

Introduction to Machine Learning - CS725

Instructor: Prof. Ganesh Ramakrishnan

Lecture 08 - Optimization Foundations Applied to Regression  
Formulations

# Building on questions on Least Squares Linear Regression

- ① Is there a probabilistic interpretation?
  - Gaussian Error, Maximum Likelihood Estimate
- ② Addressing overfitting
  - Bayesian and Maximum A posteriori Estimates, Regularization, Support Vector Regression
- ③ How to minimize the resultant and more complex error functions?
  - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality

Lagrange fn, KKT conditions, Dual formulation

- 1-norm Error, and  $L_2$  regularized:

- $\min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*)$   
s.t.  $\forall i,$   
 $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i,$   
 $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*,$   
 $\xi_i, \xi_i^* \geq 0$

} 4m constraints

- 2-norm Error, and  $L_2$  regularized:

- $\min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i^2 + \xi_i^{*2})$   
s.t.  $\forall i,$   
 $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i,$   
 $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$

} 2m constraints

- Here, the constraints  $\xi_i, \xi_i^* \geq 0$  are not necessary

# Need for Optimization so far (Equivalences we need to show)

- Unconstrained (**Penalized**) Optimization:

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 + \underbrace{\Omega(\mathbf{w})}_{\lambda}$$

- **Constrained** Optimization 1:

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2$$

such that  $\Omega(\mathbf{w}) \leq \theta$

- **Constrained** Optimization 2 ( $t = 1$  or  $2$ ):

$$\arg \min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i^t + \xi_i^{*t})$$

s.t.  $\forall i, y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i; b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$

- **Equivalence:**  $\lambda$  (**Penalized**)  $\equiv \theta$  (**Constrained**)
- **Duality:** Dual of Support Vector Regression

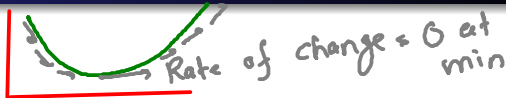
You just  
cannot  
avoid  
optimization  
in ML!

# Solving Unconstrained Minimization Problem

Derivative increases!

High school calculus

$$\frac{df}{dx} = 0$$



- Intuitively: Minimize by setting derivative (gradient) to 0 and hoping to find closed form solution.
- When is such a solution a global minimum?
- For most optimization problems, finding closed form solutions is difficult. Even for linear regression (for which closed form solution exists), are there alternative methods?
  - Eg: Consider,  $\mathbf{y} = \phi \mathbf{w}$ , where  $\phi$  is a matrix with full column rank, the least squares solution,  $\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$ . Now, imagine that  $\phi$  is a very large matrix. with say, 100,000 columns and 1,000,000 rows. Computation of closed form solution might be challenging.

$$\frac{d^2f}{dx^2} \geq 0 \text{ (Necessary)} \quad \frac{d^2f}{dx^2} > 0 \text{ (sufficient)}$$

- How about iterative methods?

$$\frac{df}{dx} \text{ increases} \Rightarrow \frac{d^2f}{dx^2} \geq 0$$

$$\frac{df}{dx} \text{ strictly increases} \Rightarrow \frac{d^2f}{dx^2} > 0$$

# Foundations: Level curves and surfaces

$\{x | f(x) \leq c\} \rightarrow$  Sublevel set  
 $\rightarrow$  set of pts encapsulated by level curve

- A level curve of a function  $f(x)$  is defined as a curve along which the value of the function remains unchanged while we change the value of its argument  $x$ .
- Formally we can define a level curve as :

$$L_c(f) = \{x | f(x) = c\} \quad (1)$$

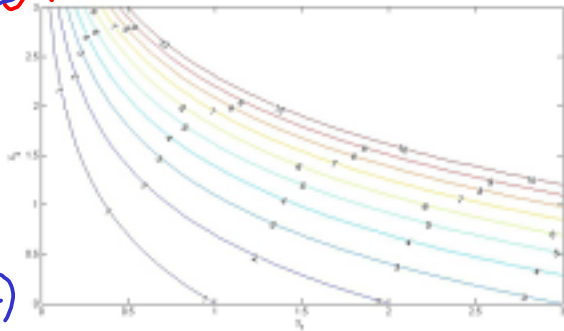
where  $c$  is a constant.

# Foundations: Level curves and surfaces

- Example of different level curves for a single function

Q: What does gap between consecutive level curves indicate?

