

# 10-701: Introduction to Machine Learning Lecture 4 – Linear Regression

Henry Chai & Zack Lipton

9/11/23

# Front Matter

- Announcements:
  - HW1 released 9/6, due 9/20 at 11:59 PM
- Recommended Readings:
  - Bishop, [Section 3.2](#)
  - Murphy, [Sections 7.1-7.3](#)

# Recall: Regression

- Learning to diagnose heart disease

as a **(supervised)**

**regression** task

features

targets

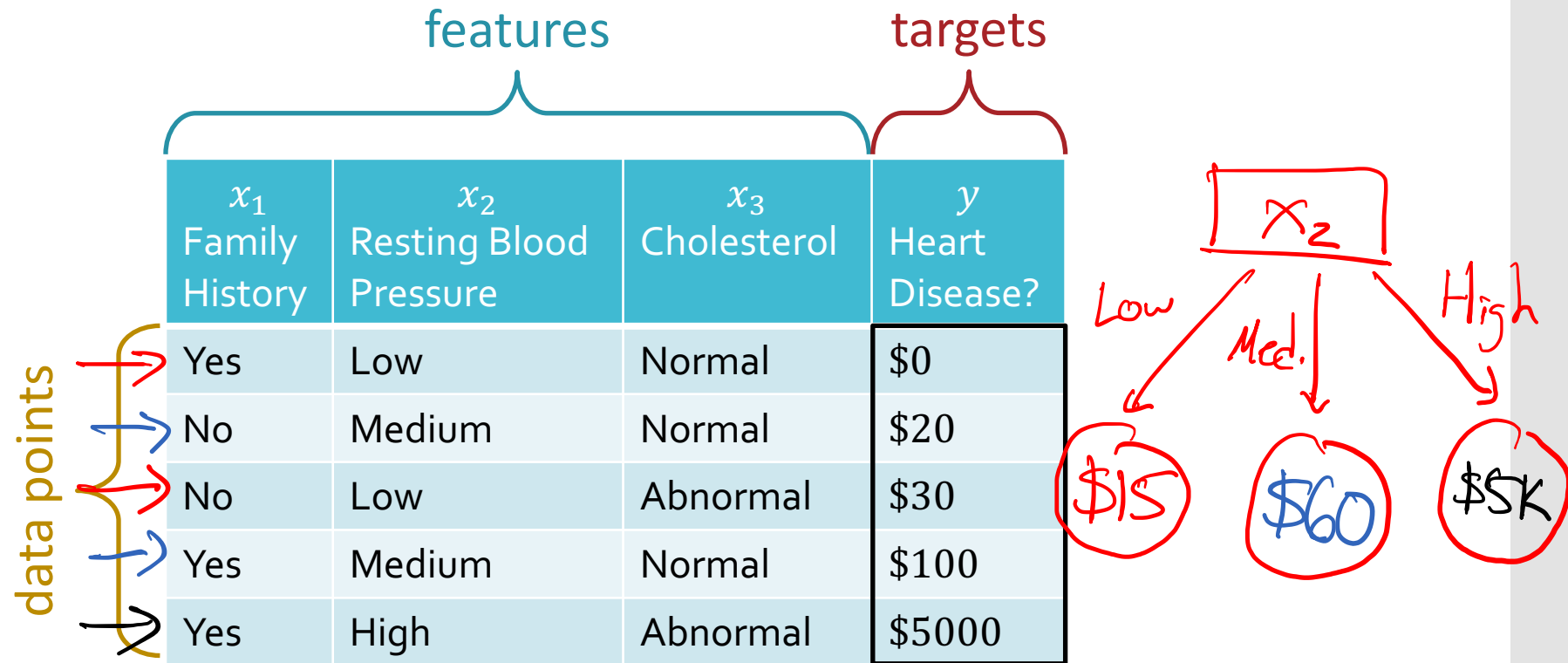
data points

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	$y$ Heart Disease?
Yes	Low	Normal	\$0
No	Medium	Normal	\$20
No	Low	Abnormal	\$30
Yes	Medium	Normal	\$100
Yes	High	Abnormal	\$5000

