}}$  to a training set with observations  $(x_1, x_2, \dots, x_n)$  and corresponding estimated responses  $(y_1, y_2, \dots, y_n)$  using least squares.

The objective function (or the loss function) to be minimized is:

$$f(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (g(x_i) - y_i)^2$$

$$f(w, b) = \frac{1}{2} \sum_{i=1}^n \left( \frac{1}{1 + e^{-(wx_i+b)}} - y_i \right)^2$$

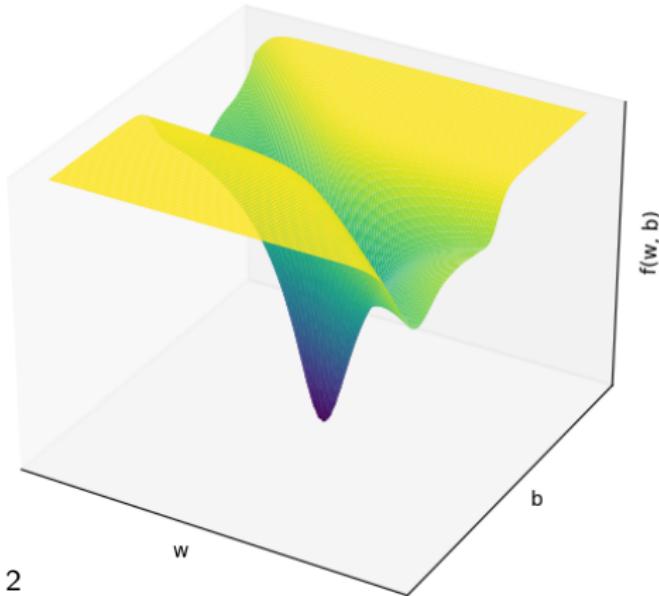


Figure: Surface plot of Loss Function  $f(w, b)$

# Gradient Descent

$$-\mathbf{w}^* \sim \mathbf{b}$$

Descent Direction:

$\mathbf{d} \in \mathbb{R}^n$  is a *descent direction* of  $f(\mathbf{w})$  at  $\mathbf{w}^* \in \mathbb{R}^n$  if  $f(\mathbf{w}^* + \eta \mathbf{d}) < f(\mathbf{w}^*)$  for all  $\eta \in (0, \delta)$ ,  $\delta > 0$

$$\delta \quad f(\mathbf{w}^* + \eta \mathbf{d}) < f(\mathbf{w}^*)$$

- In terms of minimization,  $f(\mathbf{w}_{i+1}) < f(\mathbf{w}_i)$ , where  $\mathbf{w}_{i+1} = \mathbf{w}_i + \eta_i \mathbf{d}_i$
- $\eta_i = \text{minimize } f(\mathbf{w}_i + \eta \mathbf{d}_i)$ ,  $\eta > 0$
- $\eta_i = \text{minimize } h(\eta)$ ,  $\eta > 0$
- $\boxed{\mathbf{d}_i = -\nabla f(\mathbf{w}_i)}$  in gradient descent algorithm

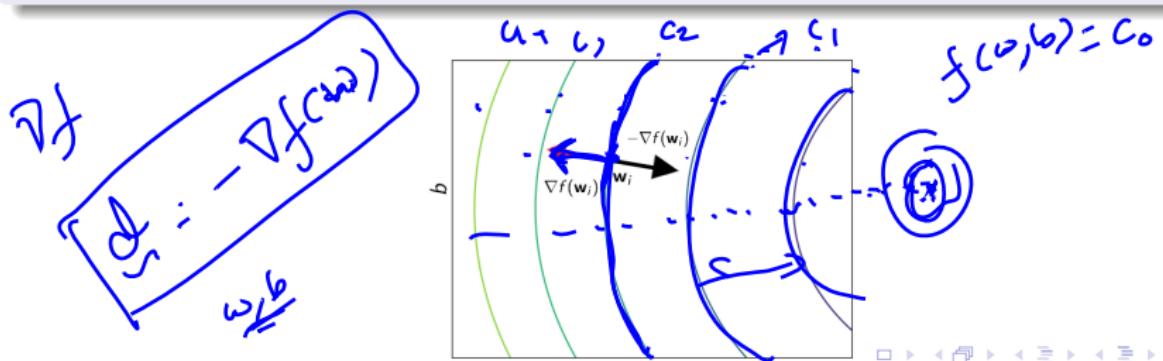
$$\mathbf{w}_{i+1} = \mathbf{w}_i + \eta_i \mathbf{d}_i$$