Region of fast rate of change (quantify using Directional derivative)

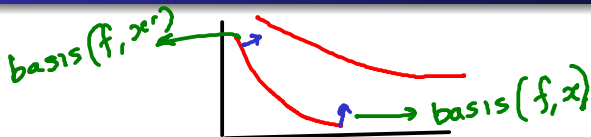


Ans: Larger the gap, smaller is rate of change!

} Region of slow rate of change

Figure 1: 10 level curves for the function  $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 e^{\mathbf{x}_2}$  (Figure 4.12 from <https://www.cse.iitb.ac.in/~CS725/notes/classNotes/BasicsOfConvexOptimization.pdf>)

# Foundations: Directional Derivatives



- Directional derivative: Rate at which the function changes at a given point  $\mathbf{x}$  in a given direction  $\mathbf{v}$
- The *directional derivative* of a function  $f$  in the direction of a unit vector  $\mathbf{v}$  at a point  $\mathbf{x}$  can be defined as :

$$D_{\mathbf{v}}(f, \mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h} \quad (2)$$

*measured at  $\mathbf{x}$*       *measured along  $\mathbf{v}$*

(Claim)

$$\text{Any } D_{\mathbf{v}}(f, \mathbf{x}) = \sum_i v_i [\text{basis}(f, \mathbf{x})]_i = \mathbf{v}^T \underbrace{\nabla f(\mathbf{x})}_{\text{Gradient}} \quad (3)$$

*s.t.  $\|\mathbf{v}\|_2 = 1$*



# Foundations: Gradient Vector

- The gradient vector of a function  $f$  at a point  $\mathbf{x}$  is defined as:

$$\nabla f_{\mathbf{x}^*} = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n$$

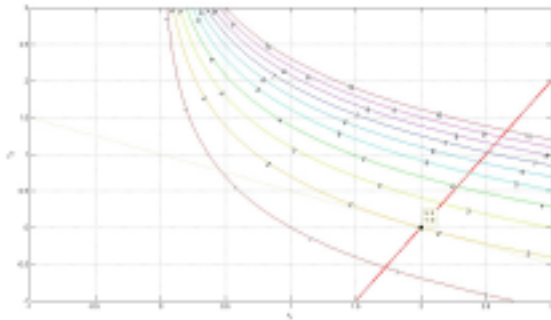
$$\begin{aligned} \max_{\mathbf{v}} D_{\mathbf{v}}(f, \mathbf{x}) \\ = \|\nabla f(\mathbf{x})\| \end{aligned} \quad (4)$$

- Magnitude (euclidean norm) of gradient vector at any point indicates maximum value of directional derivative at that point
- Direction of gradient vector indicates direction of this maximal directional derivative at that point.

$$\arg \max_{\mathbf{v}} D_{\mathbf{v}}(f, \mathbf{x}) = \frac{1}{\|\nabla f(\mathbf{x})\|} \nabla f(\mathbf{x})$$

# Foundations: Gradient Vector

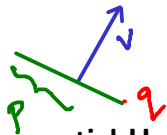
- The figure below illustrates the gradient vector for the same level curves



**Figure 2:** The level curves along with the gradient vector at  $(2, 0)$ . Note that the gradient vector is perpendicular to the level curve  $x_1 e^{x_2} = 2$  at  $(2, 0)$

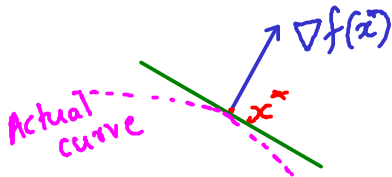
# Hyperplanes

- A hyperplane in an  $n$ -dimensional Euclidean space is a flat,  $n-1$  dimensional subset of that space that divides the space into two disjoint half-spaces.
- Technically, a hyperplane is a set of points whose direction *w.r.t.* a point  $\mathbf{q}$  is orthogonal to a vector  $\mathbf{v}$ :



$$H_{\mathbf{v},\mathbf{q}} = \left\{ \mathbf{p} \mid (\mathbf{p} - \mathbf{q})^T \mathbf{v} = 0 \right\} \quad (5)$$

- **Tangential Hyperplane:** Plane orthogonal to the gradient vector at  $\mathbf{x}^*$ .