# Decision Tree Regression

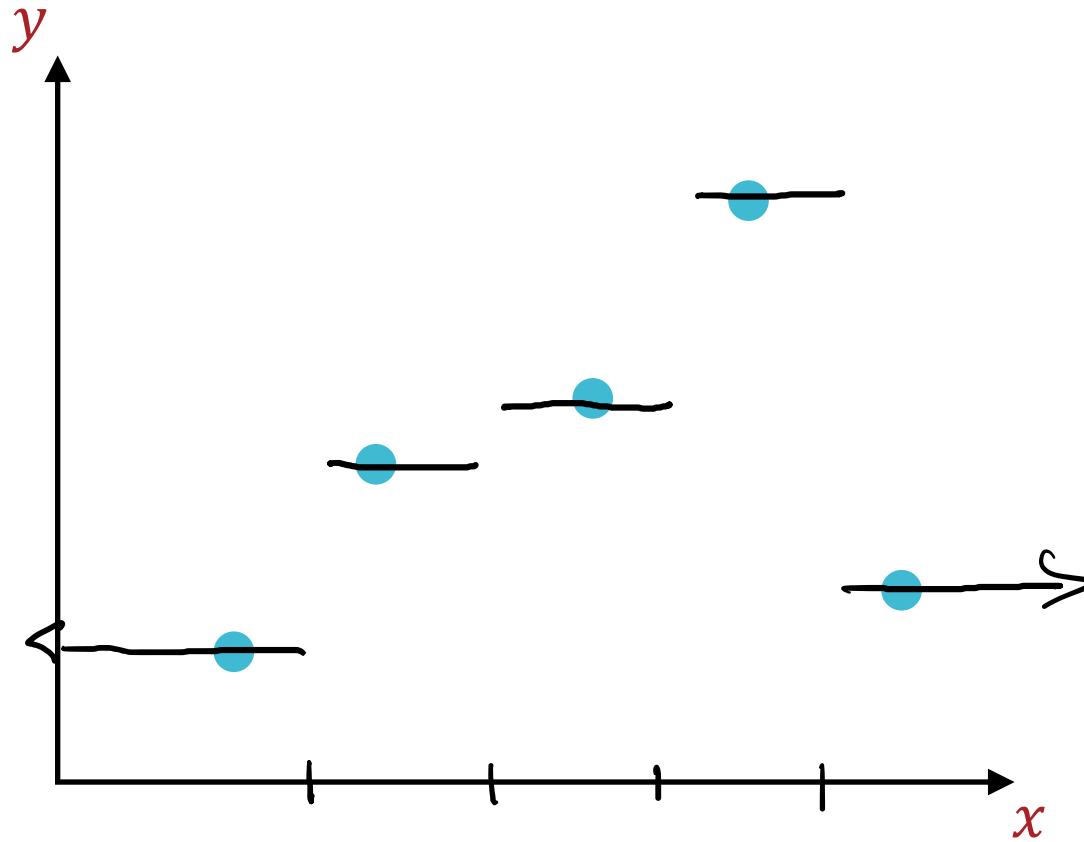
- Learning to diagnose heart disease

as a **(supervised)** regression task



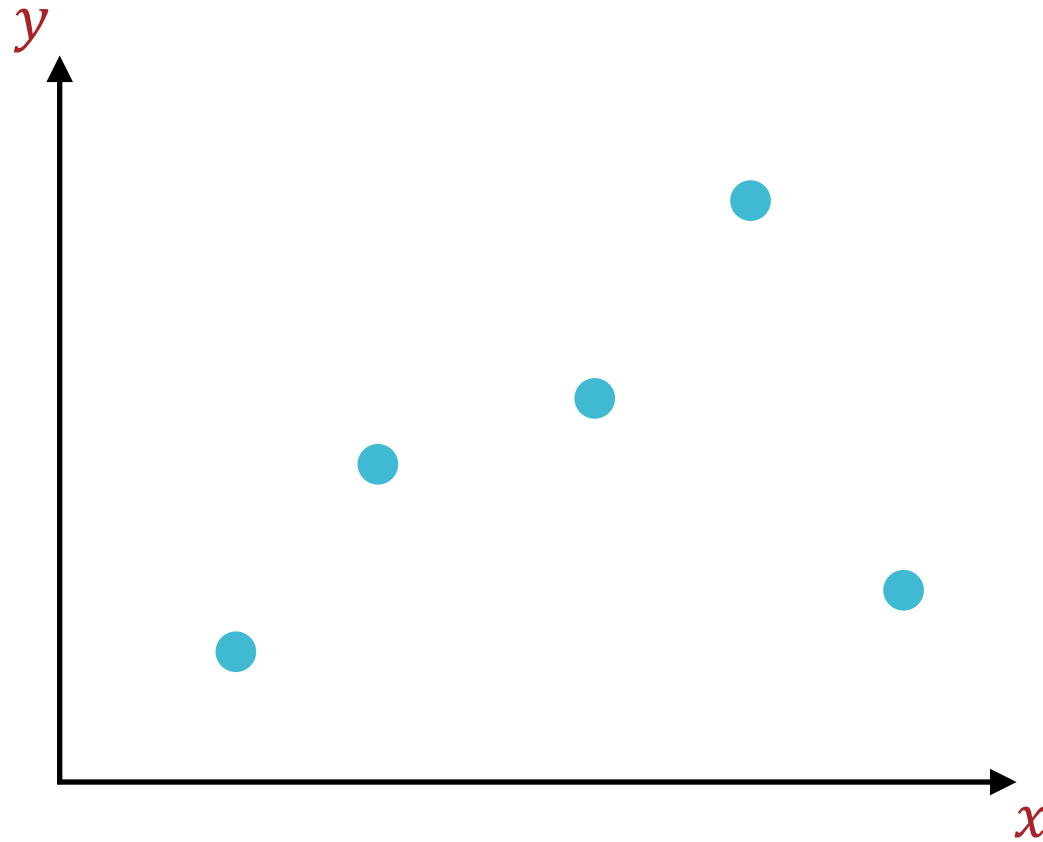
# 1-NN Regression

- Suppose we have real-valued targets  $y \in \mathbb{R}$  and one-dimensional inputs  $x \in \mathbb{R}$