$$f(\mathbf{w}_{i+1}) < f(\mathbf{w}_i)$$

# Gradient Descent

## Gradient Descent Algorithm:

- ① Initialize  $\mathbf{w}^0$  and  $k = 0$
- ② If stopping condition like  $\|\nabla f(\mathbf{w}^k)\| \leq \epsilon$  is satisfied, then stop, else continue
  - ①  $\mathbf{d}^k = -\nabla f(\mathbf{w}^k)$
  - ② Calculate  $\eta^k$  along  $\mathbf{d}^k$  so that  $f(\mathbf{w}^k + \eta^k \mathbf{d}^k) < f(\mathbf{w}^k)$
  - ③  $\mathbf{w}^{k+1} = \mathbf{w}^k + \eta^k \mathbf{d}^k$
  - ④ Update  $k$  with  $k + 1$
- ③ Final  $\mathbf{w}^k$  is the local minimum  $\mathbf{w}^*$  of  $f(\mathbf{w})$



# Gradient Descent

- The objective function (or the loss function) to be minimized is:

$$f(\mathbf{w}) = f(w, b) = \frac{1}{2} \sum_{i=1}^n (g(x_i) - y_i)^2 = \frac{1}{2} \sum_{i=1}^n \left( \frac{1}{1 + e^{-(wx_i+b)}} - y_i \right)^2$$

- The update with learning rate  $\eta$ :

$$\nabla f = \sum_i \nabla f_i$$

$$\begin{aligned} \begin{bmatrix} w^{k+1} \\ b^{k+1} \end{bmatrix} &:= \begin{bmatrix} w^k \\ b^k \end{bmatrix} - \eta \left[ \frac{\partial}{\partial w} \left( \frac{1}{2} \sum_{i=1}^n (g(x_i) - y_i)^2 \right) \right] \\ &= \begin{bmatrix} w^k \\ b^k \end{bmatrix} - \eta \left[ \frac{\sum_{i=1}^n (g(x_i) - y_i) * g(x_i) * (1 - g(x_i)) * x_i}{\sum_{i=1}^n (g(x_i) - y_i) * g(x_i) * (1 - g(x_i))} \right] \end{aligned}$$

- In Gradient Descent, gradients have to be evaluated for every training set before updating which can be expensive for large training.
- Convergence rate: Linear
- Iteration cost:  $O(n)$

cost per iteration  $O(n)$



# Gradient Descent

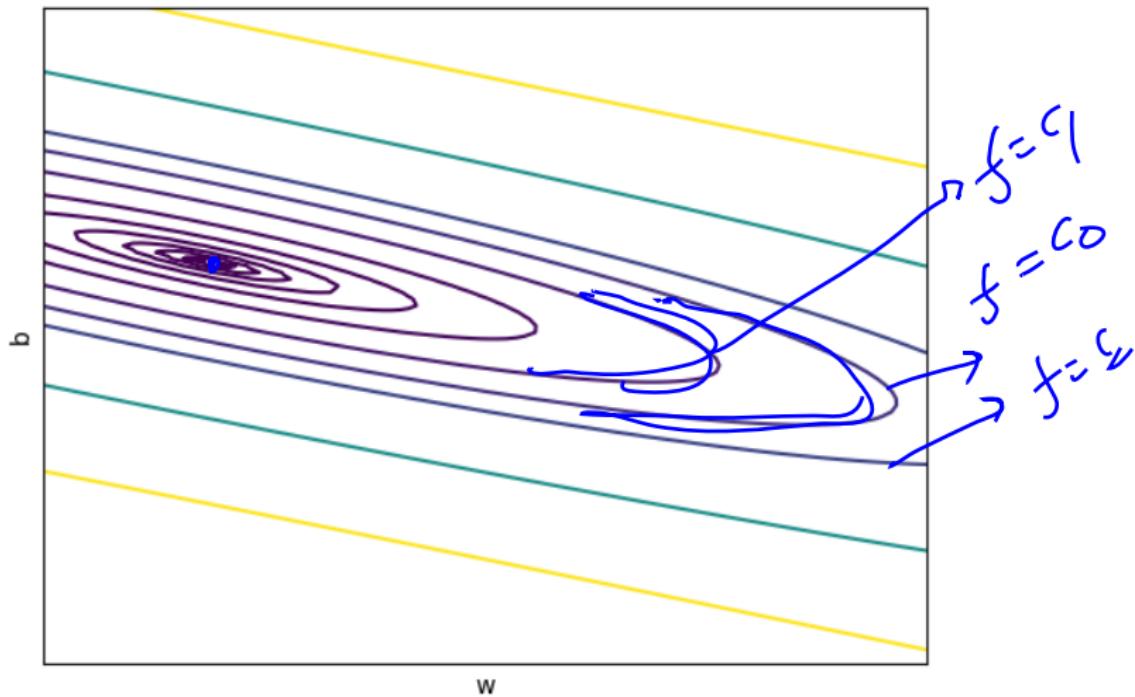
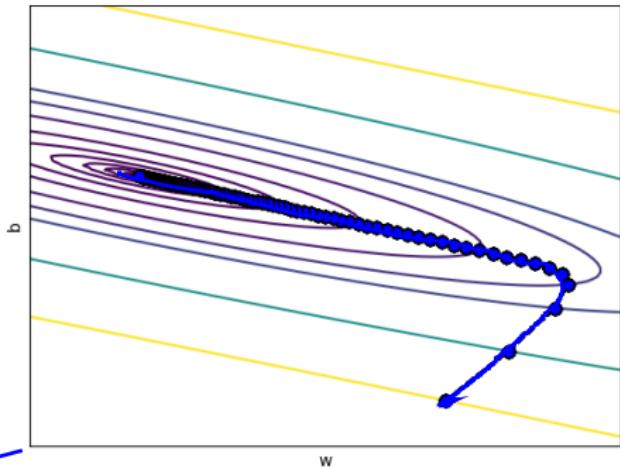