# Hyperplanes

- A hyperplane in an  $n$ -dimensional Euclidean space is a flat,  $n-1$  dimensional subset of that space that divides the space into two disjoint half-spaces.
- Technically, a hyperplane is a set of points whose direction *w.r.t.* a point  $\mathbf{q}$  is orthogonal to a vector  $\mathbf{v}$ :

$$H_{\mathbf{v},\mathbf{q}} = \left\{ \mathbf{p} \mid (\mathbf{p} - \mathbf{q})^T \mathbf{v} = 0 \right\} \quad (5)$$

- **Tangential Hyperplane:** Plane orthogonal to the gradient vector at  $\mathbf{x}^*$ .

$$TH_{\mathbf{x}^*} = \left\{ \mathbf{p} \mid (\mathbf{p} - \mathbf{x}^*)^T \nabla f(\mathbf{x}^*) = 0 \right\} \quad (6)$$

## (Ridge Regression)

We recall that the problem was to find  $\mathbf{w}$  such that

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2 \quad (7)$$

$$= \arg \min_{\mathbf{w}} (\underbrace{\mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{y} + \mathbf{y}^T \mathbf{y}}_{L_P(\mathbf{w}, \mathcal{D})} + \lambda \|\mathbf{w}\|^2) \quad (8)$$

# Foundations: Gradient Vector

- Magnitude (euclidean norm) of gradient vector at any point indicates maximum value of directional derivative at that point
- Thus, at the point of minimum of a differentiable minimization objective (such as least squares for regression), ....

If  $w^*$  is soln to Ridge Regression, it must minimize Ridge Regression Loss  $\Rightarrow \nabla_w L_R(w^*; D) = 0$

$L_R(w^*; D)$

# Foundations: Necessary condition 1

- If  $\nabla f(\mathbf{w}^*)$  is defined &  $\mathbf{w}^*$  is local minimum/maximum, then  $\nabla f(\mathbf{w}^*) = 0$  (A necessary condition) (Cite : Theorem 60) of

CS725/notes/classNotes/BasicsofConvexOptimization.pdf

- Given that

Think of  $\phi \in \mathbb{R}$  &  $w \in \mathbb{R} \Rightarrow f(w) = \phi^2 w^2 - 2w\phi y - y^2 + \lambda w^2$   
 $f'(w) = 2\phi^2 w - 2\phi y + 2\lambda w$

$$f(\mathbf{w}) = \underline{(\mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{y} - \mathbf{y}^T \mathbf{y} + \lambda \|\mathbf{w}\|^2)}$$

$\Rightarrow$  .....

$$\nabla f(w) = 2\phi^2 w - 2\phi y + 2\lambda w$$

- We would have

Note

$$\mathbf{y}^T \mathbf{A} \mathbf{x} \in \mathbb{R} \\ \Rightarrow (\mathbf{y}^T \mathbf{A} \mathbf{x})^T = (\mathbf{y}^T \mathbf{A} \mathbf{x})$$

$$\nabla f(\mathbf{w}^*) = 0 \Rightarrow$$

$$\Rightarrow \mathbf{w}^* = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

$$\Rightarrow \dots\dots\dots$$

Tricks:  $\nabla_x (x^T \mathbf{A} x) = 2\mathbf{A}x$

$$\nabla_x (x^T \mathbf{A}) = \mathbf{A}$$

Thus:  $\nabla_x (\mathbf{y}^T \mathbf{A} \mathbf{x}) = \nabla_x (x^T \mathbf{A}^T \mathbf{y}) = \mathbf{A}^T \mathbf{y}$

# Foundations: Necessary condition 1

- If  $\nabla f(\mathbf{w}^*)$  is defined &  $\mathbf{w}^*$  is local minimum/maximum, then  $\nabla f(\mathbf{w}^*) = 0$  (A necessary condition) (Cite : Theorem 60)

CS725/notes/classNotes/BasicsOfConvexOptimization.pdf

- Given that

$$f(\mathbf{w}) = (\mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{y} - \mathbf{y}^T \mathbf{y} + \lambda \|\mathbf{w}\|^2) \quad (9)$$

$$\implies \nabla f(\mathbf{w}) = 2\Phi^T \Phi \mathbf{w} - 2\Phi^T \mathbf{y} + 2\lambda \mathbf{w} \quad (10)$$

- We would have

$$\nabla f(\mathbf{w}^*) = 0 \quad (11)$$

$$\implies 2(\Phi^T \Phi + \lambda I) \mathbf{w}^* - 2\Phi^T \mathbf{y} = 0 \quad (12)$$

$$\implies \mathbf{w}^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{y} \quad (13)$$





