

10-301/601: Introduction to Machine Learning

Lecture 8 – Optimization for Machine Learning

Henry Chai

5/31/23

Front Matter

- Announcements:
 - PA2 released 5/25, due 6/01 at 11:59 PM
 - No new programming assignment this week!
- Recommended Readings:
 - None

Recall: Minimizing the Squared Error

$$\ell_{\mathcal{D}}(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y^{(n)})^2 = \sum_{n=1}^N (\mathbf{x}^{(n)T} \mathbf{w} - y^{(n)})^2$$

$$= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \text{ where } \|\mathbf{z}\|_2 = \sqrt{\sum_{d=1}^D z_d^2} = \sqrt{\mathbf{z}^T \mathbf{z}}$$

$$= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$= (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

$$\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\hat{\mathbf{w}}) = (2\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} - 2\mathbf{X}^T \mathbf{y}) = 0$$

$$\rightarrow \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$$

$$\rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

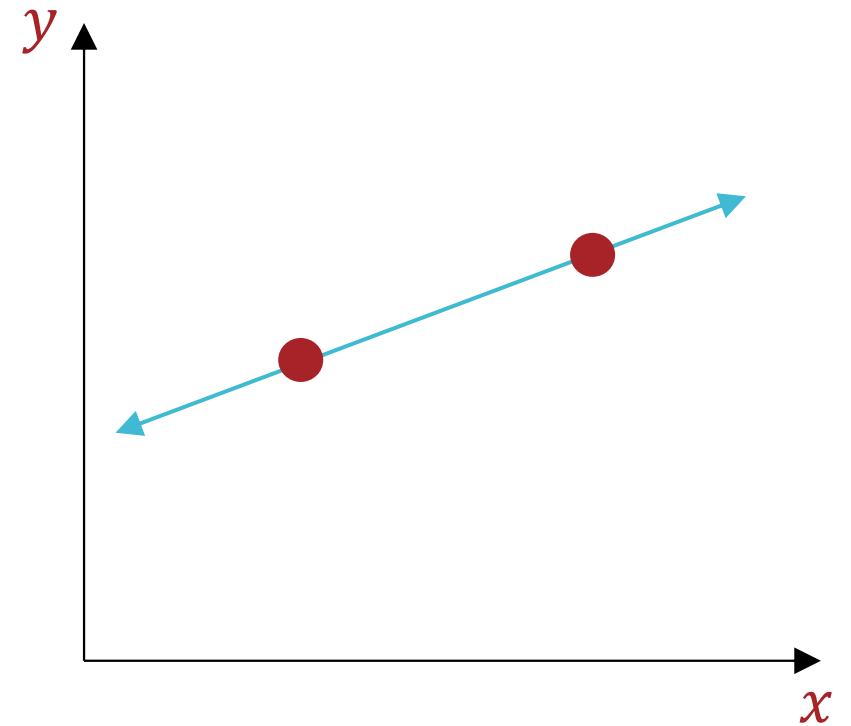
Recall: Closed Form Solution

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

1. Is $\mathbf{X}^T \mathbf{X}$ invertible?
 - When $N \gg D + 1$, $\mathbf{X}^T \mathbf{X}$ is (almost always) full rank and therefore, invertible!
 - If $\mathbf{X}^T \mathbf{X}$ is not invertible (occurs when one of the features is a linear combination of the others), what does that imply about our problem?
2. If so, how computationally expensive is inverting $\mathbf{X}^T \mathbf{X}$?
 - $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{D+1 \times D+1}$ so inverting $\mathbf{X}^T \mathbf{X}$ takes $O(D^3)$ time...
 - Computing $\mathbf{X}^T \mathbf{X}$ takes $O(ND^2)$ time
 - What alternative optimization method(s) can we use to minimize the mean squared error?

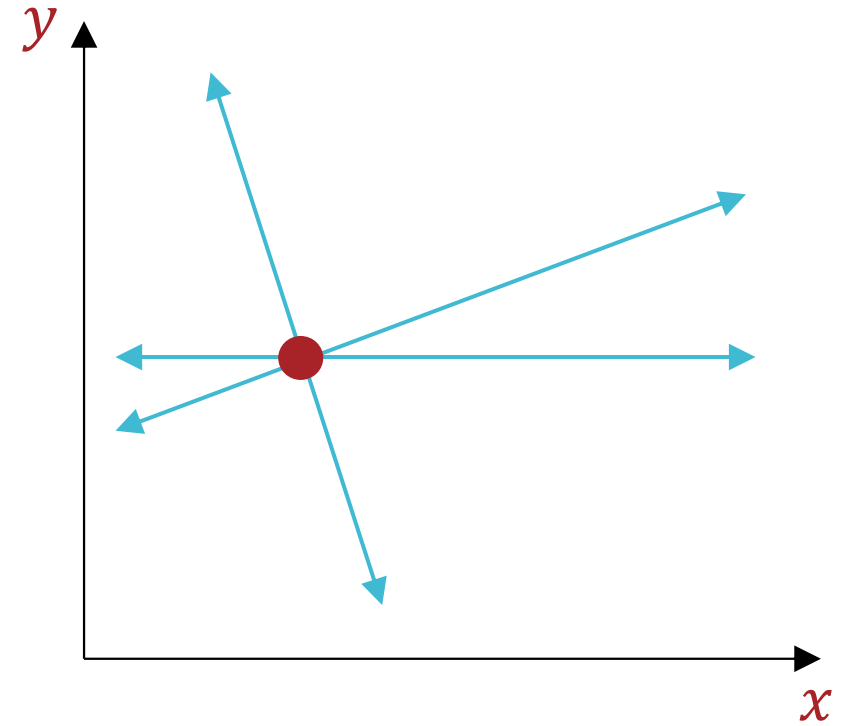
Linear Regression: Uniqueness

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters θ) are there for the given dataset?



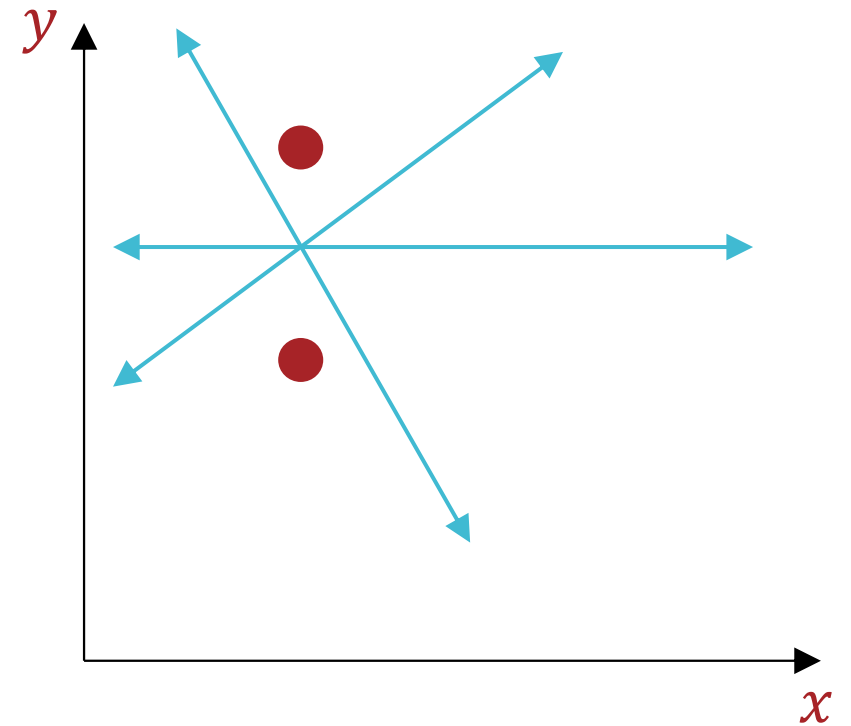
Linear Regression: Uniqueness

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters θ) are there for the given dataset?



Linear Regression: Uniqueness

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters θ) are there for the given dataset?

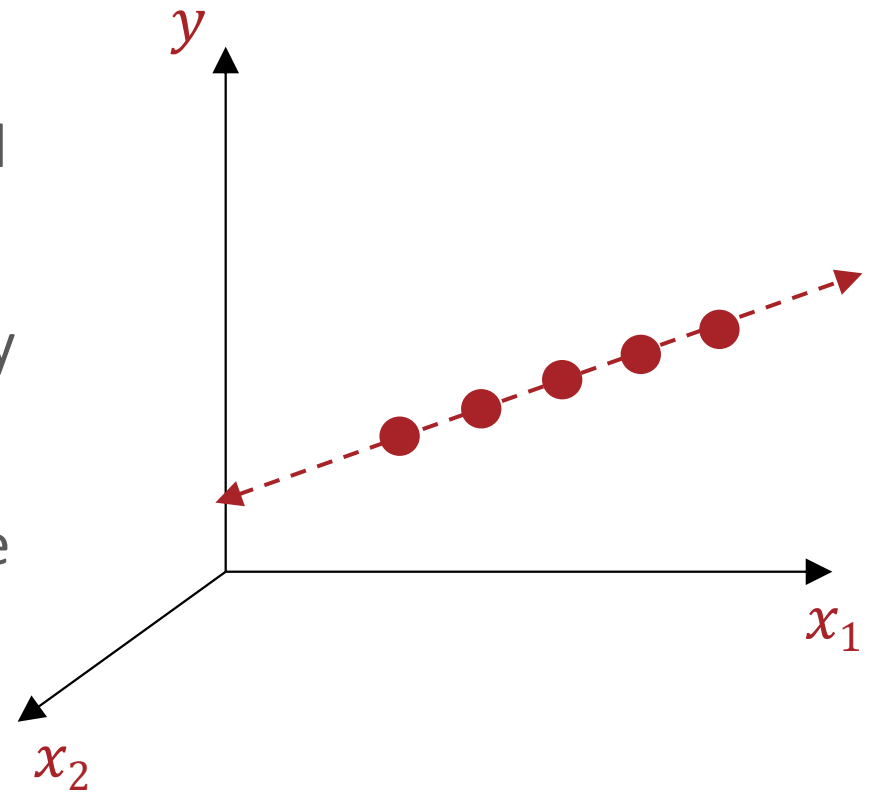


How many solutions optimal solutions are there for the given dataset?



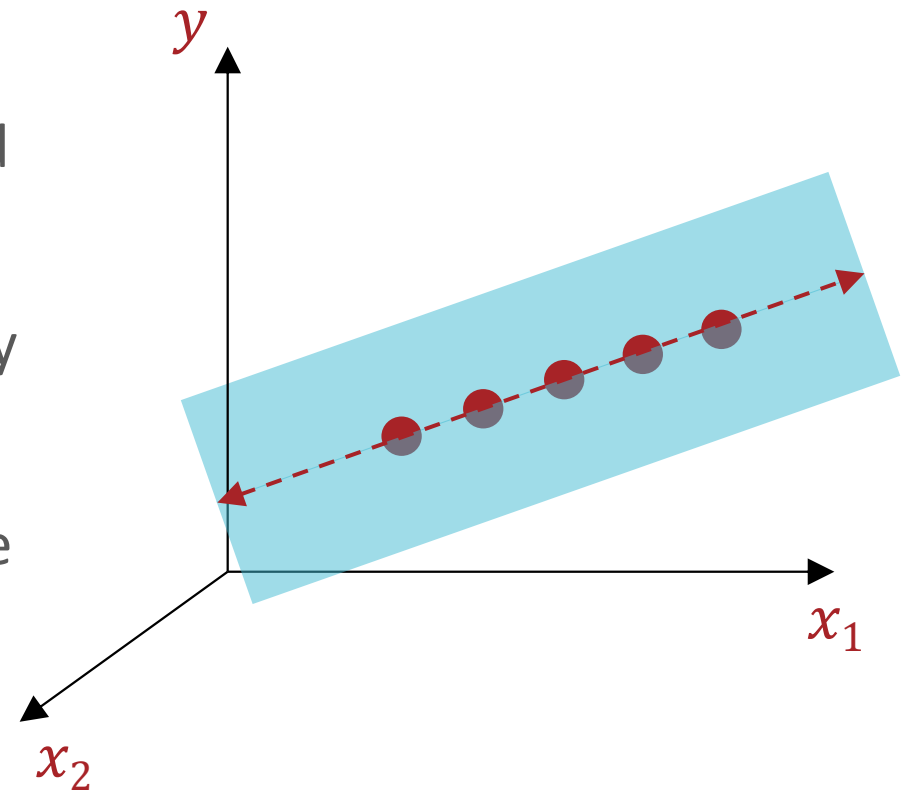
Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters θ) are there for the given dataset?