Figure: Contour plot of Loss function  $f(w, b)$

# Gradient Descent



$$\text{Epoch } \rightarrow \frac{n=1000}{\sum_{i=1}^n (g(x_i) - g_i)^2}$$

Figure: Epoch: 175 ✓  
Total Gradient Evaluations: 175000  
Number of times you go through all the training data

# Stochastic Gradient Descent

## Stochastic Gradient Descent:

- ① Initialize  $\mathbf{w}_0$
- ② Repeat until an approximate minimum is obtained:
  - ① Randomly shuffle examples in the training set.
  - ② For  $i = 1, 2, \dots, n$ , do:
    - ①  $\mathbf{w} := \mathbf{w} - \eta \nabla f_i(\mathbf{w})$
- ③ Final  $\mathbf{w}$  is the local minimum  $\mathbf{w}^*$  of  $f(\mathbf{w})$

- Here  $f_i(\mathbf{w}) = \frac{1}{2} (g(x_i) - y_i)^2$

- The update with learning rate  $\eta$ :

$$\begin{bmatrix} w \\ b \end{bmatrix} := \begin{bmatrix} w \\ b \end{bmatrix} - \eta \begin{bmatrix} (g(x_i) - y_i) * g(x_i) * (1 - g(x_i)) * x_i \\ (g(x_i) - y_i) * g(x_i) * (1 - g(x_i)) \end{bmatrix}$$

- In Stochastic Gradient Descent, update is done with gradients of just one data point at a time. This significantly reduces computation compared to Gradient Descent.
- Convergence rate: Sub-Linear; Iteration cost:  $O(1)$

# Stochastic Gradient Descent



$\eta = 10^{-6}$

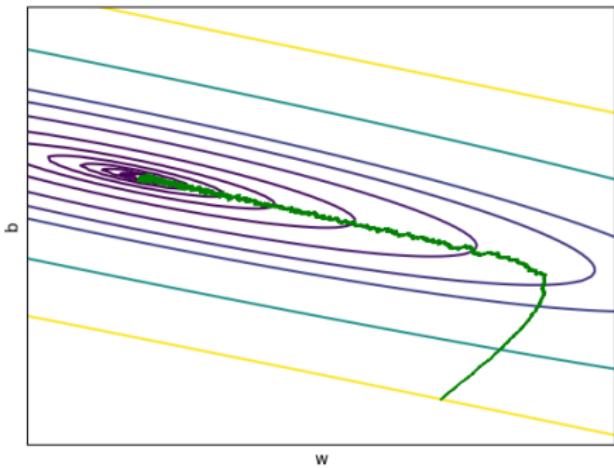
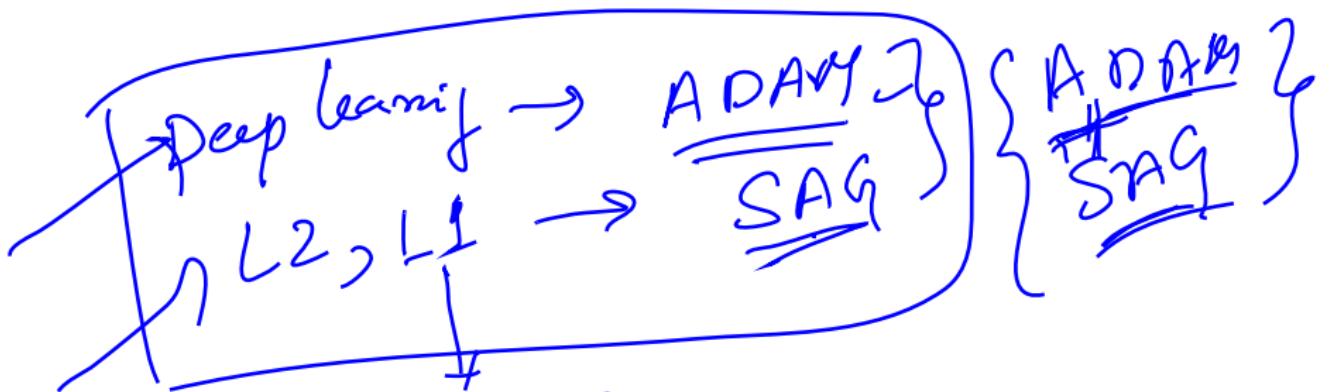