# 2-NN Regression?

- Suppose we have real-valued targets  $y \in \mathbb{R}$  and one-dimensional inputs  $x \in \mathbb{R}$



# Linear Regression

- Suppose we have real-valued targets  $y \in \mathbb{R}$  and  $D$ -dimensional inputs  $\mathbf{x} = [x_1, \dots, x_D]^T \in \mathbb{R}^D$
- **Assume**

$$y = \mathbf{w}^T \mathbf{x} + w_0$$

# Linear Regression

- Suppose we have real-valued targets  $y \in \mathbb{R}$  and  $D$ -dimensional inputs  $\mathbf{x} = [1, x_1, \dots, x_D]^T \in \mathbb{R}^{D+1}$

- Assume

$$\rightarrow \mathbf{w} = [w_0, w_1, \dots, w_D]^T$$
$$y = \mathbf{w}^T \mathbf{x}$$



# Linear Regression

- Suppose we have real-valued targets  $y \in \mathbb{R}$  and  $D$ -dimensional inputs  $\mathbf{x} = [1, x_1, \dots, x_D]^T \in \mathbb{R}^{D+1}$

- **Assume**

$$y = \mathbf{w}^T \mathbf{x}$$

- Notation: given training data  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$

$$\bullet \mathbf{X} = \begin{bmatrix} 1 & \mathbf{x}^{(1)T} \\ 1 & \mathbf{x}^{(2)T} \\ \vdots & \vdots \\ 1 & \mathbf{x}^{(N)T} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_D^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_D^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \cdots & x_D^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times D+1}$$

is the *design matrix*

- $\hat{\mathbf{y}} = [y^{(1)}, \dots, y^{(N)}]^T \in \mathbb{R}^N$  is the *target vector*

# General Recipe for Machine Learning

1. Define a model and model parameters
2. Write down an objective function
3. Optimize the objective w.r.t. the model parameters

# Recipe for Linear Regression

1. Define a model and model parameters

Assume  $y = w^T x$   
→ parameters  $w = [w_0, w_1, \dots, w_D]$

2. Write down an objective function

Minimize squared loss  
$$l_D(w) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - w^T x^{(n)})^2$$

3. Optimize the objective w.r.t. the model parameters

1. Solve in closed-form: take derivatives, set them equal to 0 and solve!

## Minimizing the Squared Error

gradient

$$l_D(w) = \frac{1}{N} \sum_{n=1}^N (w^T x^{(n)} - y^{(n)})^2 = \frac{1}{N} \sum_{n=1}^N (x^{(n)T} w - y^{(n)})^2$$

$$Xw - y \in \mathbb{R}^N$$

$$\begin{aligned} l_D(w) &= \frac{1}{N} (Xw - y)^T (Xw - y) = \|Xw - y\|_2^2 \\ &= \frac{1}{N} (w^T X^T X w - 2w^T X^T y + y^T y) \end{aligned}$$

$$\nabla_w l_D(w) = \frac{1}{N} (2X^T X w - 2X^T y)$$

$$\Rightarrow \frac{1}{N} (2X^T X \hat{w} - 2X^T y) = 0$$

$$\begin{aligned} \Rightarrow X^T X \hat{w} - X^T y &= 0 \Rightarrow X^T X \hat{w} = X^T y \\ \Rightarrow \hat{w} &= (X^T X)^{-1} X^T y \end{aligned}$$

$$H_w l_D(w) = \frac{1}{N} (2X^T X)$$

# Closed Form Solution

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

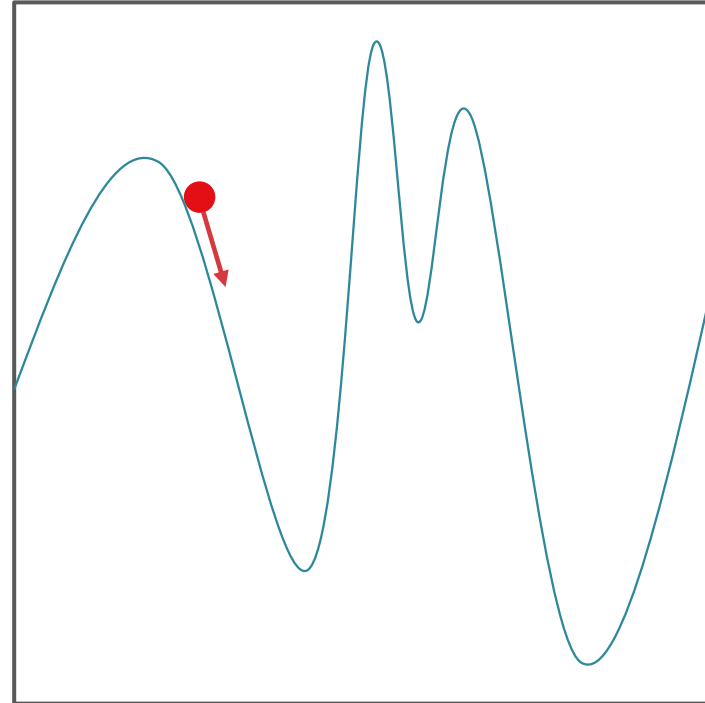
1. Is  $\mathbf{X}^T \mathbf{X}$  invertible?

