

# Gradient Descent Solutions to Least-Square Problems

# Objectives

- explain need for iterative algorithms
- derive gradient descent algorithm
- consider impact of step size on convergence
- introduce notion of convex functions

Iterative solution methods play an important role<sup>2</sup>

Features/labels:  $\underline{x}_i, d_i, i=1, 2, \dots, N$  N samples

Classifier or model error:  $e^2 = \sum_{i=1}^N (\underline{x}_i^T \underline{w} - d_i)^2$

$$\underline{A} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} \quad \underline{d} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_N \end{bmatrix} \quad e^2 = \|\underline{A} \underline{w} - \underline{d}\|_2^2$$

Regularized least squares:  $\arg \min_{\underline{w}} \|\underline{A} \underline{w} - \underline{d}\|_2^2 + \lambda r(\underline{w})$

1. Computational cost  $(\underline{A}^T \underline{A})^{-1}$
2. Closed form solution maybe unavailable
3. Adapt  $\underline{w}$  to new features/labels

The reason for

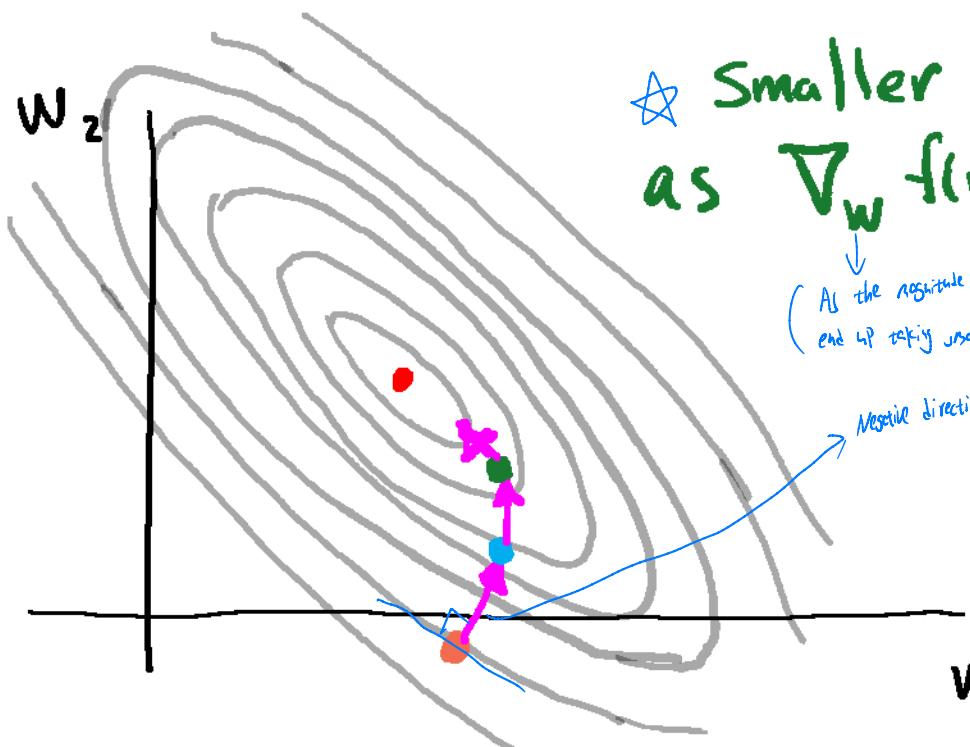
develop iterative approach

# Gradient descent finds the minimum

$$f(\underline{w}) = \|\underline{A}\underline{w} - \underline{d}\|_2^2$$

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \tau' \nabla_{\underline{w}} f(\underline{w})$$

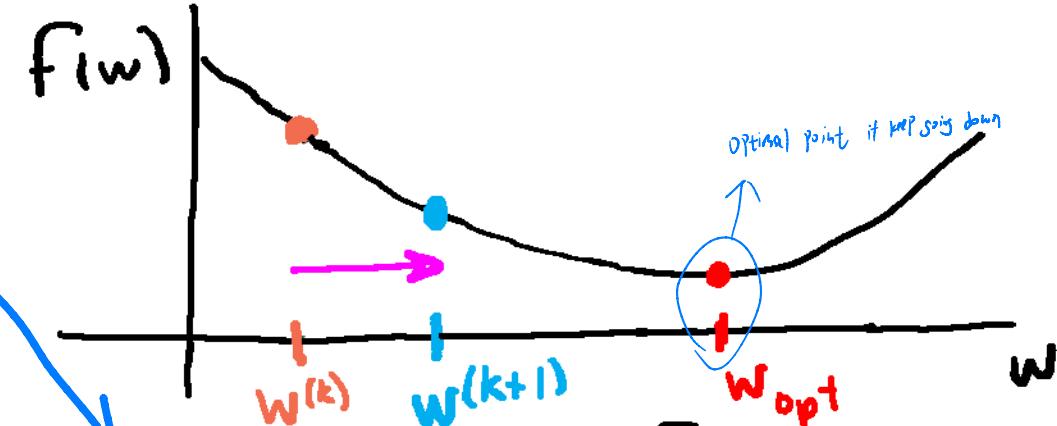
New weight  
( $\tau' > 0$ )      Previous weight  
Step size  $\tau'$       gradient