Figure: Epoch: 110  
Total Gradient Evaluations: 110000



$\overline{\text{LBFGS}}$   
second-order  
method

# Mini-Batched Gradient Descent

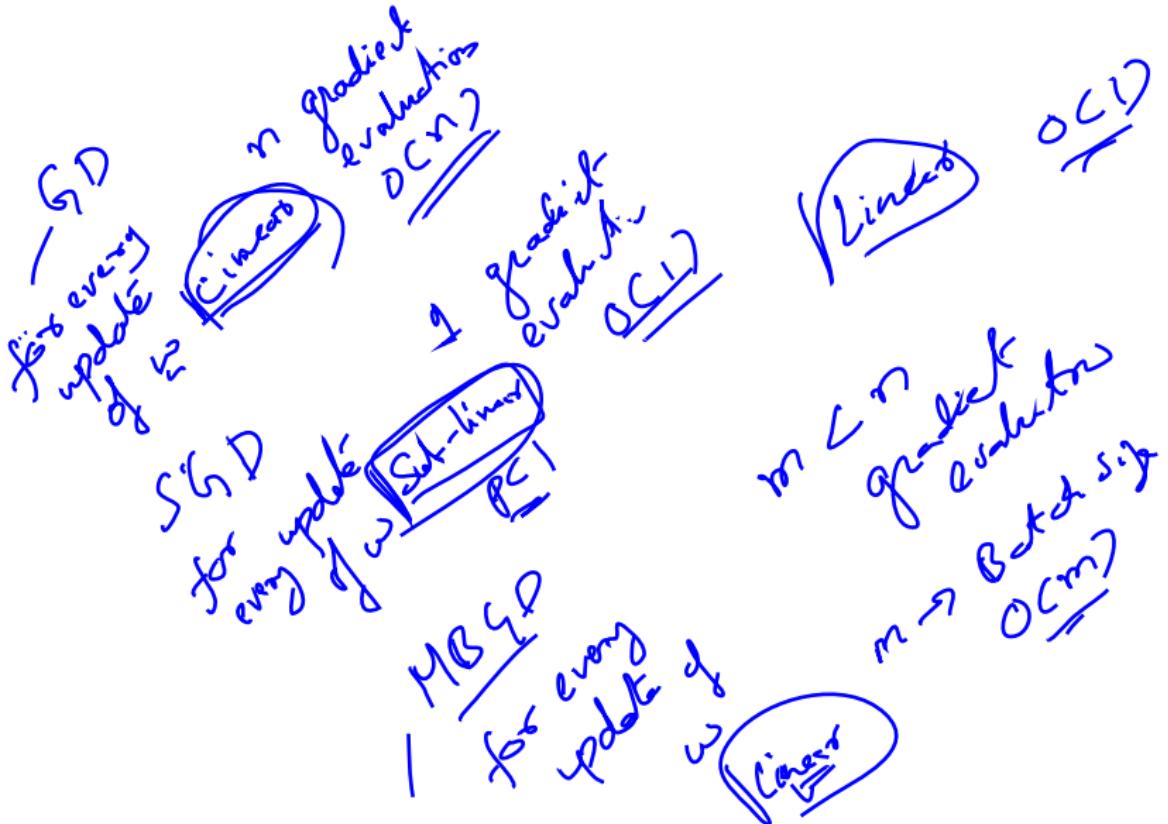
## Algorithm:

- ① Initialize  $\mathbf{w}_0$
- ② Repeat until an approximate minimum is obtained:
  - ① Randomly shuffle examples in the training set.
  - ② Repeat for  $n/m$  times:
    - ① Calculate  $\sum_i \nabla f_i(\mathbf{w})$  for a batch size of  $m$
    - ②  $\mathbf{w} := \mathbf{w} - \eta * \sum_i \nabla f_i(\mathbf{w})$
  - ③ Final  $\mathbf{w}$  is the local minimum  $\mathbf{w}^*$  of  $f(\mathbf{w})$

- Here  $f_i(\mathbf{w}) = \frac{1}{2} (g(x_i) - y_i)^2$
- The update with learning rate  $\eta$ :

$$\begin{bmatrix} w \\ b \end{bmatrix} := \begin{bmatrix} w \\ b \end{bmatrix} - \eta \left[ \begin{array}{l} \sum_i (g(x_i) - y_i) * g(x_i) * (1 - g(x_i)) * x_i \\ \sum_i (g(x_i) - y_i) * g(x_i) * (1 - g(x_i)) \end{array} \right]$$

- In Mini-Batched Gradient Descent, gradients are evaluated for a batch of data training (size  $m$ ),  $m < n$ .
- Convergence rate: Linear, Iteration cost:  $O(m)$



# Mini-Batched Gradient Descent



$n=1000$   
 $m=50$

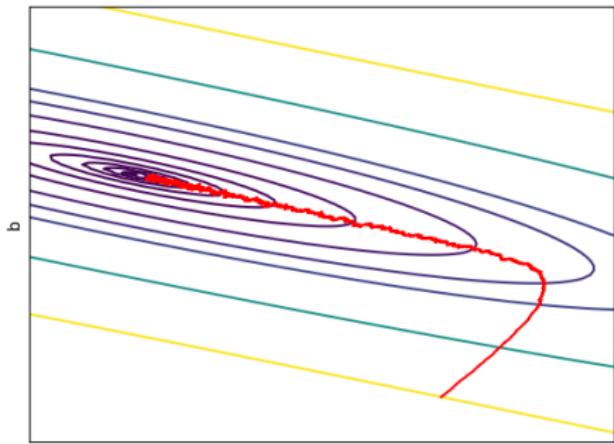


Figure: Epoch: 95  
Total Gradient Evaluations: 95000

# Stochastic Average Gradient (SAG)

- Like stochastic gradient (SG) methods, the SAG method's iteration cost is independent of the number of terms in the sum.
- However, by incorporating a memory of previous gradient values the SAG method achieves a faster convergence rate than black-box SG methods.
- The update at  $k^{th}$  iteration for  $i^{th}$  training data:

$$\mathbf{v}_i^k = \begin{cases} \nabla f_i(\mathbf{w}^k) & \text{if } i = i_k \\ \mathbf{v}_i^{k-1} & \text{otherwise} \end{cases}$$

$$\begin{bmatrix} w \\ b \end{bmatrix} := \begin{bmatrix} w \\ b \end{bmatrix} - \eta \begin{bmatrix} \sum_i (\mathbf{v}_i)_w \\ \sum_i (\mathbf{v}_i)_b \end{bmatrix}$$

Diagram illustrating the SAG update rule:

The diagram shows a large bracket labeled  $\nabla f_i(\mathbf{w})$  above a summation symbol. The summation symbol has two parts: one from  $i=1$  to  $i=k$  labeled  $\nabla f_i(\mathbf{w})$ , and another from  $i=k+1$  to  $i=n$  labeled  $\nabla f_i(\mathbf{w})$ . A blue arrow points from the term  $\nabla f_i(\mathbf{w})$  in the first summation to the update equation. A red arrow points from the term  $\nabla f_i(\mathbf{w})$  in the second summation to the update equation.