# Foundations: Necessary Condition 2

- Is  $\nabla^2 f(\mathbf{w}^*)$  positive definite?

i.e.  $\forall \mathbf{x} \neq 0$ , is  $\mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} > 0$ ? (A sufficient condition for local minimum)

(Note: Any positive definite matrix is also positive semi-definite)

(Cite: Section 3.12 & 3.12.1)<sup>1</sup>

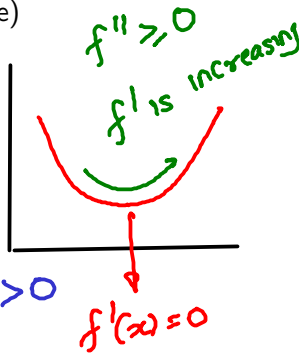
Hessian Matrix

$$\nabla^2 f(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 \partial w_2} & \dots & \frac{\partial^2 f}{\partial w_1 \partial w_p} \\ \frac{\partial^2 f}{\partial w_2 \partial w_1} & \frac{\partial^2 f}{\partial w_2^2} & \dots & \frac{\partial^2 f}{\partial w_2 \partial w_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_p \partial w_1} & \frac{\partial^2 f}{\partial w_p \partial w_2} & \dots & \frac{\partial^2 f}{\partial w_p^2} \end{bmatrix}$$

Trick:  $\nabla^2 f(\mathbf{x}^T \mathbf{A} \mathbf{x}) = 2\mathbf{A}$

$$\nabla^2 f(\mathbf{w}^*) = 2\Phi^T \Phi + 2\lambda \mathbf{I}$$

$\Phi^T \Phi$  is p.s.d and with  $\lambda > 0$ ,  $(\Phi^T \Phi + \lambda \mathbf{I}) > 0$



- And if  $\Phi$  has full column rank,

$$\therefore \text{If } \mathbf{x} \neq 0, \quad \mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} > 0$$

<sup>1</sup>CS725/notes/classNotes/LinearAlgebra.pdf

## Foundations: Necessary Condition 2

- Is  $\nabla^2 f(\mathbf{w}^*)$  positive definite ?

i.e.  $\forall \mathbf{x} \neq 0$ , is  $\mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} > 0$ ? (A sufficient condition for local minimum)

(Any positive definite matrix is also positive semi-definite)

(Cite : Section 3.12 & 3.12.1)<sup>2</sup>

$$\nabla^2 f(\mathbf{w}^*) = 2\Phi^T \Phi + 2\lambda I \quad (14)$$

$$\implies \mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} = 2\mathbf{x}^T (\Phi^T \Phi + \lambda I) \mathbf{x} \quad (15)$$

$$= 2 \left( (\Phi + \sqrt{\lambda} I) \mathbf{x} \right)^T \Phi \mathbf{x} \quad (16)$$

$$= 2 \left\| (\Phi + \sqrt{\lambda} I) \mathbf{x} \right\|^2 \geq 0 \quad (17)$$

- And with  $\lambda = 0$ , if  $\Phi$  has full column rank ,

$$\Phi \mathbf{x} = 0 \quad \text{iff} \quad \mathbf{x} = 0 \quad (18)$$

$\therefore$  If  $\mathbf{x} \neq 0$ ,  $\mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} > 0$

## Example of linearly correlated features

- Example where  $\Phi$  doesn't have a full column rank,

$$\Phi = \begin{bmatrix} x_1 & x_1^2 & x_1^2 & x_1^3 \\ x_2 & x_2^2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ x_n & x_n^2 & x_n^2 & x_n^3 \end{bmatrix}$$

→ Not full column rank

$$\Phi^T \Phi \geq 0$$

(19)

- This is the simplest form of linear correlation of features, and it is not at all desirable.
- Effect of a nonzero  $\lambda$  with such  $\Phi$  is that

$$\Phi^T \Phi + \lambda I > 0$$

← yet

## Example of linearly correlated features

- Example where  $\Phi$  doesn't have a full column rank,

$$\Phi = \begin{bmatrix} x_1 & x_1^2 & x_1^2 & x_1^3 \\ x_2 & x_2^2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ x_n & x_n^2 & x_n^2 & x_n^3 \end{bmatrix} \quad (19)$$

- This is the simplest form of linear correlation of features, and it is not at all desirable.
- Effect of a nonzero  $\lambda$  with such  $\Phi$  is that it tends to make the Hessian more positive definite