★ Smaller steps as  $\nabla_{\underline{w}} f(\underline{w}) \downarrow$

(As the magnitude of gradient decreases, we end up taking smaller steps since  $\tau' \propto f(\underline{w})$ )

Negative direction to minimize the  $f(\underline{w})$



$$f(\underline{w}) = (\underline{A}\underline{w} - \underline{d})^T (\underline{A}\underline{w} - \underline{d})$$

$$= \underline{w}^T \underline{A}^T \underline{A} \underline{w} - 2 \underline{w}^T \underline{A}^T \underline{d} + \underline{d}^T \underline{d}$$

$$\nabla_{\underline{w}} f(\underline{w}) = 2 \underline{A}^T \underline{A} \underline{w} - 2 \underline{A}^T \underline{d}$$

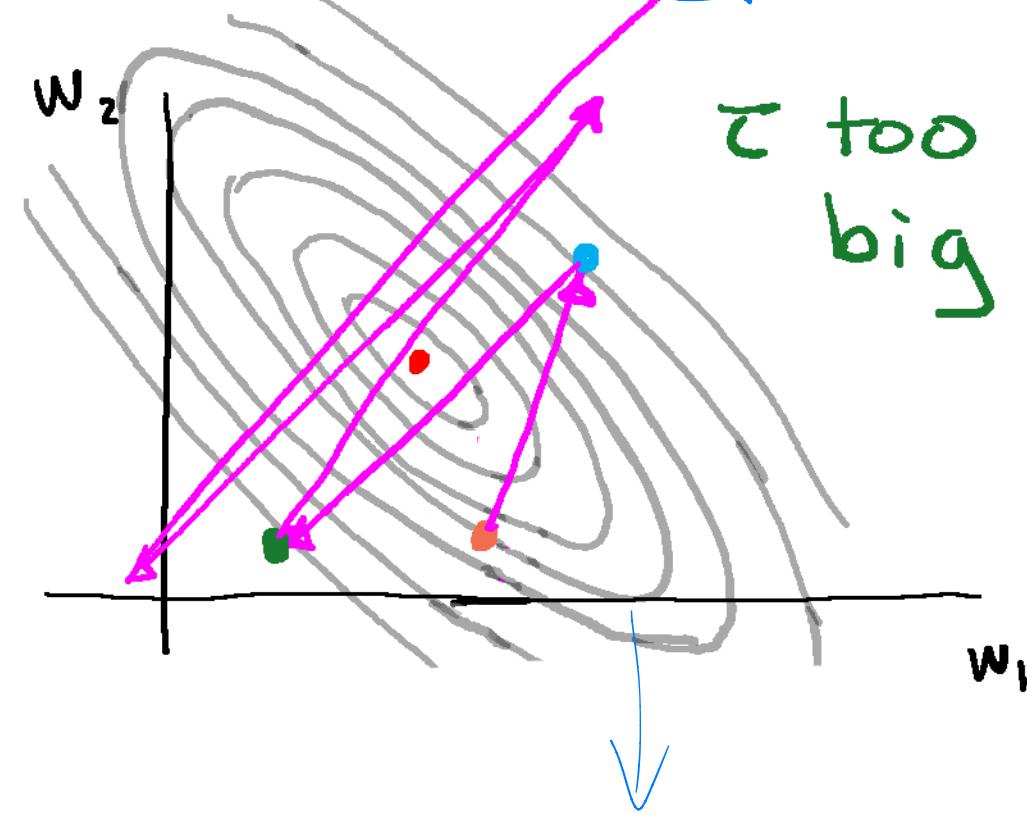
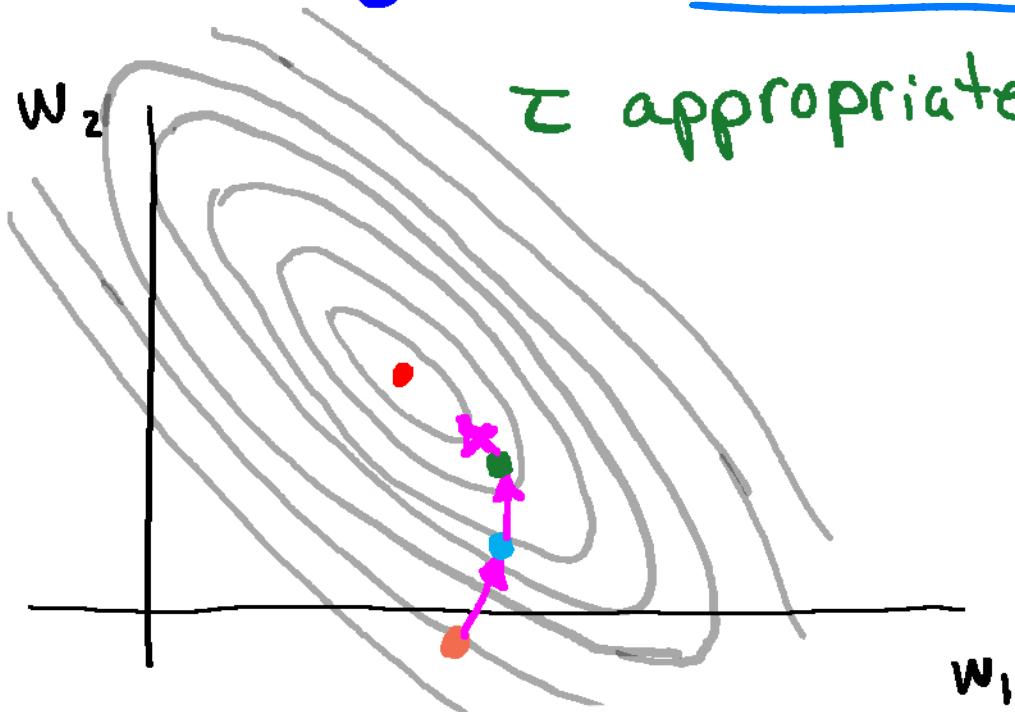
$$= 2 \underline{A}^T (\underline{A} \underline{w} - \underline{d})$$