# Stochastic Average Gradient (SAG)

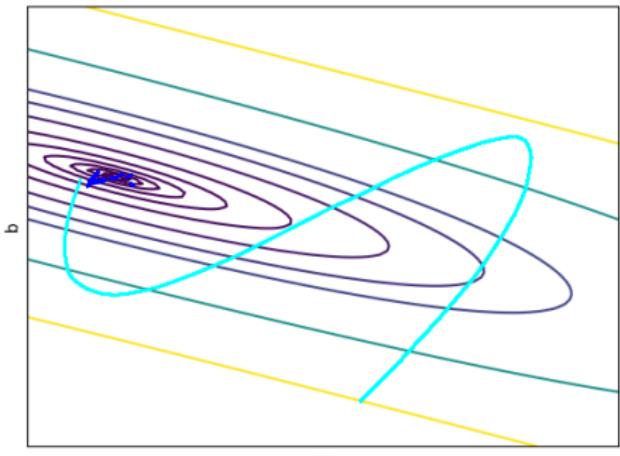


Figure: Epoch: 365  
Total Gradient Evaluations: 365

$\theta(365)$   
 $x \theta(8)$   
 $\theta(365) x$

# Challenges in gradient descent type methods

$$y = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

- Choice of proper learning rate can be difficult. (Too small → slow convergence, too large → can hinder convergence)
- Having a same learning rate for all parameter updates is problematic when the data is sparse.
- Lot of research has been focused on developing adaptive learning rate for Gradient Descent type methods widely used in Deep learning community.

# Adaptive learning rate method: AdaGrad (Adaptive Gradient Descent)

- Perform larger updates (higher learning rate) for parameters corresponding to infrequent features and smaller updates (smaller learning rate) for parameters corresponding to frequent features.
- Update rule:

$$\begin{bmatrix} v_w \\ v_b \end{bmatrix} := \begin{bmatrix} v_w \\ v_b \end{bmatrix} + \left[ \begin{array}{c} \left( \frac{\partial f(w, b)}{\partial w} \right)^2 \\ \left( \frac{\partial f(w, b)}{\partial b} \right)^2 \end{array} \right]$$
$$\rightarrow \begin{bmatrix} w \\ b \end{bmatrix} := \begin{bmatrix} w \\ b \end{bmatrix} - \left[ \begin{array}{c} \frac{\eta}{\sqrt{v_w + \epsilon}} * \frac{\partial f(w, b)}{\partial w} \\ \frac{\eta}{\sqrt{v_b + \epsilon}} * \frac{\partial f(w, b)}{\partial b} \end{array} \right]$$

Diagram illustrating the AdaGrad update rule. A coordinate system shows axes for  $w$  and  $b$ . A point  $(w, b)$  is shown with a vector  $v_w + \epsilon$  originating from it. The update rule adds a diagonal matrix (the Hessian of the function) to the current velocity vector  $v_w$  and  $v_b$ , resulting in a new velocity vector where the diagonal elements are the squares of the partial derivatives. This leads to the final update step, where the parameters  $w$  and  $b$  are updated by subtracting scaled gradients. The scaling factor  $\eta$  is divided by the square root of the sum of the squared partial derivatives plus a small constant  $\epsilon$ .

- For sparse training data,  $w$  corresponding to infrequent feature will have higher learning rate and  $b$  will have smaller updates. Over time the effective learning rate for  $b$  will decay to an extent that there will be no further updates to  $b$ .

$$\begin{pmatrix} \omega^{k+1} \\ b^{k+1} \end{pmatrix} := \begin{pmatrix} \omega^k \\ b^k \end{pmatrix} - \eta \begin{pmatrix} \frac{\partial f}{\partial \omega} \\ \frac{\partial f}{\partial b} \end{pmatrix} \quad \boxed{f = \frac{1}{2} \sum_{i=1}^n (g(x_i) - y_i)^2}$$

$$= \begin{pmatrix} \omega^k \\ b^k \end{pmatrix} - \eta \begin{pmatrix} \sum_{i=1}^n (g(x_i) - y_i) g(x_i) (1 - g(x_i)) \\ \sum_{i=1}^n (g(x_i) - y_i) g(x_i) (1 - g(x_i)) \end{pmatrix}$$

$$x = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix}$$

$$\begin{aligned} \eta_\omega &= \frac{\eta}{\sqrt{V_\omega + \Sigma}} & i=3 \\ \eta_b &= \frac{\eta}{\sqrt{V_b + \Sigma}} & i=1 \end{aligned}$$