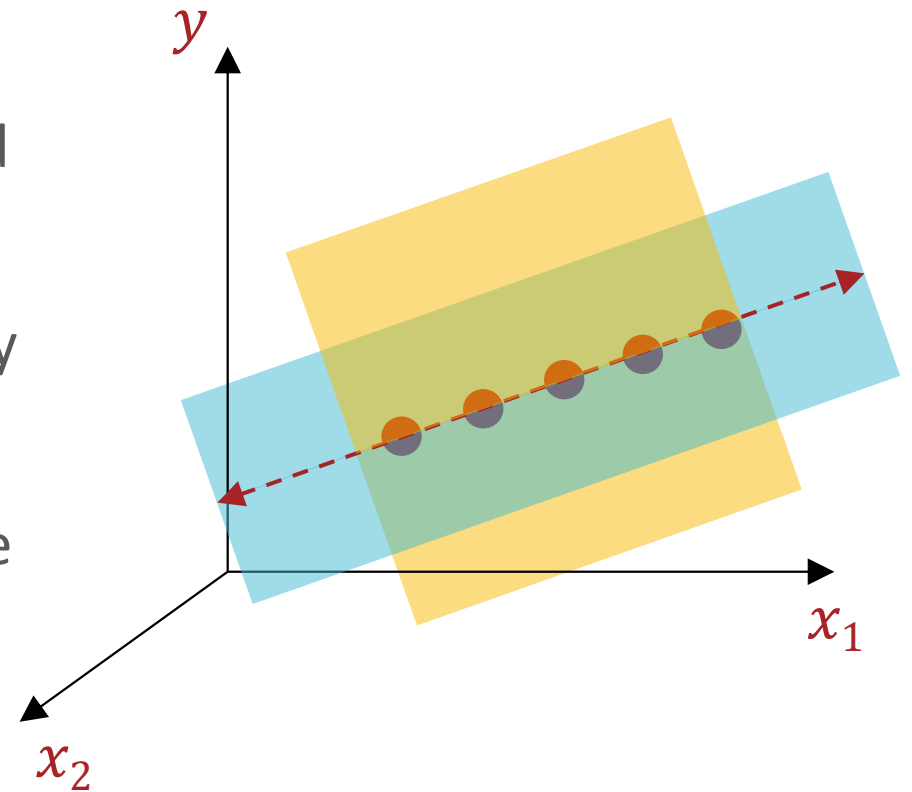
Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters θ) are there for the given dataset?



Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters θ) are there for the given dataset?



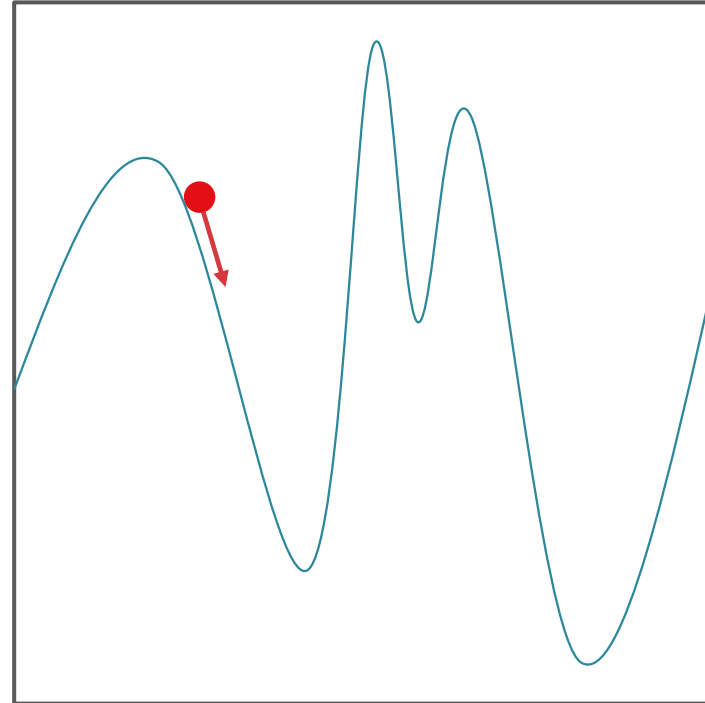
Closed Form Solution

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

1. Is $\mathbf{X}^T \mathbf{X}$ invertible?
 - When $N \gg D + 1$, $\mathbf{X}^T \mathbf{X}$ is (almost always) full rank and therefore, invertible!
 - If $\mathbf{X}^T \mathbf{X}$ is not invertible (occurs when one of the features is a linear combination of the others) then there are infinitely many solutions.
2. If so, how computationally expensive is inverting $\mathbf{X}^T \mathbf{X}$?
 - $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{D+1 \times D+1}$ so inverting $\mathbf{X}^T \mathbf{X}$ takes $O(D^3)$ time...
 - Computing $\mathbf{X}^T \mathbf{X}$ takes $O(ND^2)$ time
 - Can use gradient descent to (potentially) speed things up when N and D are large!

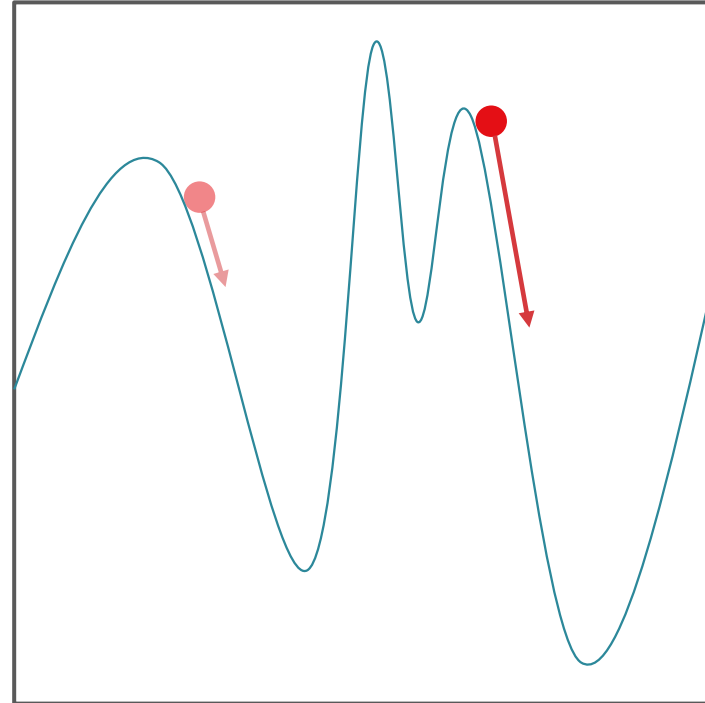
Gradient Descent: Intuition

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



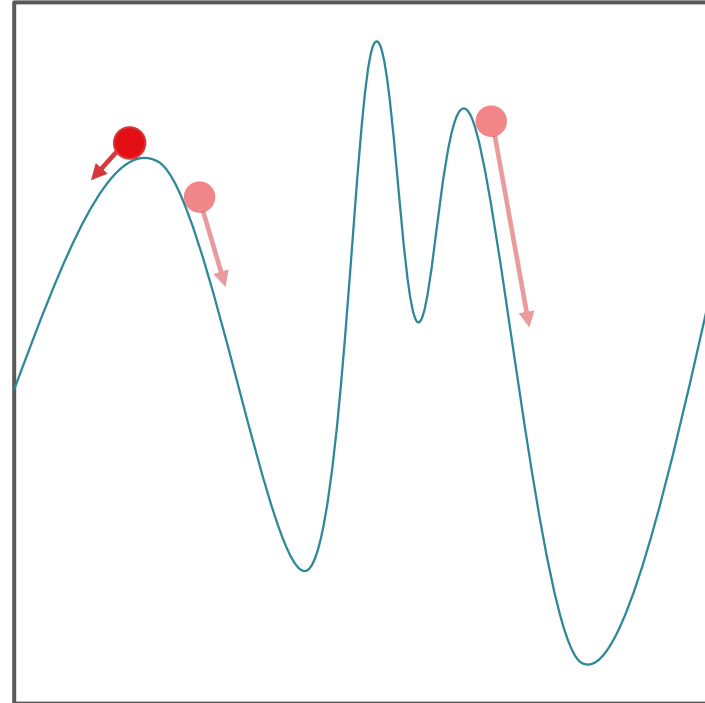
Gradient Descent: Intuition

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



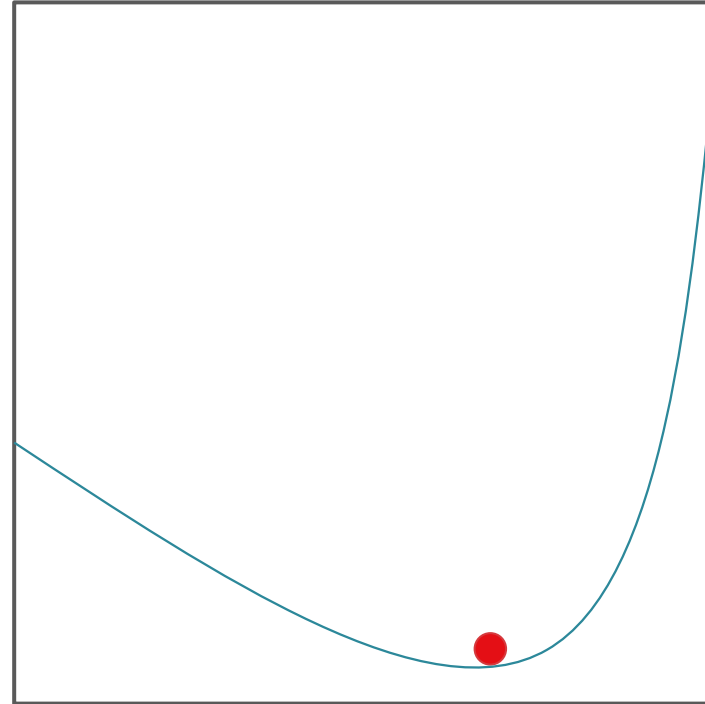
Gradient Descent: Intuition

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



Gradient Descent: Intuition

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



- Good news: the squared error is convex!