$$\star \underline{w}^{(k+1)} = \underline{w}^{(k)} - \tau \underline{A}^T (\underline{A} \underline{w}^{(k)} - \underline{d})$$

(This is absorbed with  $\tau$ )

(Landweber iteration)

# Convergence behavior depends on $\tau$



\*  $\tau$  too small: slow convergence  
 $\tau$  too big: no convergence  
unstable!

Require  $0 < \tau < 2/\|\underline{A}\|_{\text{op}}^2$  for convergence 5

$$\text{Recall } \|\underline{A}\|_{\text{op}} = \|\underline{A}\|_2 = \sigma_{\max}(\underline{A})$$

$\|\$   
(Largest singular value of  $\underline{A}$ )

↓(Put bounds on step size  $\tau$ )

★(Note: For matrix  $\underline{A}$ ,  
largest singular value  $\sigma_1 = \sqrt{\lambda_{\max}}$ )  
 $\|\underline{A}^T \underline{A}\|$  is the estimate of square root

Convergence:  $f(\underline{w}^{(k+1)}) < f(\underline{w}^{(k)})$  cost decreases as  $k$  increases

$$\|\underline{A}\underline{w}^{(k+1)} - \underline{d}\|_2^2 < \|\underline{A}\underline{w}^{(k)} - \underline{d}\|_2^2$$

↓(Definition of convergence, the cost needs to be smaller each iteration. If it goes up it does not go down)

↓(Ensures you are actually going downhill)  $(k+1)$  round's cost <  $k$  round's cost

Notes - guaranteed convergence for

$$0 < \tau < 2/\|\underline{A}\|_{\text{op}}^2$$

$$\underline{w}^{(0)} = \underline{0}, \quad \underline{w}^{(k+1)} = \underline{w}^{(k)} - \tau \underline{A}^T (\underline{A}\underline{w}^{(k)} - \underline{d}) \xrightarrow{k} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