# AdaGrad



$$g^{(t)} + = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

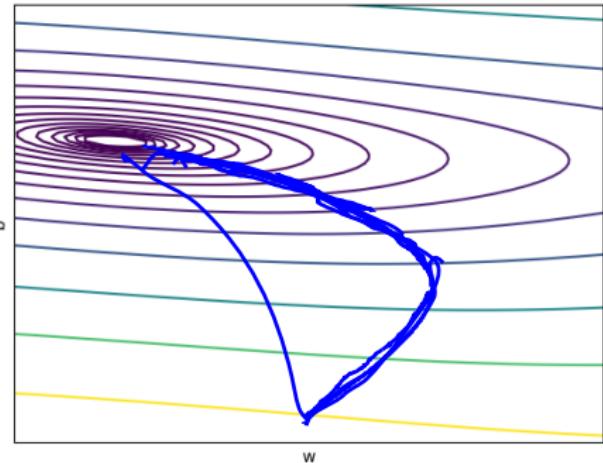


Figure: Epoch: 250  
Total Gradient Evaluations: 250000

# Adaptive Learning Rate Method: RMSProp

- In AdaGrad, the parameters corresponding to frequent features will start receiving very small updates because of the decayed learning rate. To avoid this the denominator is decayed in RMSProp to prevent its rapid growth.
- Update rule:

$$\begin{aligned} \left[ \begin{array}{c} v_w \\ v_b \end{array} \right]^{(k+1)} &:= \beta \cdot \left[ \begin{array}{c} v_w \\ v_b \end{array} \right]^{(k)} + (1 - \beta) \left[ \begin{array}{c} \left( \frac{\partial f(w, b)}{\partial w} \right)^2 \\ \left( \frac{\partial f(w, b)}{\partial b} \right)^2 \end{array} \right]^{(k)} \\ \left[ \begin{array}{c} w \\ b \end{array} \right] &:= \left[ \begin{array}{c} w \\ b \end{array} \right] - \left[ \begin{array}{c} \frac{\eta}{\sqrt{v_w + \epsilon}} * \frac{\partial f(w, b)}{\partial w} \\ \frac{\eta}{\sqrt{v_b + \epsilon}} * \frac{\partial f(w, b)}{\partial b} \end{array} \right] \end{aligned}$$

- RMSProp overcomes the problem of AdaGrad by being less aggressive on the decay.

$$\begin{pmatrix} v_{w0} \\ v_b \end{pmatrix} := \beta \begin{pmatrix} v_{w0} \\ v_b \end{pmatrix} + (1-\beta) \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix}$$

$$k=0 \quad \begin{pmatrix} v_{w0} \\ v_b \end{pmatrix} := (1-\beta) \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix} \quad (0)$$

$$k=1 \quad \begin{pmatrix} v_{w0} \\ v_b \end{pmatrix} = \beta \left[ (1-\beta) \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix} \right] + (1-\beta) \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix} \quad (1)$$

$$k=2 \quad \begin{pmatrix} v_{w0} \\ v_b \end{pmatrix} = \beta \left[ \beta (1-\beta) \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix} \right] + (1-\beta) \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix} \quad (2)$$

$$\boxed{\beta(1-\beta)} \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix} + \boxed{\beta(1-\beta)} \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix} + \boxed{(1-\beta)} \begin{pmatrix} \left(\frac{\partial f}{\partial w}\right)^2 \\ \left(\frac{\partial f}{\partial b}\right)^2 \end{pmatrix}$$