Recall: Minimizing the Squared Error

$$\ell_{\mathcal{D}}(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y^{(n)})^2 = \sum_{n=1}^N (\mathbf{x}^{(n)T} \mathbf{w} - y^{(n)})^2$$

$$= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \text{ where } \|\mathbf{z}\|_2 = \sqrt{\sum_{d=1}^D z_d^2} = \sqrt{\mathbf{z}^T \mathbf{z}}$$

$$= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$= (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

$$\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}) = (2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y})$$

$$H_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}) = 2\mathbf{X}^T \mathbf{X}$$

$H_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w})$ is positive semi-definite

Gradient Descent: Step Direction

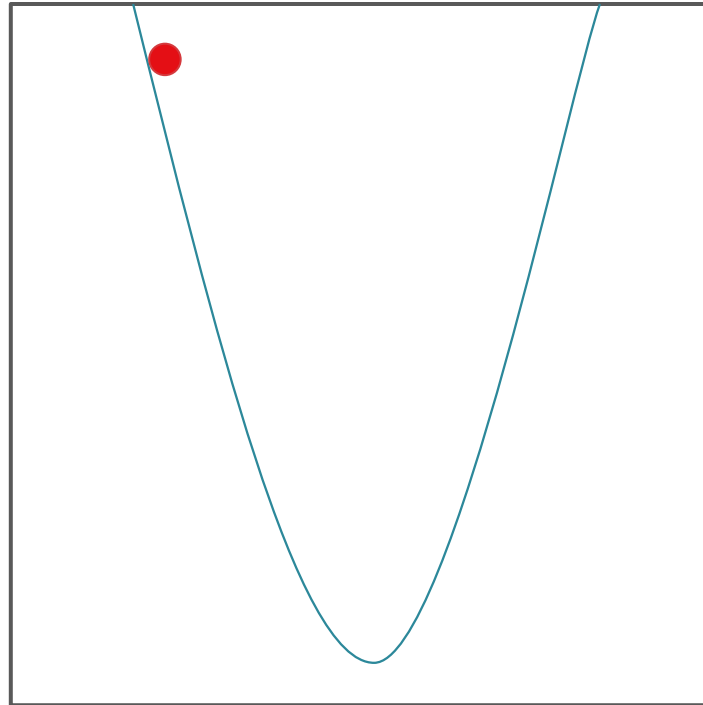
- Suppose the current weight vector is $\mathbf{w}^{(t)}$
- Move some distance, η , in the “most downhill” direction, $\hat{\mathbf{v}}$:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta \hat{\mathbf{v}}$$

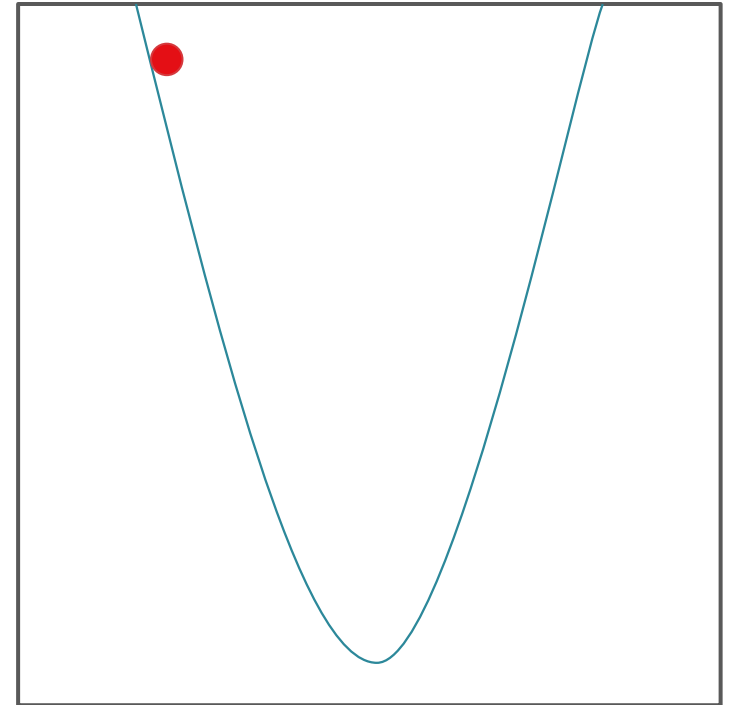
- The gradient points in the direction of steepest *increase* ...
- ... so $\hat{\mathbf{v}}$ is a unit vector pointing in the opposite direction:

$$\hat{\mathbf{v}}^{(t)} = - \frac{\nabla_{\mathbf{w}} \ell_{\mathcal{D}} (\mathbf{w}^{(t)})}{\|\nabla_{\mathbf{w}} \ell_{\mathcal{D}} (\mathbf{w}^{(t)})\|}$$

Gradient Descent: Step Size

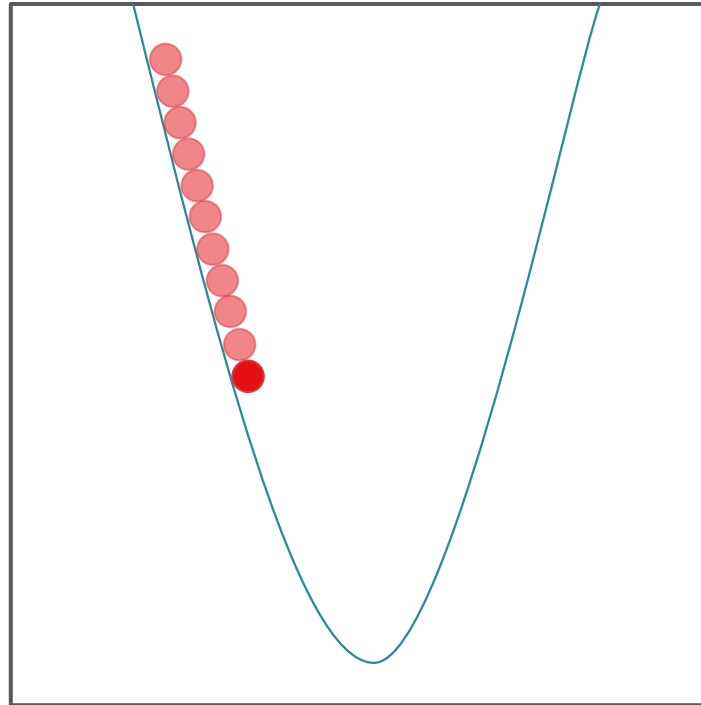


Small η

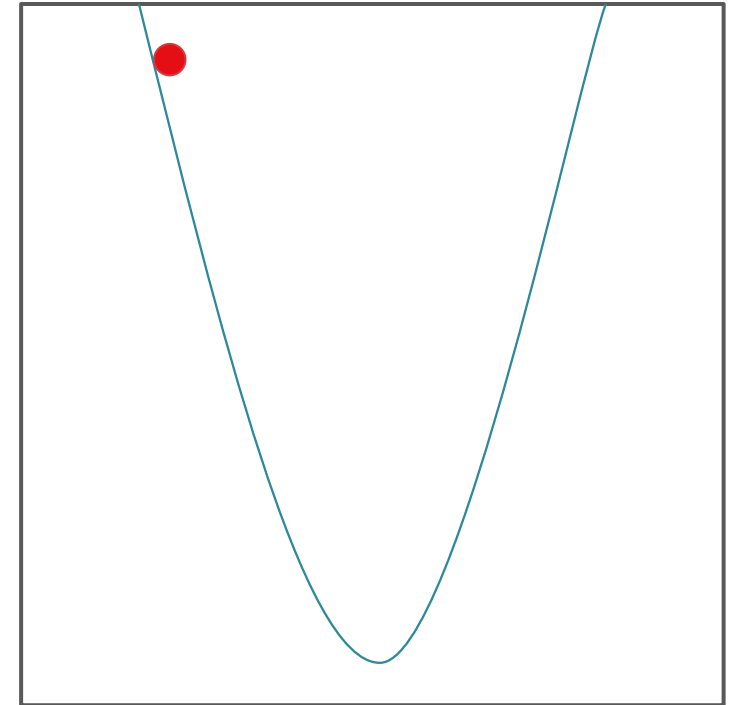


Large η

Gradient Descent: Step Size

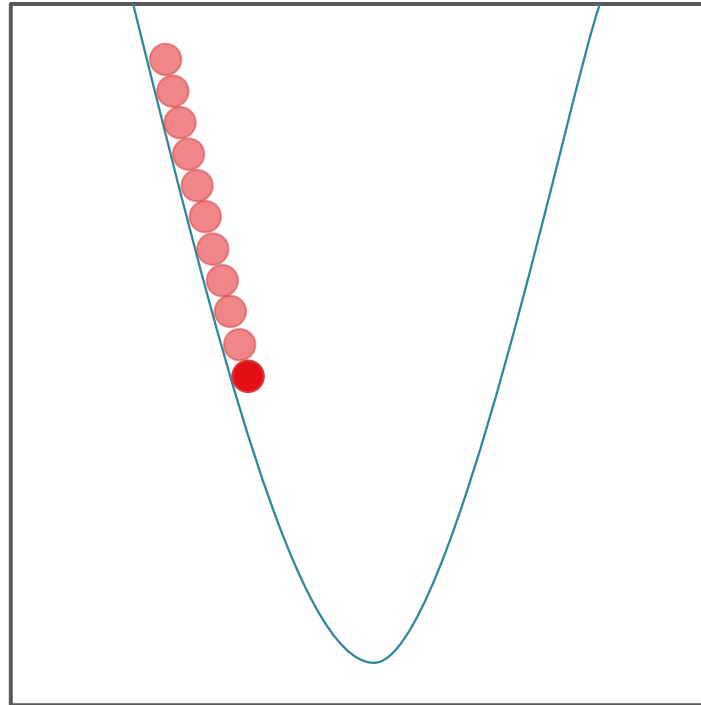


Small η

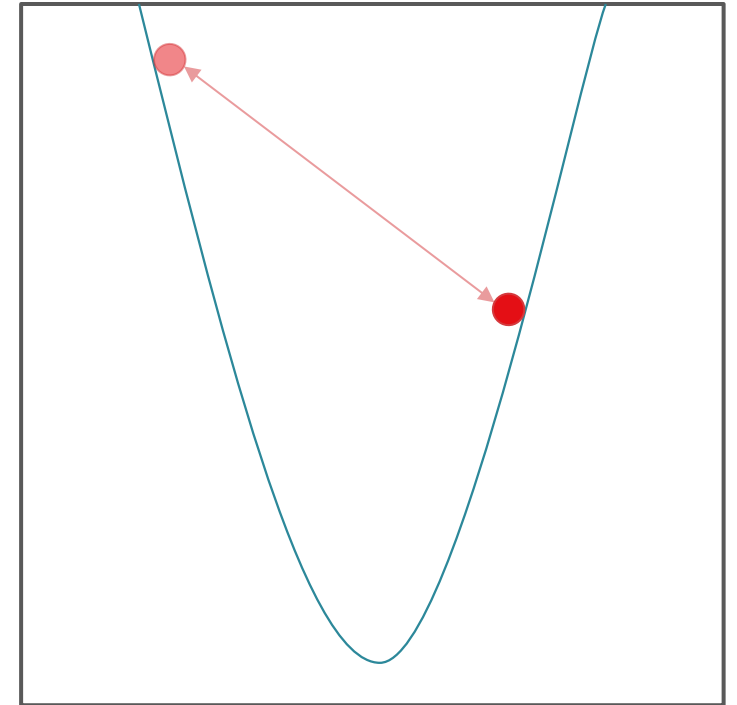


Large η

Gradient Descent: Step Size



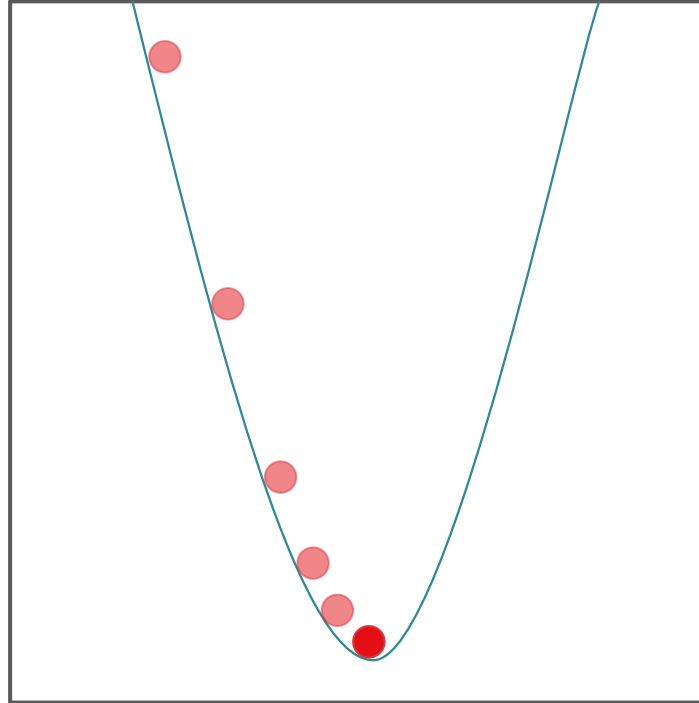
Small η



Large η

Gradient Descent: Step Size

- Use a variable $\eta^{(t)}$ instead of a fixed η !