↓  
The iteration algorithm

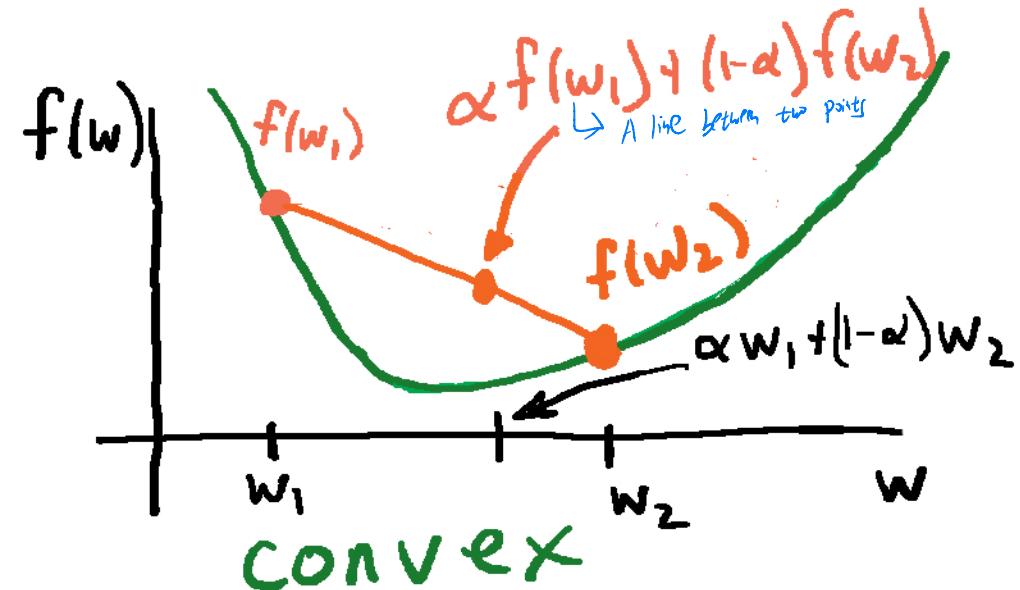
↓  
After  $k$  rounds

↓  
(converge to the least squares solution)

↓  
(Let's initialize the weights to be 0 at first for convenience)

# Gradient descent is effective for convex cost functions

(★ If non-convex, GD will not guarantee to find the global minimum) 6



① Line segment test

$$f(\alpha w_1 + (1-\alpha)w_2) \leq \alpha f(w_1) + (1-\alpha) f(w_2);$$

(The cost function needs to be below the line segment for any two points  $\rightarrow$  convex)

$$0 < \alpha < 1, \text{ all } w_1, w_2$$

Multidimensional case

\* Quick way for testing  $H \geq 0$ :

$H$  is symmetric (known)

PSD: All eigenvalues  $\lambda \geq 0$

$$H(\underline{w}) \geq 0 \quad \rightarrow \text{Needs to be semi-definite}$$



② Second derivative test

$$\frac{d^2}{dw^2} f(w) \geq 0 \quad \begin{array}{l} (\text{if function also does not have inflection points, then convex}) \\ \downarrow \end{array}$$

$$[H(\underline{w})]_{ij} = \frac{\partial^2}{\partial w_i \partial w_j} f(\underline{w})$$

Hessian matrix

★ (Second derivative respect to few cost functions in various coordinate directions)

**Copyright 2019  
Barry Van Veen**

# Gradient Descent for Solving Least-Squares Problems

## Proof: Bounds on Step Size for Guaranteed Convergence

©Barry Van Veen 2019

Gradient descent minimizes the cost function

$$f(\mathbf{w}) = \|\mathbf{A}\mathbf{w} - \mathbf{d}\|_2^2$$

using the iterative algorithm

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \tau \mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}), \quad k = 0, 1, 2, 3, \dots \quad (1)$$

where  $\tau > 0$  so that we modify the current iterate in the negative gradient direction. Often this algorithm is initialized with  $\mathbf{w}^{(0)} = \mathbf{0}$ . The initialization does not affect the convergence behavior because  $f(\mathbf{w})$  is convex.

The iteration is guaranteed to converge to the minimum of the cost function if the squared error decreases with each iteration, that is, if

$$f(\mathbf{w}^{(k+1)}) = \|\mathbf{A}\mathbf{w}^{(k+1)} - \mathbf{d}\|_2^2 < f(\mathbf{w}^{(k)}) = \|\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}\|_2^2$$

Begin by substituting  $\mathbf{w}^{(k)} - \tau \mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d})$  for  $\mathbf{w}^{(k+1)}$  in  $f(\mathbf{w}^{(k+1)})$  to write

$$f(\mathbf{w}^{(k+1)}) = \|\mathbf{A}(\mathbf{w}^{(k)} - \tau \mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d})) - \mathbf{d}\|_2^2 \quad (2)$$

$$= \|(\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}) - \tau (\mathbf{A}\mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}))\|_2^2 \quad (3)$$