$0.02$        $0.08$        $0.1$

$$\beta = 0.9$$

# RMSProp

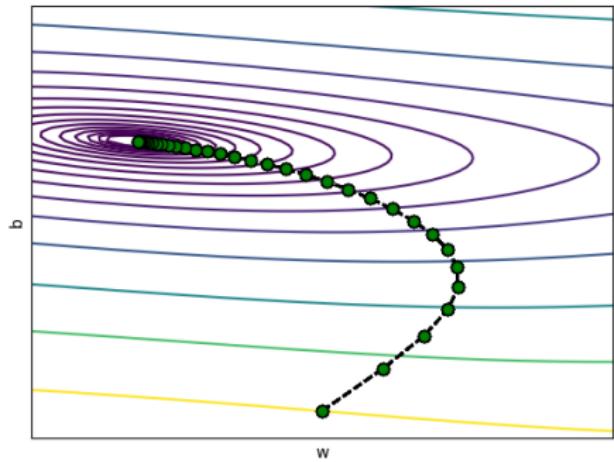


Figure: Epoch: 50  
Total Gradient Evaluations: 50000

# AdaGrad & RMSProp

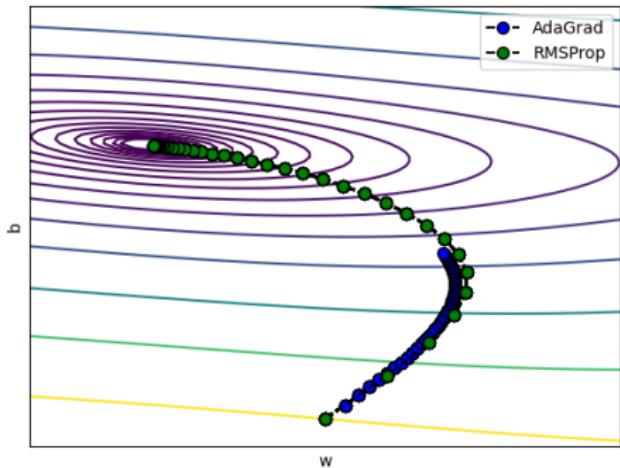


Figure: AdaGrad & RMSProp  
Epoch: 50  
Total Gradient Evaluations: 50000

# Adam (Adaptive Moment Estimation)

- Adam in addition to taking care of denominator decay problem as RMSProp, uses a cumulative history of gradients.
- Update rule:

$$\begin{bmatrix} m_w \\ m_b \end{bmatrix} := \beta_1 \begin{bmatrix} m_w \\ m_b \end{bmatrix} + (1 - \beta_1) \begin{bmatrix} \frac{\partial f(w, b)}{\partial w} \\ \frac{\partial f(w, b)}{\partial b} \end{bmatrix}$$

$w$   $\leftarrow$   $w - \hat{m}_w / \sqrt{\hat{v}_w + \epsilon}$

$$\begin{bmatrix} v_w \\ v_b \end{bmatrix} := \beta_2 \begin{bmatrix} v_w \\ v_b \end{bmatrix} + (1 - \beta_2) \begin{bmatrix} \left( \frac{\partial f(w, b)}{\partial w} \right)^2 \\ \left( \frac{\partial f(w, b)}{\partial b} \right)^2 \end{bmatrix}$$

$$\begin{bmatrix} \hat{m}_w \\ \hat{m}_b \end{bmatrix} := \frac{1}{1 - \beta_1^i} \begin{bmatrix} m_w \\ m_b \end{bmatrix} \quad \begin{bmatrix} \hat{v}_w \\ \hat{v}_b \end{bmatrix} := \frac{1}{1 - \beta_2^i} \begin{bmatrix} v_w \\ v_b \end{bmatrix} \quad \text{for } i^{\text{th}} \text{ iteration}$$

$$\begin{bmatrix} w \\ b \end{bmatrix} := \begin{bmatrix} w \\ b \end{bmatrix} - \begin{bmatrix} \frac{\eta}{\sqrt{\hat{v}_w + \epsilon}} * \hat{m}_w \\ \frac{\eta}{\sqrt{\hat{v}_b + \epsilon}} * \hat{m}_b \end{bmatrix}$$

- Adam takes a cumulative history of gradients which speeds it up and converges faster than AdaGrad and RMSProp.

$A \in \mathbb{R}^{n \times n}$   
 $x \in \mathbb{R}^n$   
 $\lambda \in \mathbb{C}$   
 $x \neq 0$   
 $x \in \mathbb{R}^n$   
 $x = v$   
 $x^T A x = v^T A v$   
 $= v^T \lambda v$   
 $= \lambda v^T v$   
 $= \lambda \|v\|^2$   
 $\lambda \neq 0$   
 $v \neq 0$   
 $v$  is eigenvector for  $A$   
 $A v = \lambda v$   
 $v^T A^T = \lambda v^T$   
 $v^T A^T v = \lambda v^T v$   
 $= \lambda \|v\|^2$   
 $\lambda \neq 0$   
 $v^T v \neq 0$   
 $\lambda > 0$   
 $\lambda$  for a positive definite matrix



125000  
125000

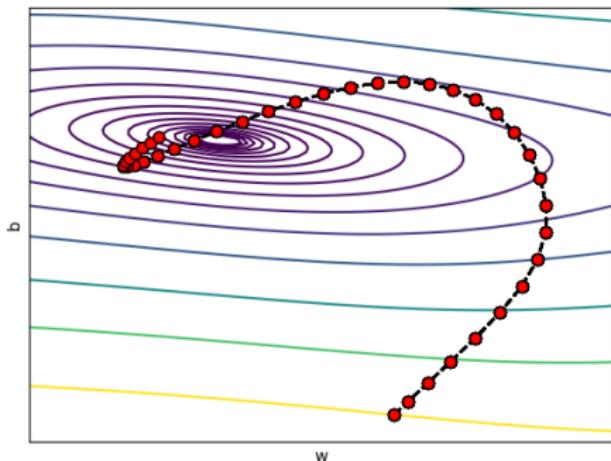


Figure: Epoch 40  
Total Gradient Evaluations: 40000