# Do Closed-form solutions Always Exist?

- Linear regression and Ridge regression both have closed-form solutions

- For linear regression,

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T y$$

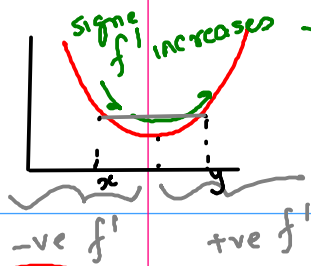
- For ridge regression,

$$w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

(for linear regression,  $\lambda = 0$ )

- What about optimizing the formulations (constrained/penalized) of Lasso ( $L_1$  norm)? And support-based penalty ( $L_0$  norm)? *Also requires tools of Optimization/duality*

$\nabla (w^T \Phi^T \Phi w + 2w^T \Phi y + y^T y + \lambda \|w\|_1)$  does not exist!



$\nabla^2 f > 0 \leftrightarrow f$  is strictly convex

$\nabla^2 f \geq 0 \leftrightarrow f$  is convex

Cup shaped

$f$  is always below the gray line segment

If  $\nabla^2 f(x) \geq 0 \quad \forall x \in D$  then  $f$  is

|||

strictly convex in  $D$

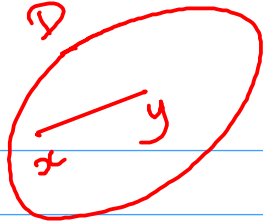
If  $x, y \in D, \theta x + (1-\theta)y \in D$

If  $f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$

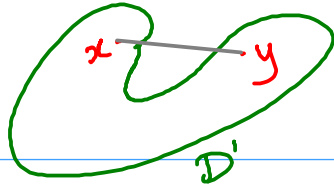
for  $\theta \in [0,1]$

is a pt on line segment  $[x,y]$

$D$  is also convex set

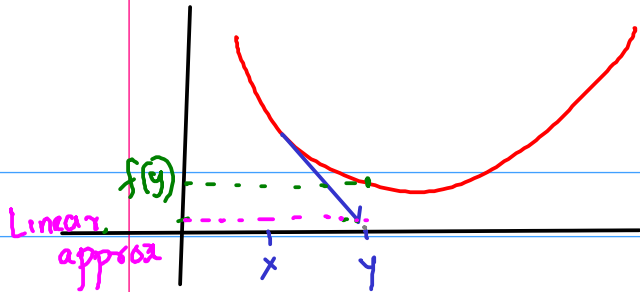


Convex  $D$



Line segment  $[x, y]$   
is NOT completely  
contained in  $D'$

$\Rightarrow \underline{D'}$  is NOT convex



Observation for convex fn:

(Linear approx at  $x$   
evaluated at  $y$ )  $\leq$   $f(y)$

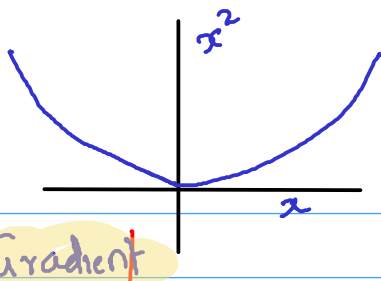
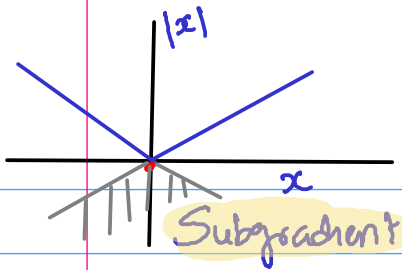
strict convexity

convex

$$f(x) + \nabla^T f(x)(y-x)$$

Accounts for slope





③  $\nabla$  is direction that gives lower bnd linear approx

①  $\nabla$  is direction of max rate of change

②  $\perp$  to tangent hyperplane

Think of a vector that behaves like gradient vec:

# Gradient Descent Algorithm

**Find** starting point  $\mathbf{w}^{(0)} \in \mathcal{D}$

- $\Delta \mathbf{w}^k = -\nabla \varepsilon(\mathbf{w}^{(k)})$
- Choose a step size  $t^{(k)} > 0$  using exact or backtracking ray search.
- Obtain  $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mathbf{t}^{(k)} \Delta \mathbf{w}^{(k)}$ .
- Set  $k = k + 1$ . **until** stopping criterion (such as  $\|\nabla \varepsilon(\mathbf{w}^{(k+1)})\| \leq \epsilon$ ) is satisfied