- Set $\eta^{(t)} = \eta^{(0)} \|\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\|$
- $\|\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\|$ decreases as $\ell_{\mathcal{D}}$ approaches its minimum $\rightarrow \eta^{(t)}$ (hopefully) decreases over time

Gradient Descent

- $\hat{\mathbf{v}}^{(t)} = -\frac{\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})}{\|\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\|}$
- $\eta^{(t)} = \eta^{(0)} \|\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\|$
- $\begin{aligned}\mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} + \eta^{(t)} \hat{\mathbf{v}}^{(t)} \\ &= \mathbf{w}^{(t)} + (\eta^{(0)} \|\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\|) \left(-\frac{\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})}{\|\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\|} \right) \\ &= \mathbf{w}^{(t)} - \eta^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\end{aligned}$

Gradient Descent

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \eta^{(0)}$
 1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
 2. While TERMINATION CRITERION is not satisfied
 - a. Compute the gradient:
 $\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$
 - b. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$
 - c. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{w}^{(t)}$

Gradient Descent

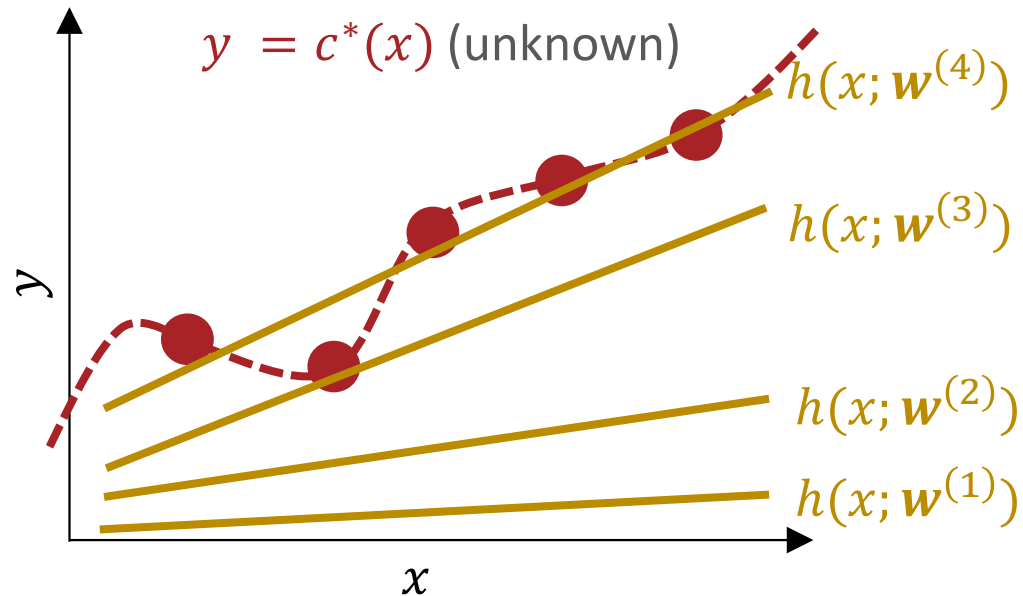
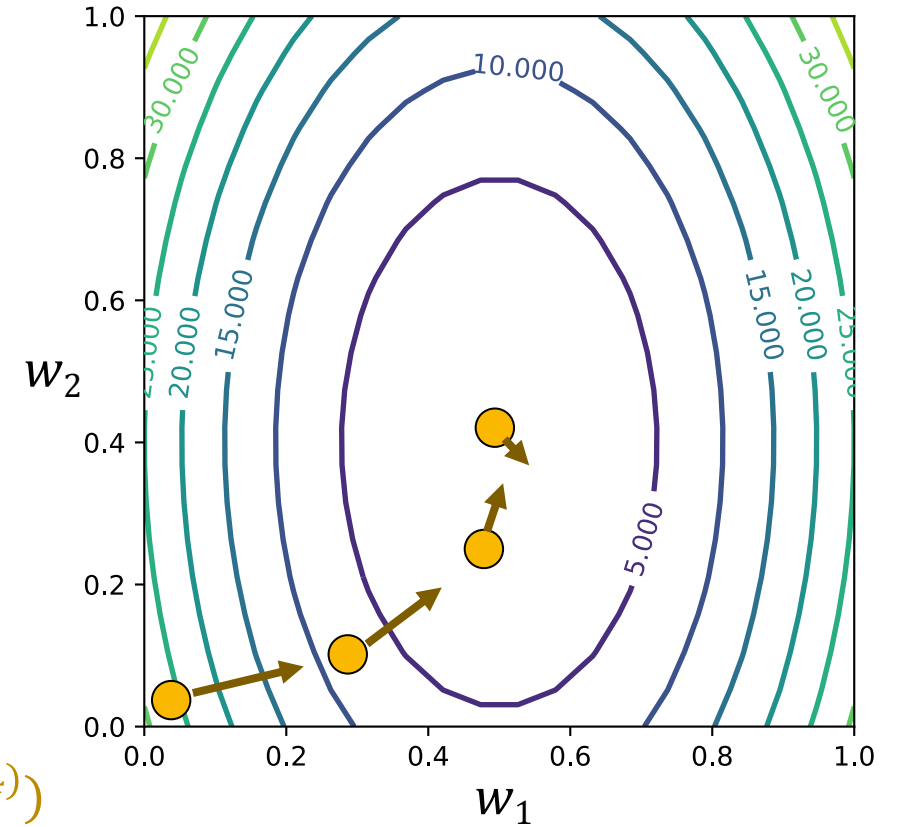
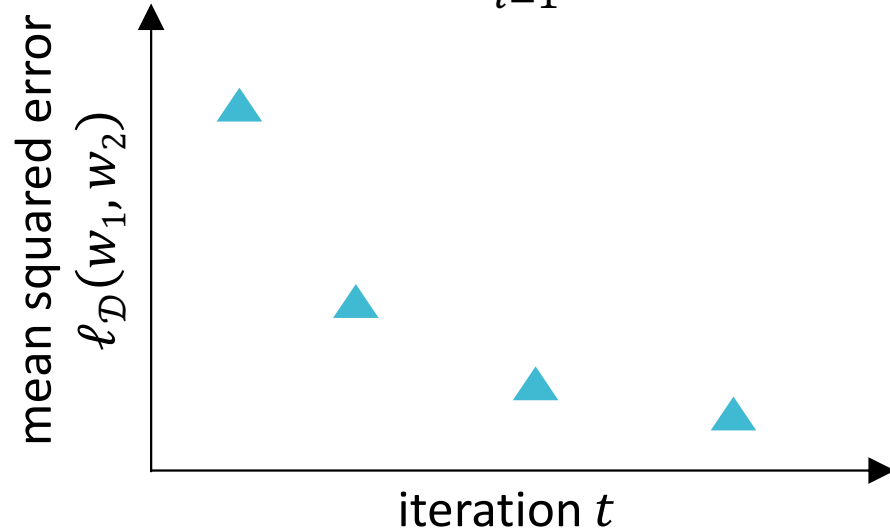
- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \eta^{(0)}, \epsilon$
 1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
 2. While $\|\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\| > \epsilon$
 - a. Compute the gradient:
 $\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$
 - b. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$
 - c. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{w}^{(t)}$

Gradient Descent

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \eta^{(0)}, T$
- 1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
- 2. While $t < T$
 - a. Compute the gradient:
 $\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$
 - b. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$
 - c. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{w}^{(t)}$

Gradient Descent for Linear Regression

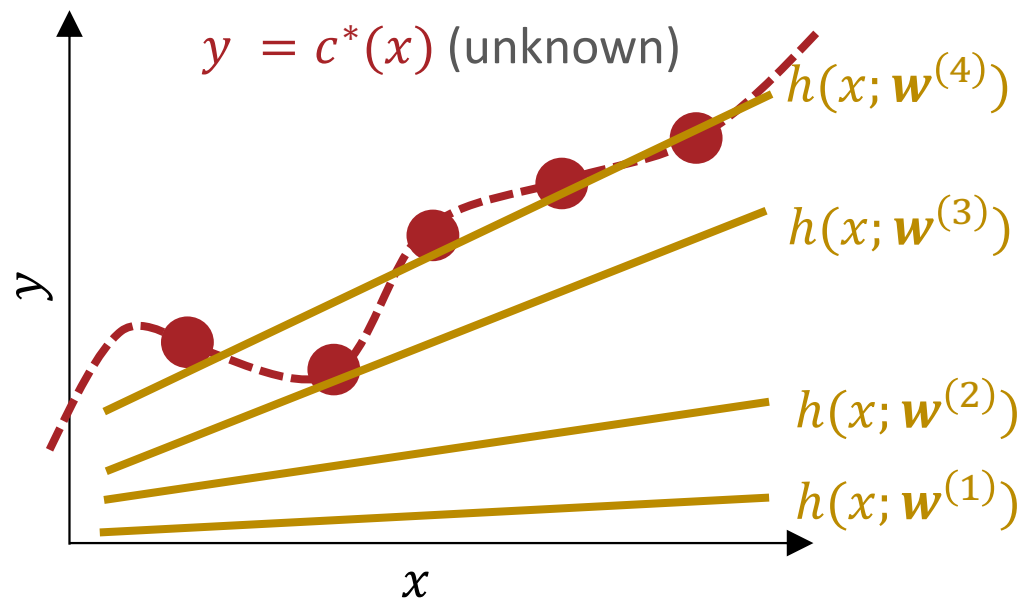
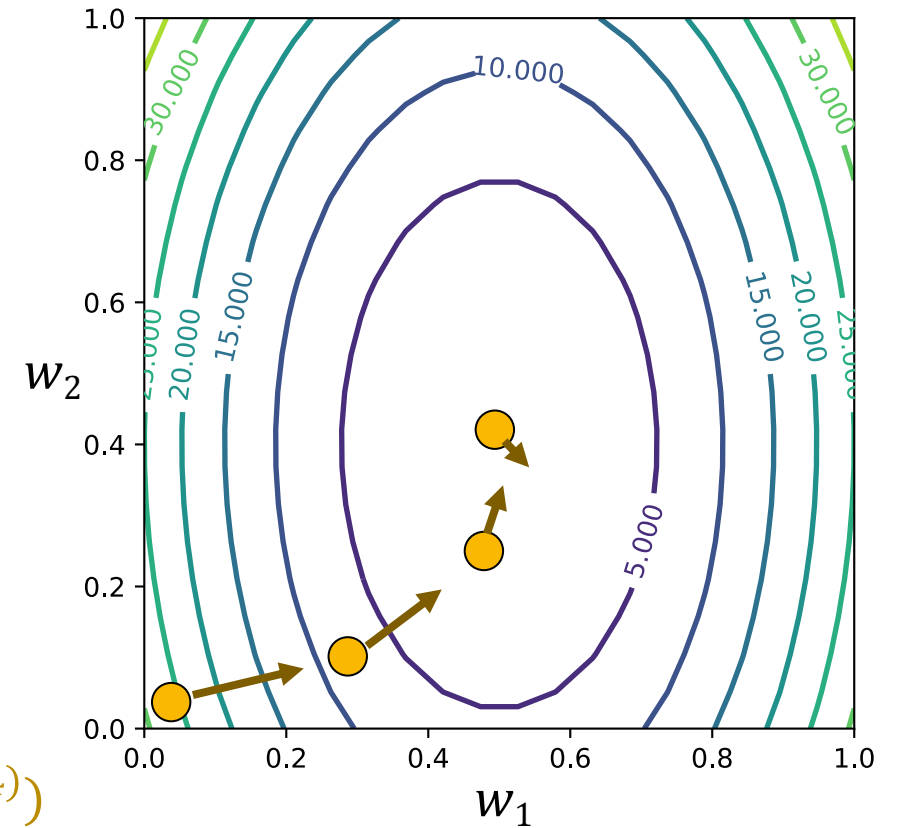
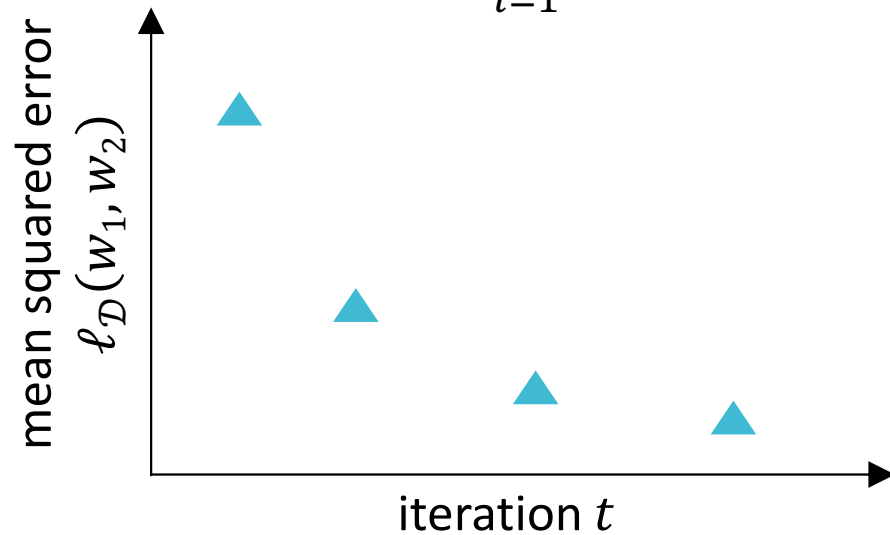
$$\ell_{\mathcal{D}}(w_1, w_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$