Now let  $\mathbf{c} = \mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}$  and  $\mathbf{e} = \tau (\mathbf{A}\mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}))$  be the first and second terms in parentheses so  $f(\mathbf{w}^{(k+1)}) = \|\mathbf{c} - \mathbf{e}\|_2^2 = (\mathbf{c} - \mathbf{e})^T (\mathbf{c} - \mathbf{e})$ . Expand the product to write  $f(\mathbf{w}^{(k+1)}) = \|\mathbf{c}\|_2^2 + \|\mathbf{e}\|_2^2 - 2\mathbf{e}^T \mathbf{c}$ . Substituting for  $\mathbf{c}$  and  $\mathbf{e}$  we thus obtain

$$f(\mathbf{w}^{(k+1)}) = \|\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}\|_2^2 + \tau^2 \|\mathbf{A}\mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d})\|_2^2 - 2\tau ((\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d})^T \mathbf{A}\mathbf{A}^T) (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}) \quad (4)$$

$$= f(\mathbf{w}^{(k)}) + \tau^2 \|\mathbf{A}(\mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}))\|_2^2 - 2\tau ((\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d})^T \mathbf{A}) (\mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d})) \quad (5)$$

Define  $\mathbf{v} = \mathbf{A}^T (\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d})$  to simplify the expression and rewrite Eq. 5 as

$$f(\mathbf{w}^{(k+1)}) = f(\mathbf{w}^{(k)}) + \tau^2 \|\mathbf{A}\mathbf{v}\|_2^2 - 2\tau \mathbf{v}^T \mathbf{v}$$

Note that  $\mathbf{v}$  does not depend on  $\tau$ . Thus, to prove  $f(\mathbf{w}^{(k+1)}) < f(\mathbf{w}^{(k)})$ , we must find the condition for which

$$q(\tau) = \tau^2 \|\mathbf{A}\mathbf{v}\|_2^2 - 2\tau \mathbf{v}^T \mathbf{v}$$

is less than zero.

Recall the operator norm of a matrix  $\mathbf{X}$  satisfies  $\max_{\mathbf{g}} \|\mathbf{X}\mathbf{g}\|_2 \leq \|\mathbf{X}\|_{op} \|\mathbf{g}\|_2$ , so the first term in  $q(\tau)$  may be upper bounded as

$$\tau^2 \|\mathbf{A}\mathbf{v}\|_2^2 \leq \tau^2 \|\mathbf{A}\|_{op}^2 \|\mathbf{v}\|_2^2$$

We may rewrite the second term in  $q(\tau)$  as

$$-2\tau \mathbf{v}^T \mathbf{v} = -2\tau \|\mathbf{v}\|_2^2$$

Hence, we obtain an upper bound on  $q(\tau)$

$$q(\tau) \leq \tau^2 \|\mathbf{A}\|_{op}^2 \|\mathbf{v}\|_2^2 - 2\tau \|\mathbf{v}\|_2^2$$

Factoring out the common terms we write

$$q(\tau) \leq (\tau \|\mathbf{A}\|_{op}^2 - 2) \tau \|\mathbf{v}\|_2^2$$

The second term in  $q(\tau)$ ,  $\tau \|\mathbf{v}\|_2^2$ , is positive provided  $\mathbf{v} \neq \mathbf{0}$ , so we obtain  $q(\tau) < 0$  by requiring

$$(\tau \|\mathbf{A}\|_{op}^2 - 2) < 0$$

which indicates  $\tau$  must satisfy

$$\tau < \frac{2}{\|\mathbf{A}\|_{op}^2}$$

Note that  $\mathbf{v} = \mathbf{0}$  implies  $\mathbf{A}^T(\mathbf{A}\mathbf{w}^{(k)} - \mathbf{d}) = \mathbf{0}$ , or  $\mathbf{A}^T \mathbf{A} \mathbf{w}^{(k)} = \mathbf{A}^T \mathbf{d}$ , or  $\mathbf{w}^{(k)} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{d}$ . Thus, if  $\mathbf{v} = \mathbf{0}$ , then the iteration has converged to the minimum of the squared error and the update term in Eq. 1 is zero.

Hence, the gradient descent algorithm will converge to the minimum of the squared error cost function provided the step-size  $\tau$  satisfies  $\tau < \frac{2}{\|\mathbf{A}\|_{op}^2}$ .