2. If so, how computationally expensive is inverting  $\mathbf{X}^T \mathbf{X}$ ?

$\mathbf{X} \in \mathbb{R}^{N \times (D+1)} \Rightarrow \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{(D+1) \times (D+1)}$   
classically inverting is  $O(D^3)$  (but we can get  $O(D^{2.373})$ )  
we need to store  $\mathbf{X}$ ,  $O(ND)$

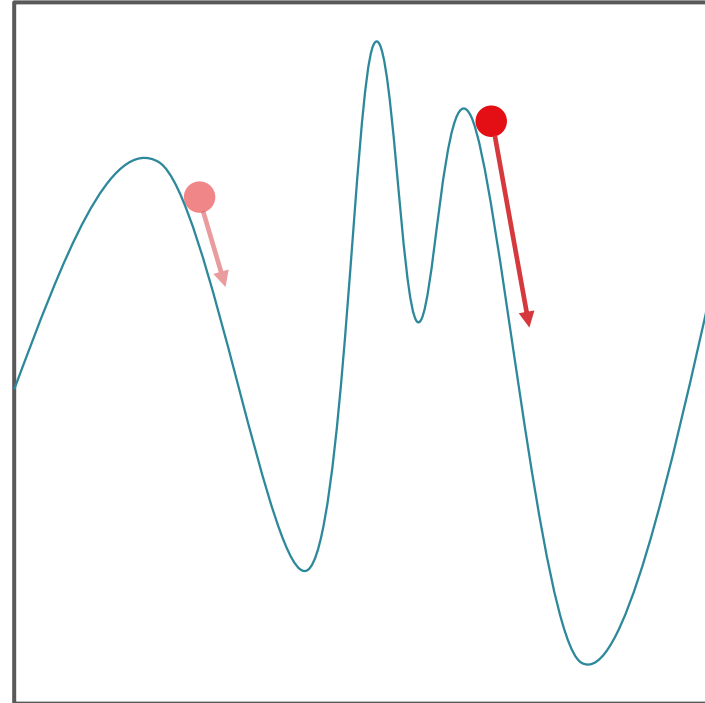
# Gradient Descent: Intuition

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



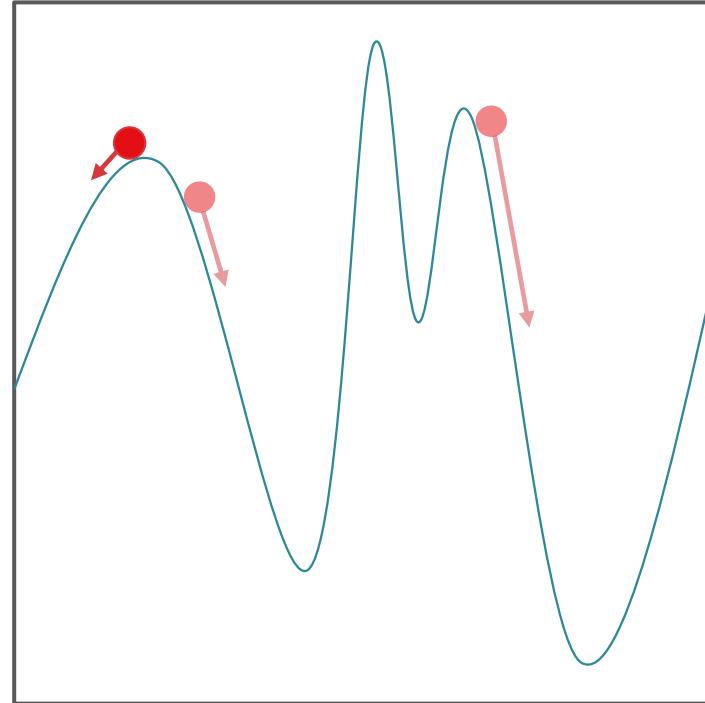
# Gradient Descent: Intuition

- 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

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


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

# Gradient Descent: Step Direction

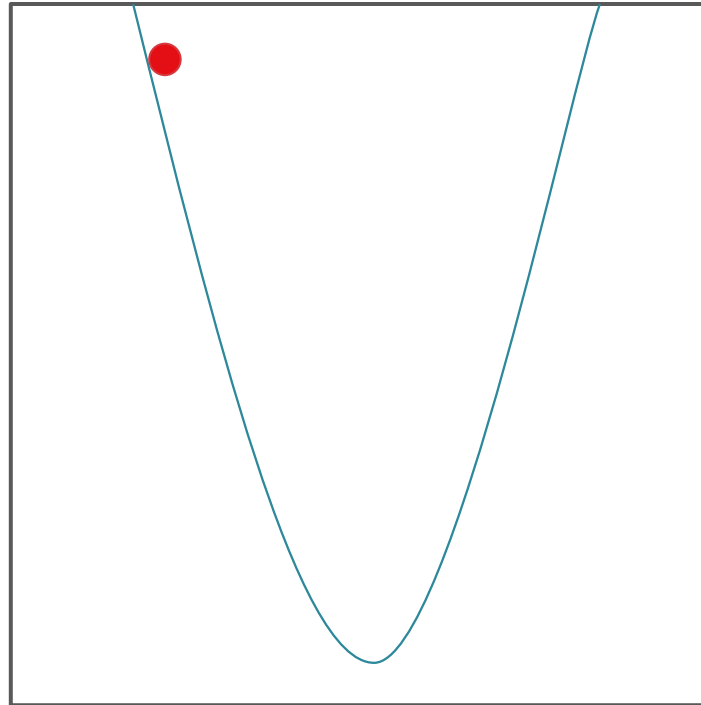
- Suppose the current weight vector is  $\mathbf{w}^{(t)}$
- Move some distance,  $\eta$ , in the “most downhill” direction,  $\hat{\mathbf{v}}$ :

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

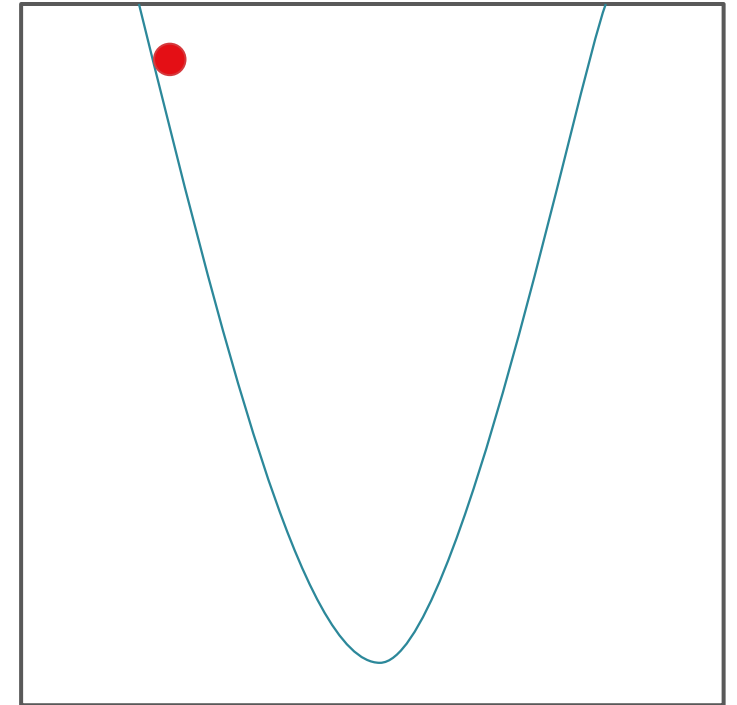
- The gradient points in the direction of steepest *increase* ...
- ... so  $\hat{\mathbf{v}}$  should point in the opposite direction:

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


# Gradient Descent: Step Size

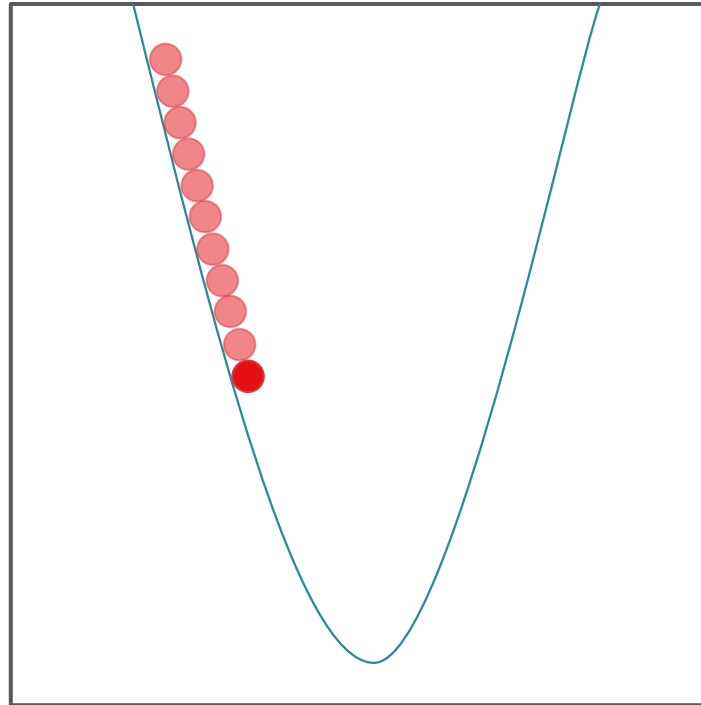


Small  $\eta$

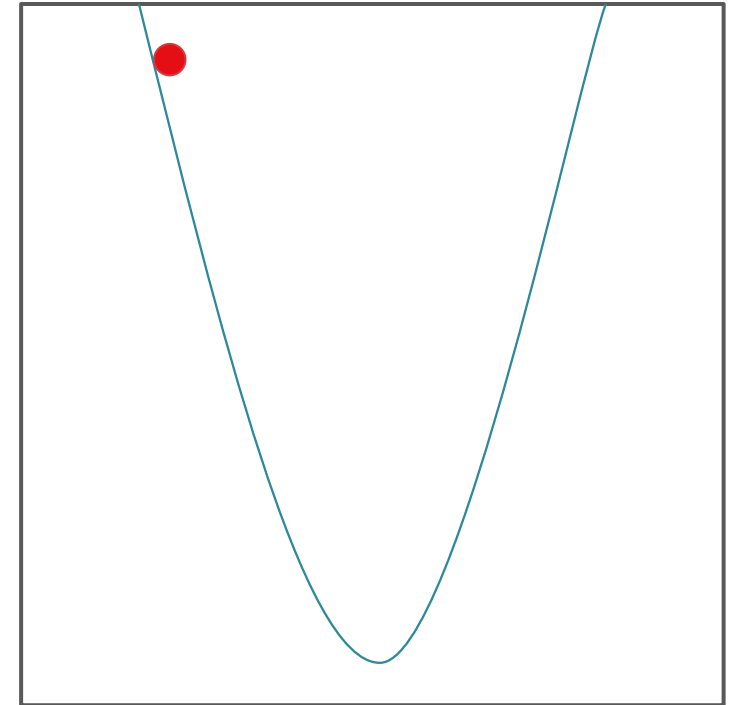


Large  $\eta$

# Gradient Descent: Step Size

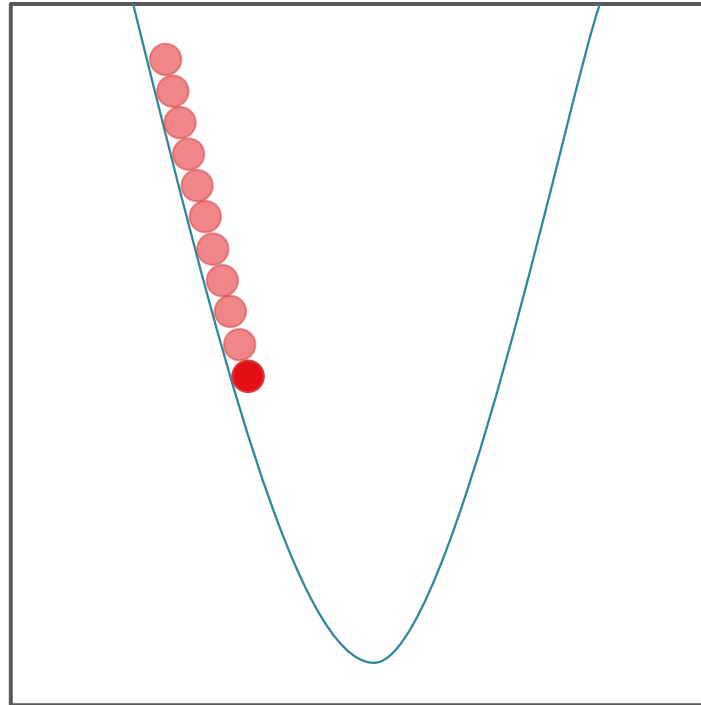


Small  $\eta$

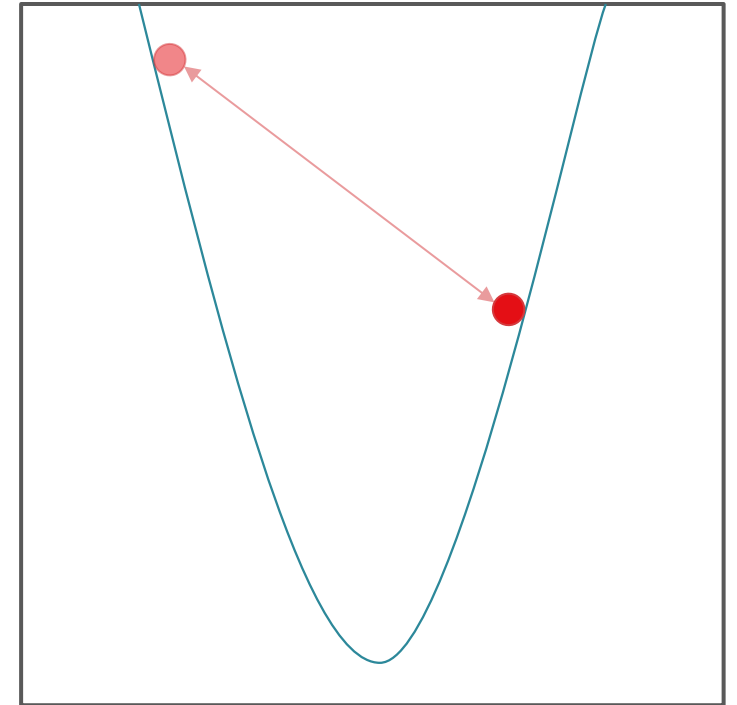


Large  $\eta$

# Gradient Descent: Step Size



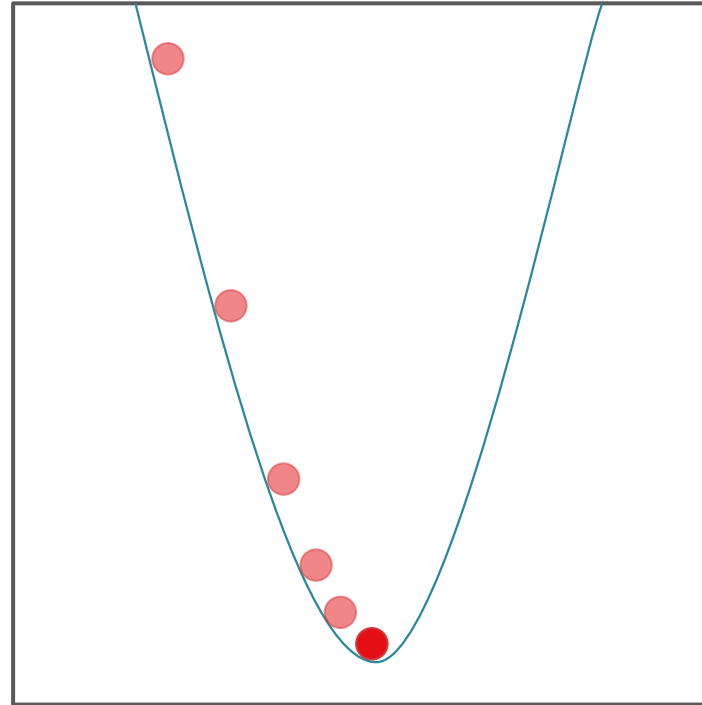
Small  $\eta$



Large  $\eta$