t	w_1	w_2	$\ell_{\mathcal{D}}(w_1, w_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

Why Gradient Descent for Linear Regression?

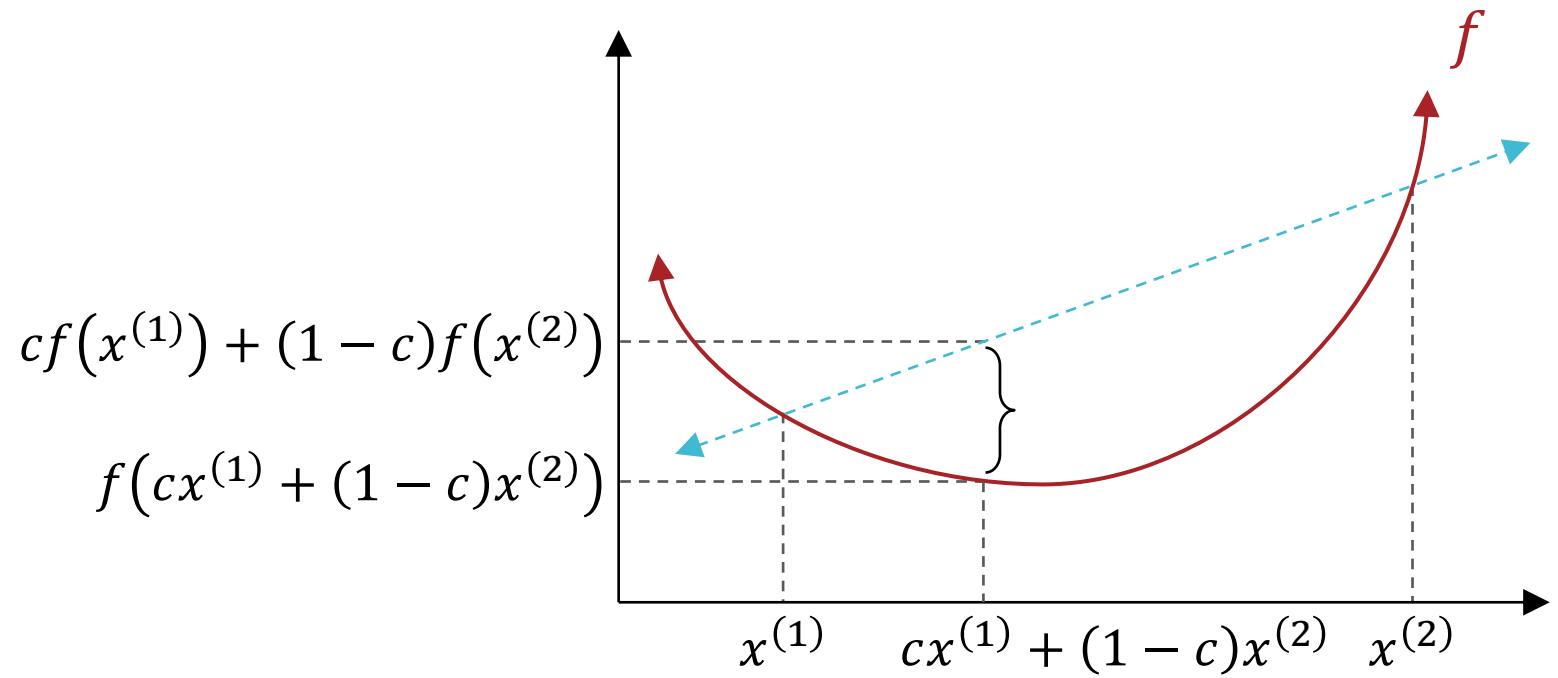
$$\ell_{\mathcal{D}}(w_1, w_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$



t	w_1	w_2	$\ell_{\mathcal{D}}(w_1, w_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

Convexity

- A function $f: \mathbb{R}^D \rightarrow \mathbb{R}$ is convex if
 $\forall \mathbf{x}^{(1)} \in \mathbb{R}^D, \mathbf{x}^{(2)} \in \mathbb{R}^D$ and $0 \leq c \leq 1$
 $f(c\mathbf{x}^{(1)} + (1-c)\mathbf{x}^{(2)}) \leq cf(\mathbf{x}^{(1)}) + (1-c)f(\mathbf{x}^{(2)})$

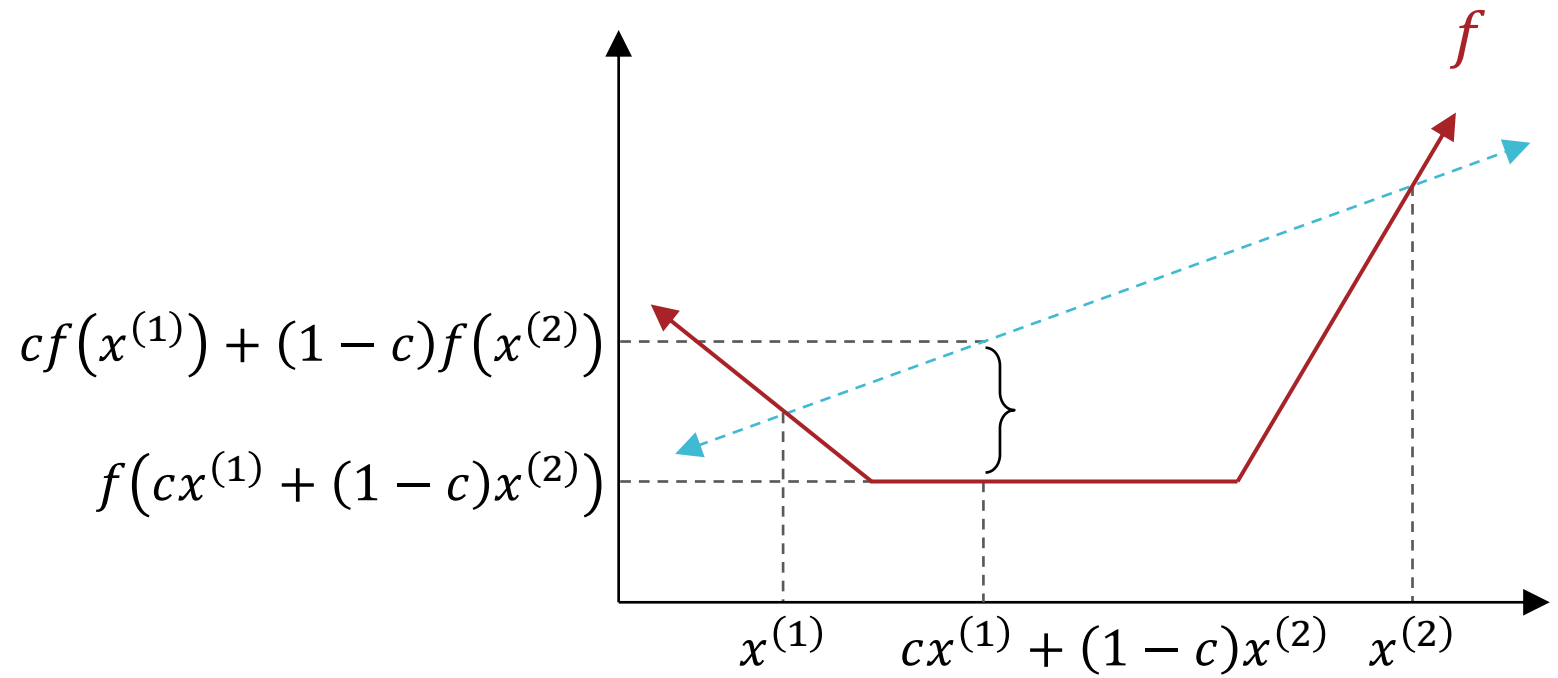


Convexity

- A function $f: \mathbb{R}^D \rightarrow \mathbb{R}$ is convex if

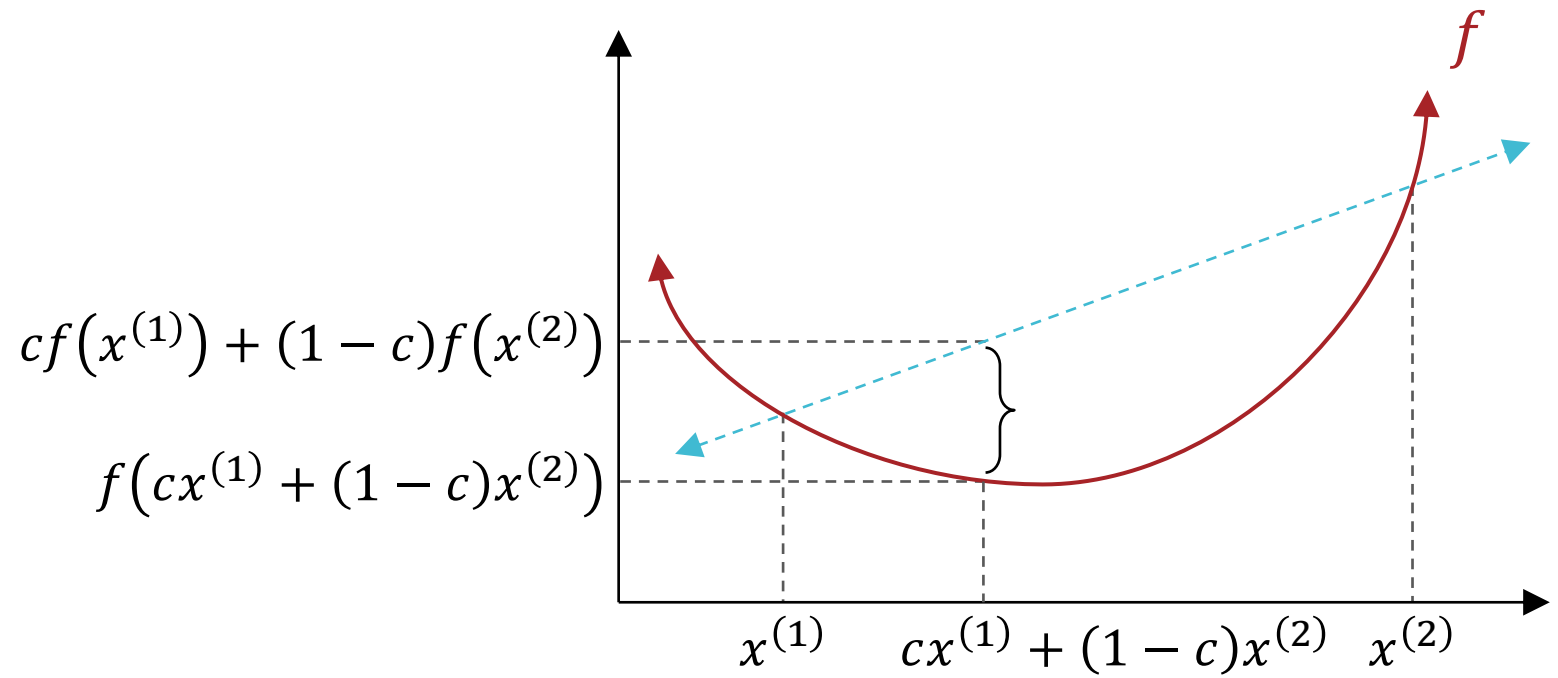
$$\forall \mathbf{x}^{(1)} \in \mathbb{R}^D, \mathbf{x}^{(2)} \in \mathbb{R}^D \text{ and } 0 \leq c \leq 1$$

$$f(c\mathbf{x}^{(1)} + (1-c)\mathbf{x}^{(2)}) \leq cf(\mathbf{x}^{(1)}) + (1-c)f(\mathbf{x}^{(2)})$$

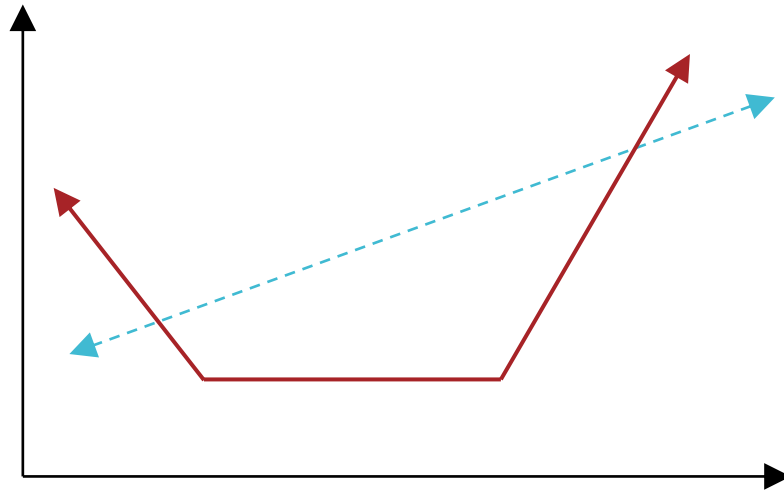


Convexity

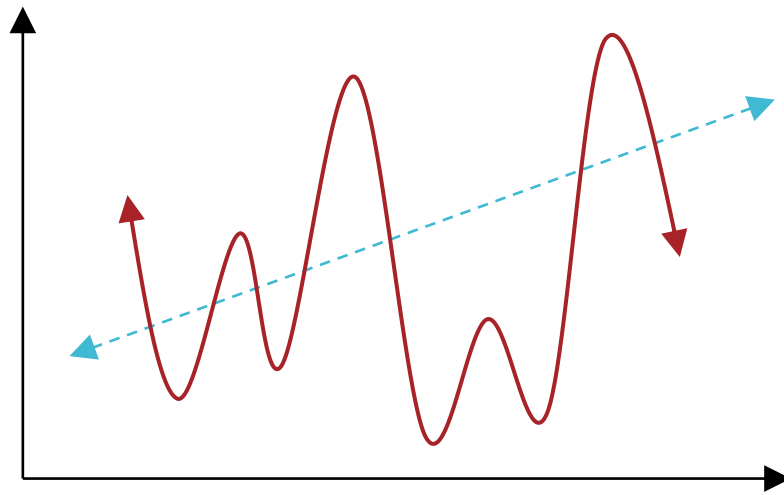
- A function $f: \mathbb{R}^D \rightarrow \mathbb{R}$ is *strictly convex* if
 $\forall \mathbf{x}^{(1)} \in \mathbb{R}^D, \mathbf{x}^{(2)} \in \mathbb{R}^D$ and $0 < c < 1$
 $f(c\mathbf{x}^{(1)} + (1 - c)\mathbf{x}^{(2)}) < cf(\mathbf{x}^{(1)}) + (1 - c)f(\mathbf{x}^{(2)})$



Convexity

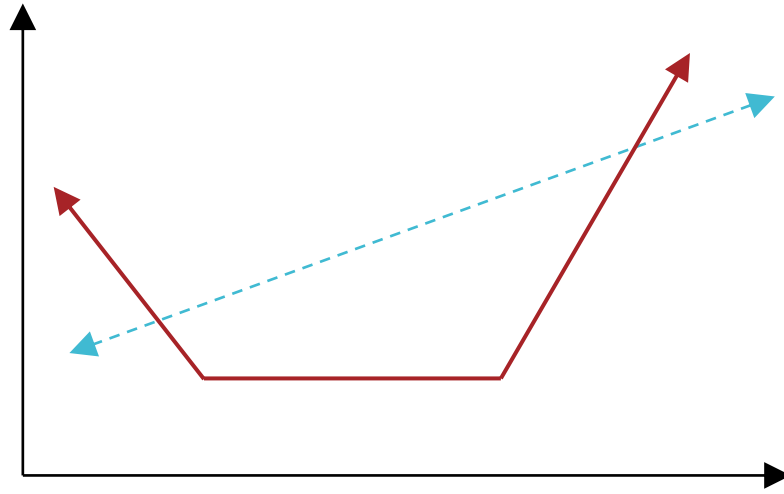


Convex functions



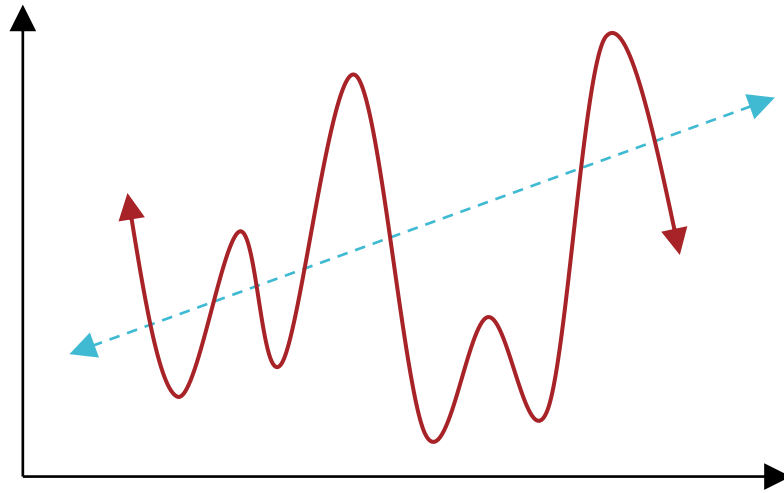
Non-convex functions

Convexity



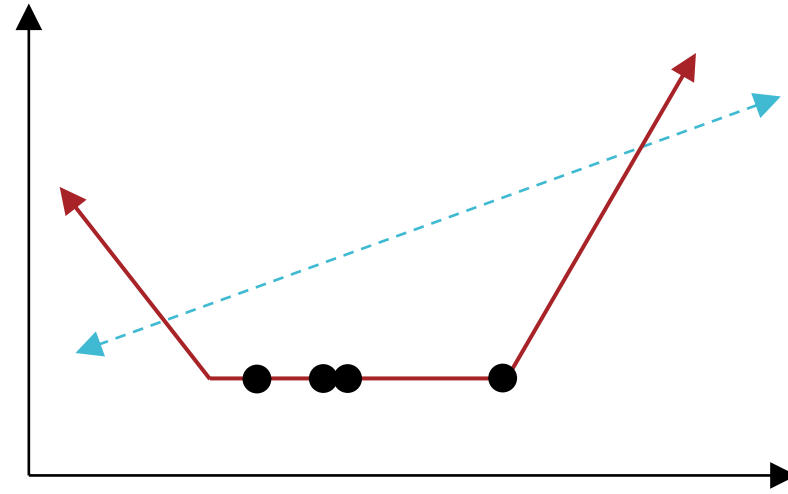
Given a function $f: \mathbb{R}^D \rightarrow \mathbb{R}$

- \mathbf{x}^* is a global minimum iff $f(\mathbf{x}^*) \leq f(\mathbf{x}) \forall \mathbf{x} \in \mathbb{R}^D$

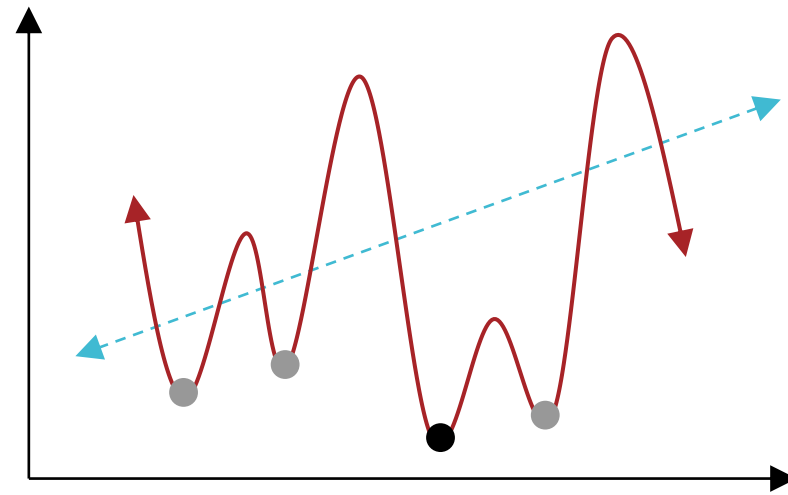


- \mathbf{x}^* is a local minimum iff $\exists \epsilon$ s.t. $f(\mathbf{x}^*) \leq f(\mathbf{x}) \forall \mathbf{x}$ s.t. $\|\mathbf{x} - \mathbf{x}^*\|_2 < \epsilon$

Convexity

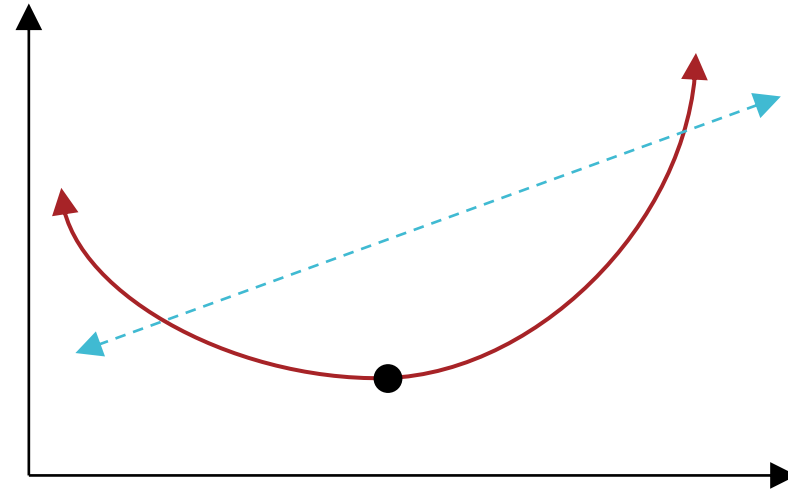


Convex functions:
Each local minimum is a
global minimum!