# Gradient Descent: Step Size

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



- 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  
→  $\eta^{(t)}$  (hopefully) decreases over time

# Gradient Descent

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

- $\eta^{(t)} = \eta^{(0)} \|\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})\|$

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

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

# 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:  $\mathbf{w}^{(t)}$

# Why Gradient Descent for linear regression?

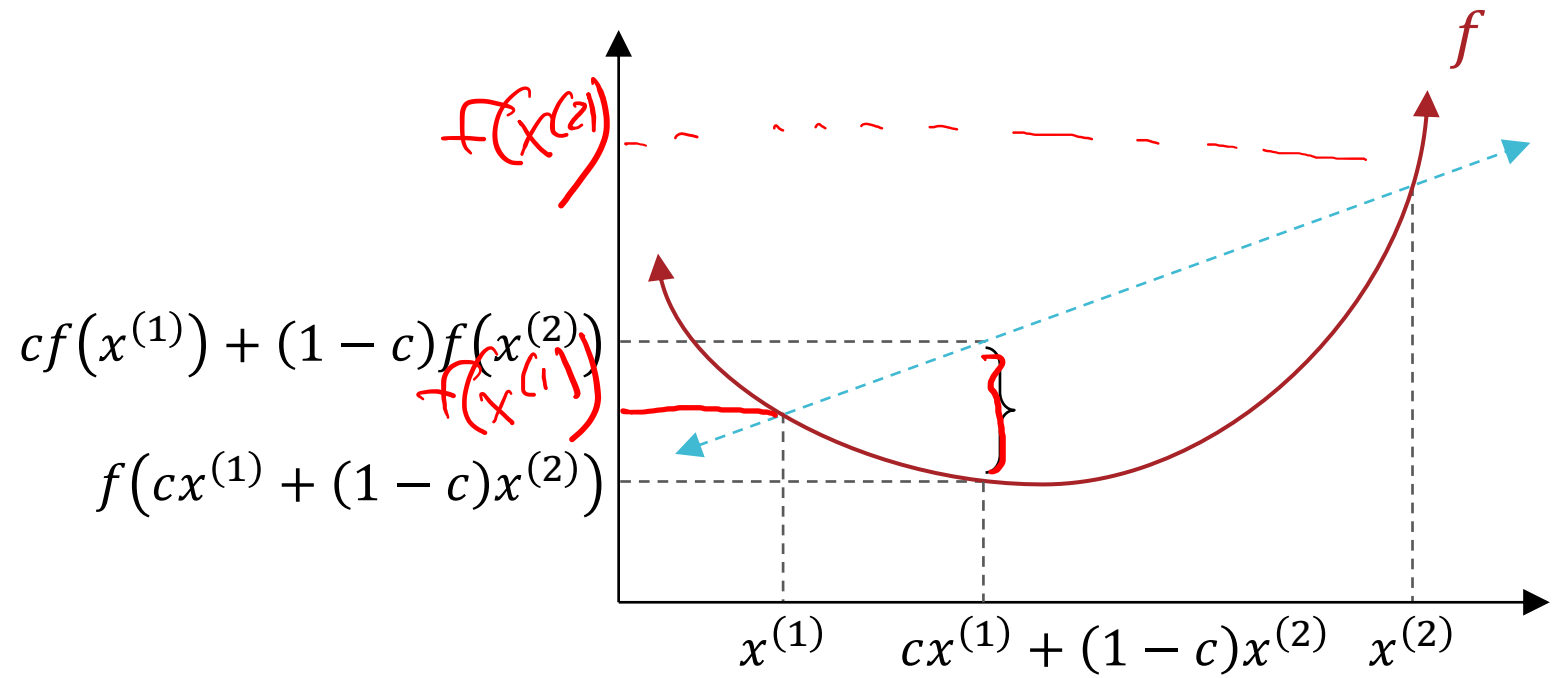
- Input:  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \eta^{(0)}, T$ 
  1. Initialize  $\mathbf{w}^{(0)}$  to all zeros and set  $t = 0$
  2. While TERMINATION CRITERION is not satisfied
    - a. Compute the gradient:
$$\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)}) = \frac{1}{N}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