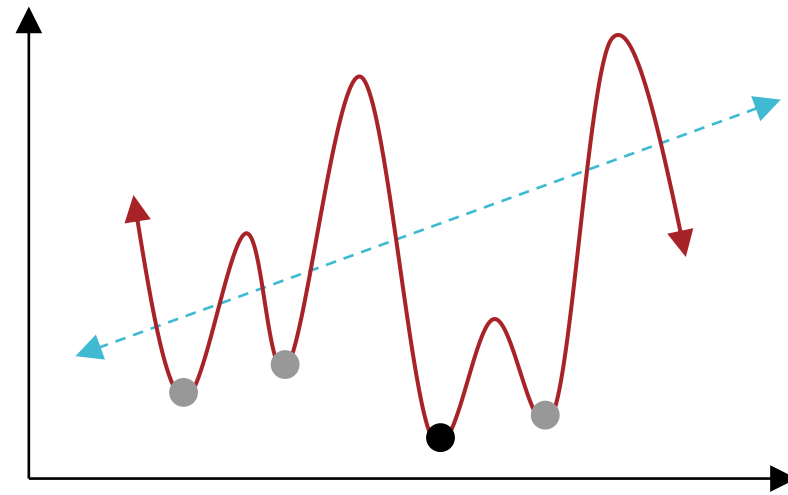


Non-convex functions:
A local minimum may or may
not be a global minimum...

Convexity



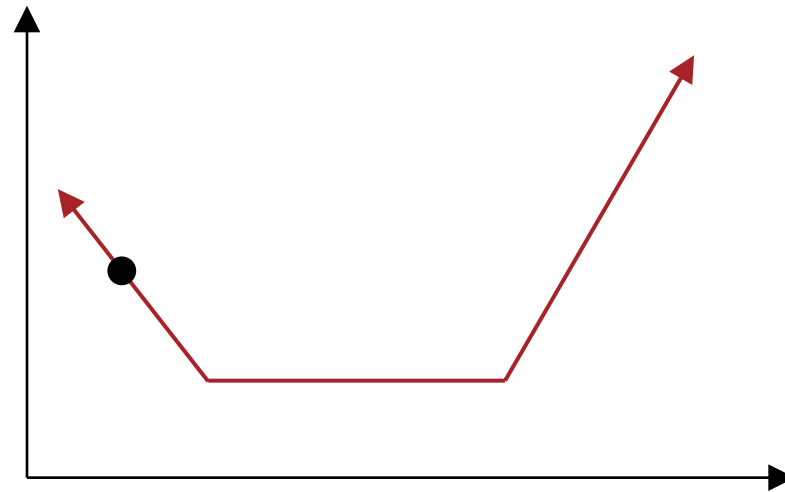
Strictly convex functions:
There exists a unique global minimum!



Non-convex functions:
A local minimum may or may not be a global minimum...

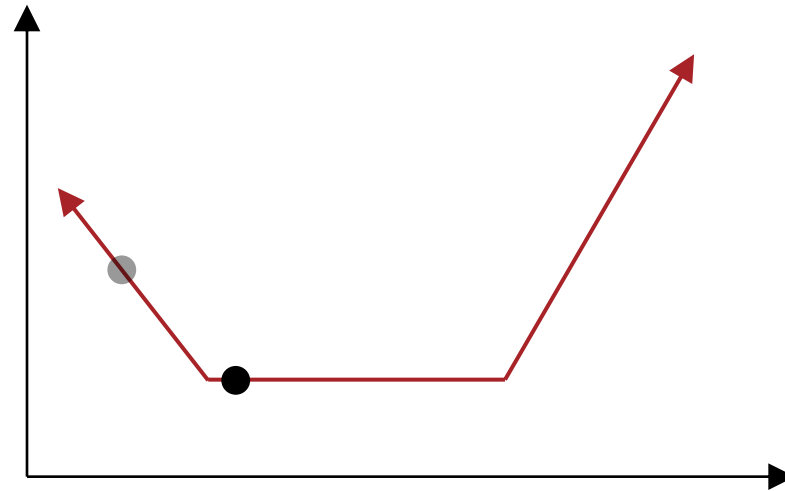
Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
 - Works great if the objective function is convex!



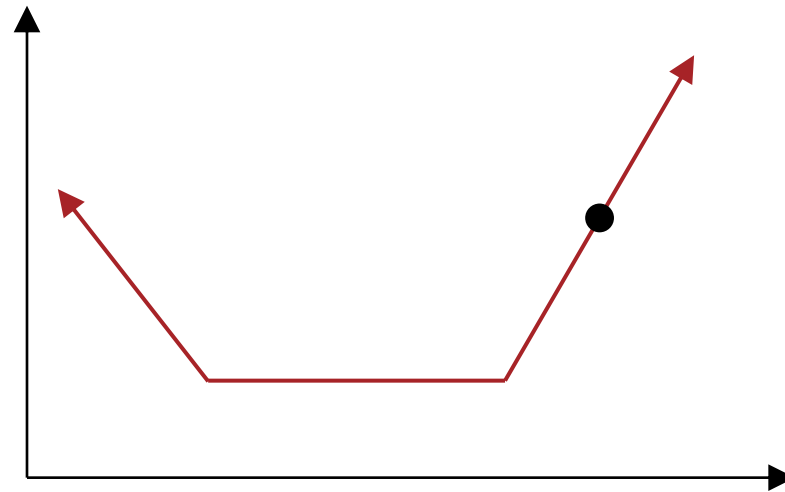
Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
 - Works great if the objective function is convex!



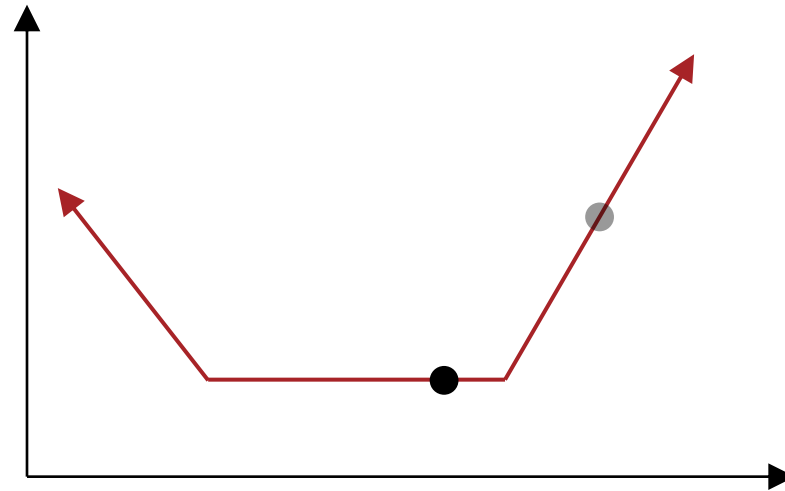
Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
 - Works great if the objective function is convex!



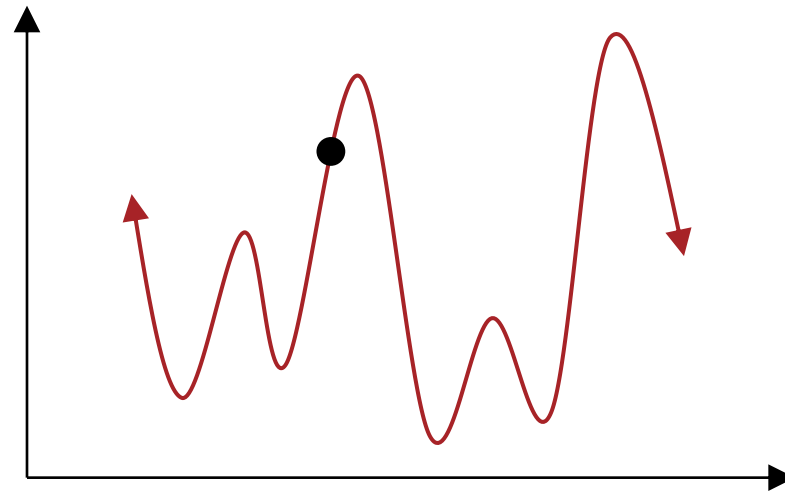
Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
 - Works great if the objective function is convex!



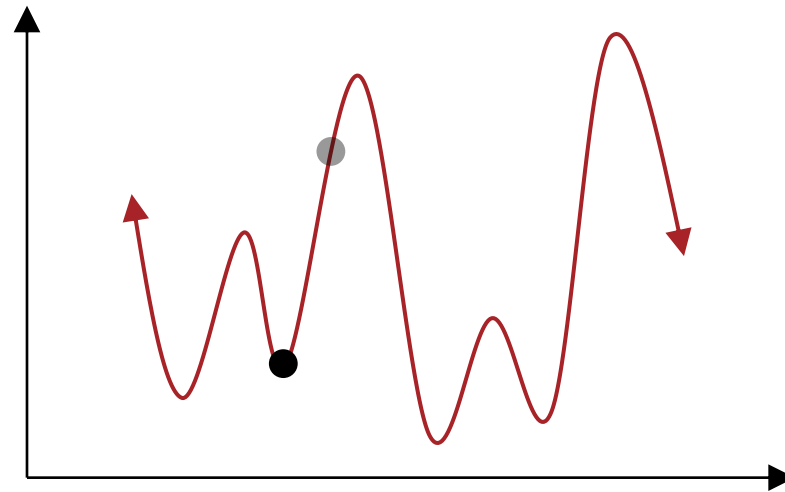
Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
 - Not ideal if the objective function is non-convex...



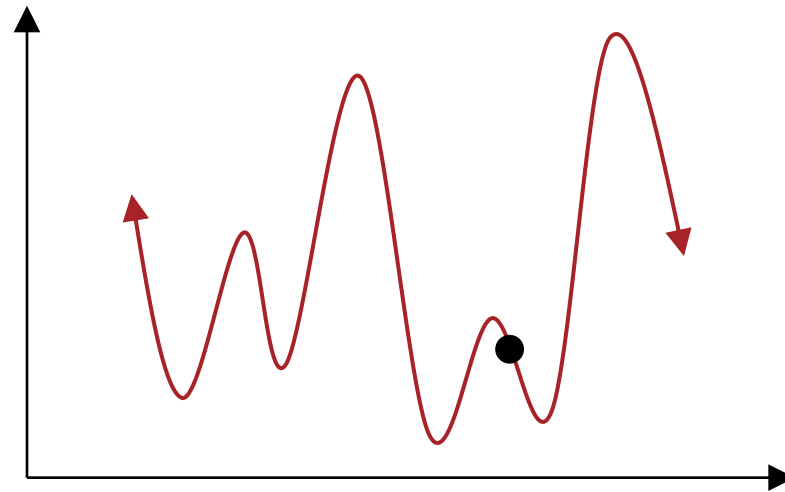
Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
 - Not ideal if the objective function is non-convex...



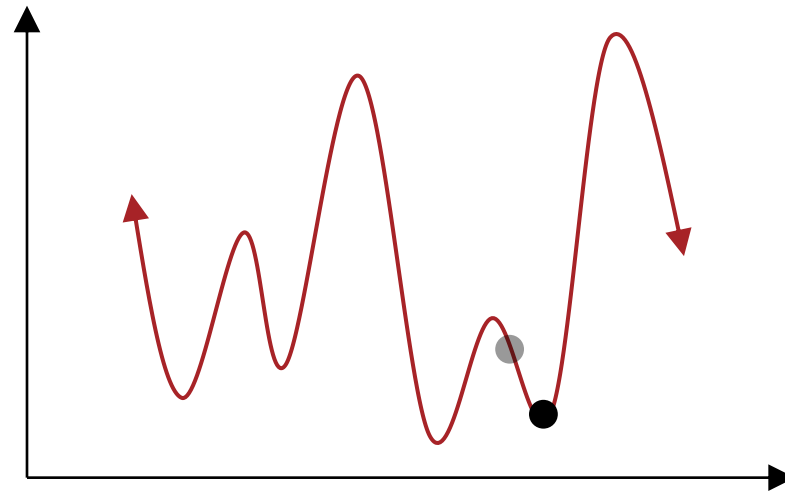
Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
 - Not ideal if the objective function is non-convex...



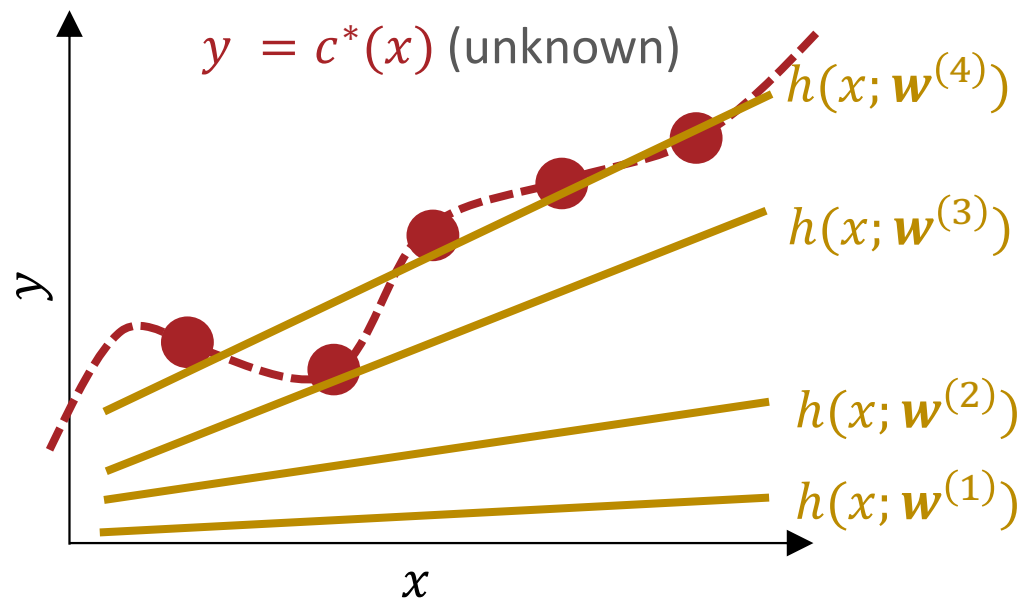
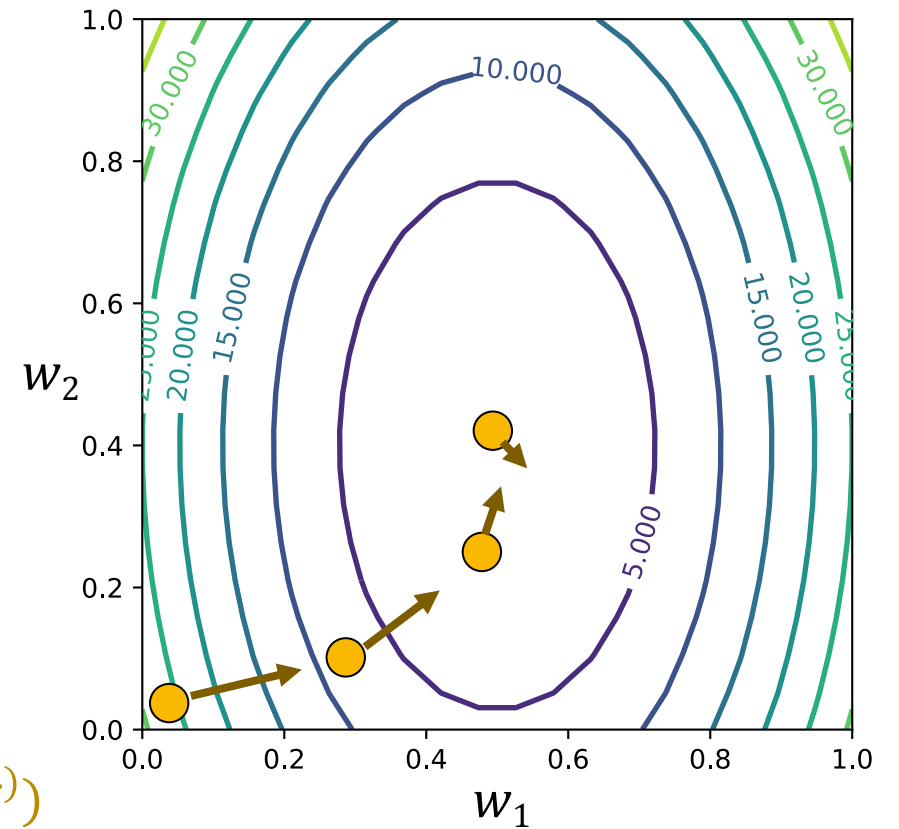
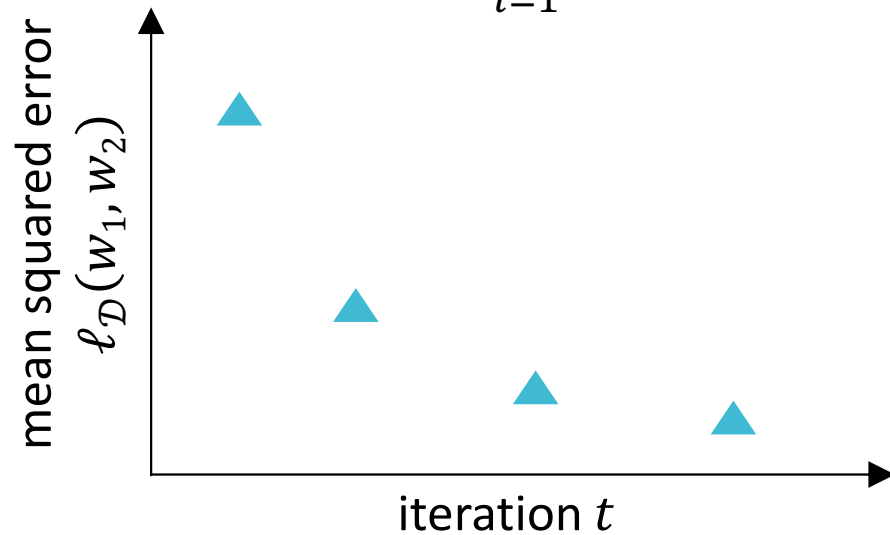
Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
 - Not ideal if the objective function is non-convex...



The mean squared error is convex (but not always strictly convex)

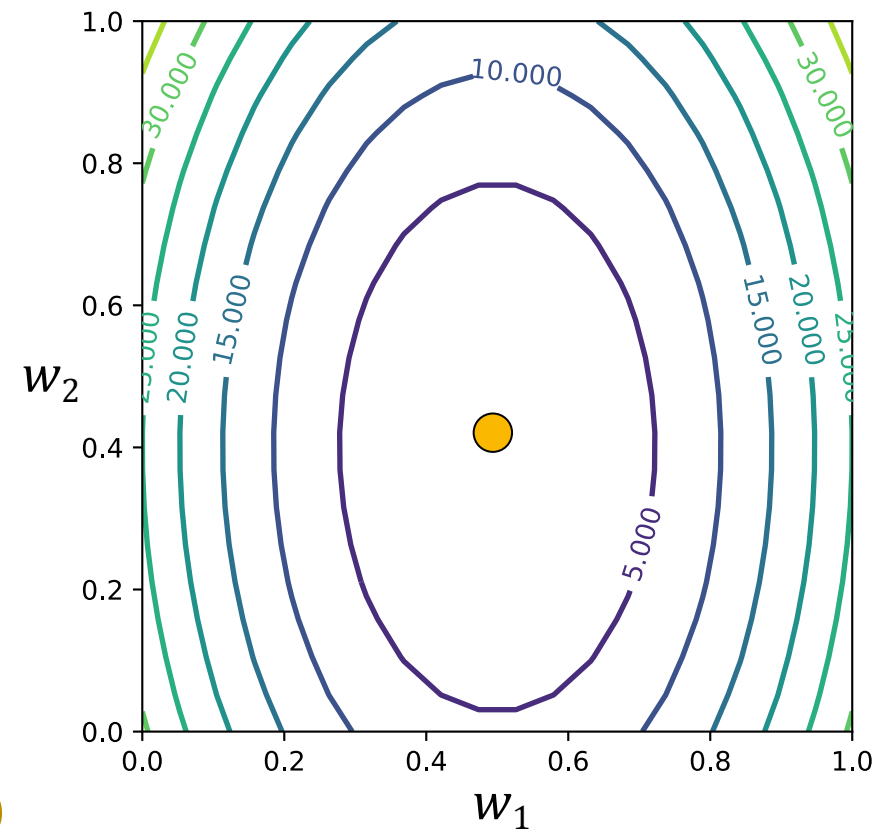
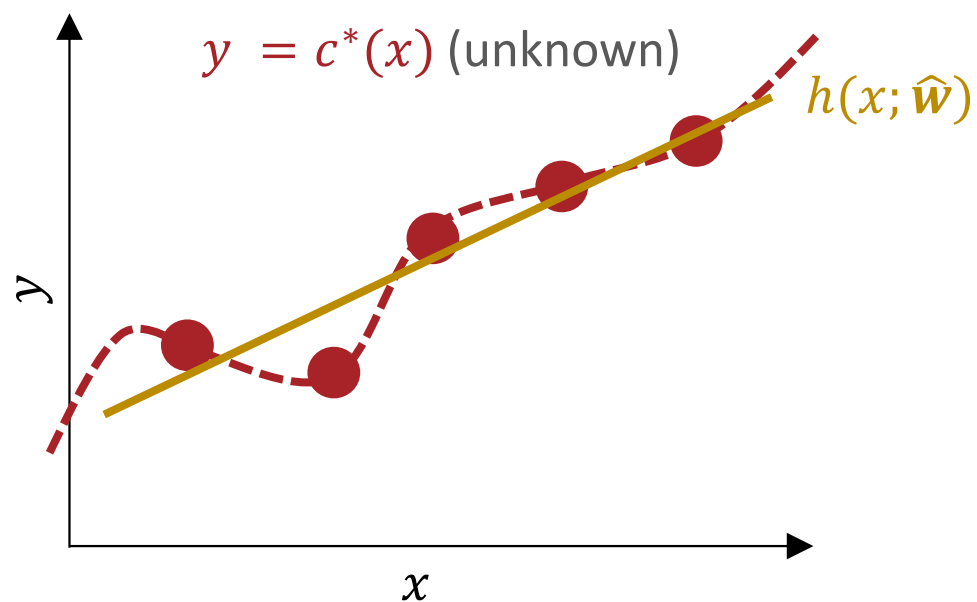
$$\ell_{\mathcal{D}}(w_1, w_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$



t	w_1	w_2	$\ell_{\mathcal{D}}(w_1, w_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

Closed Form Optimization

$$\hat{w} = (X^T X)^{-1} X^T y$$



t	w_1	w_2	$l_{\mathcal{D}}(w_1, w_2)$
1	0.59	0.43	0.2

Key Takeaways

- Convexity vs. non-convexity
 - Strong vs. weak convexity
 - Implications for local, global and unique optima
- Gradient descent
 - Effect of step size
 - Termination criteria