# Convexity

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

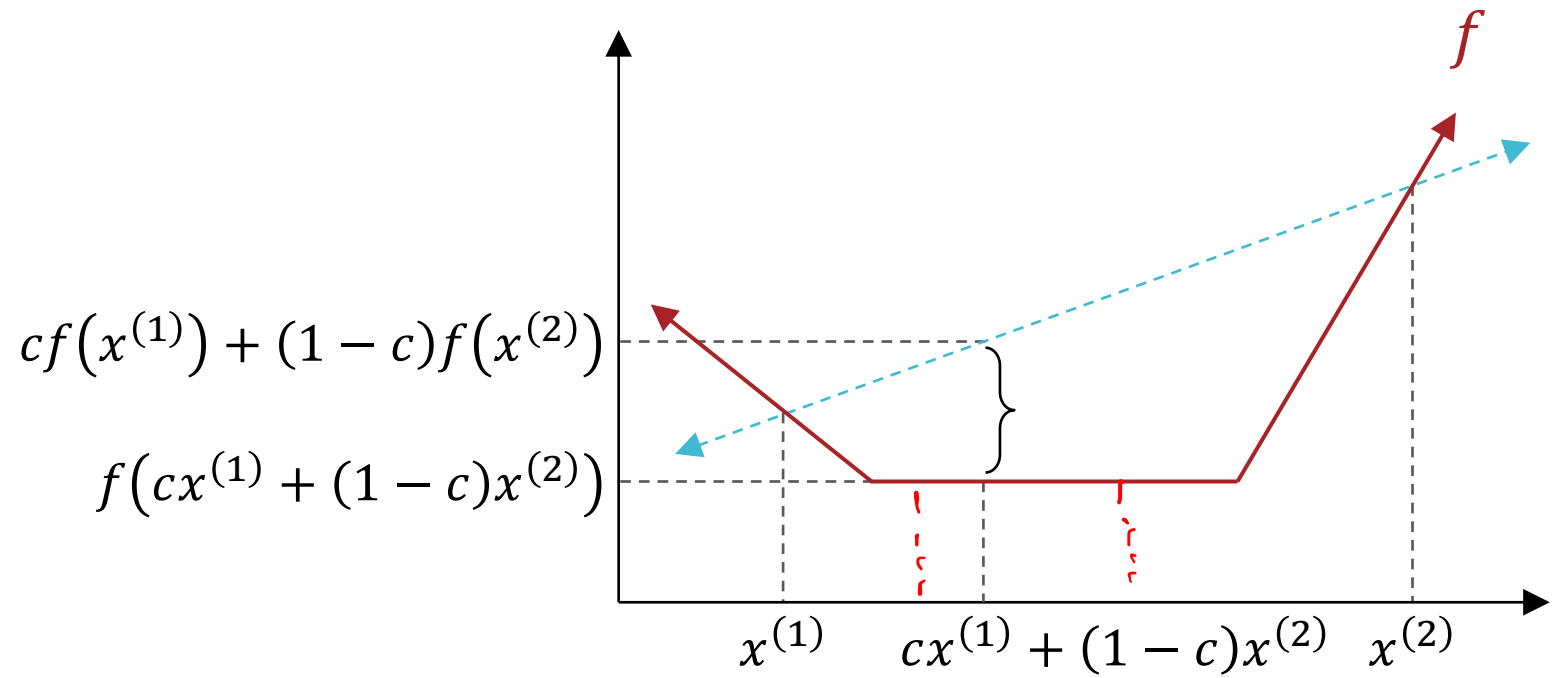
$$\forall \mathbf{x}^{(1)} \in \mathbb{R}^D, \mathbf{x}^{(2)} \in \mathbb{R}^D \text{ and } 0 \leq c \leq 1$$

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



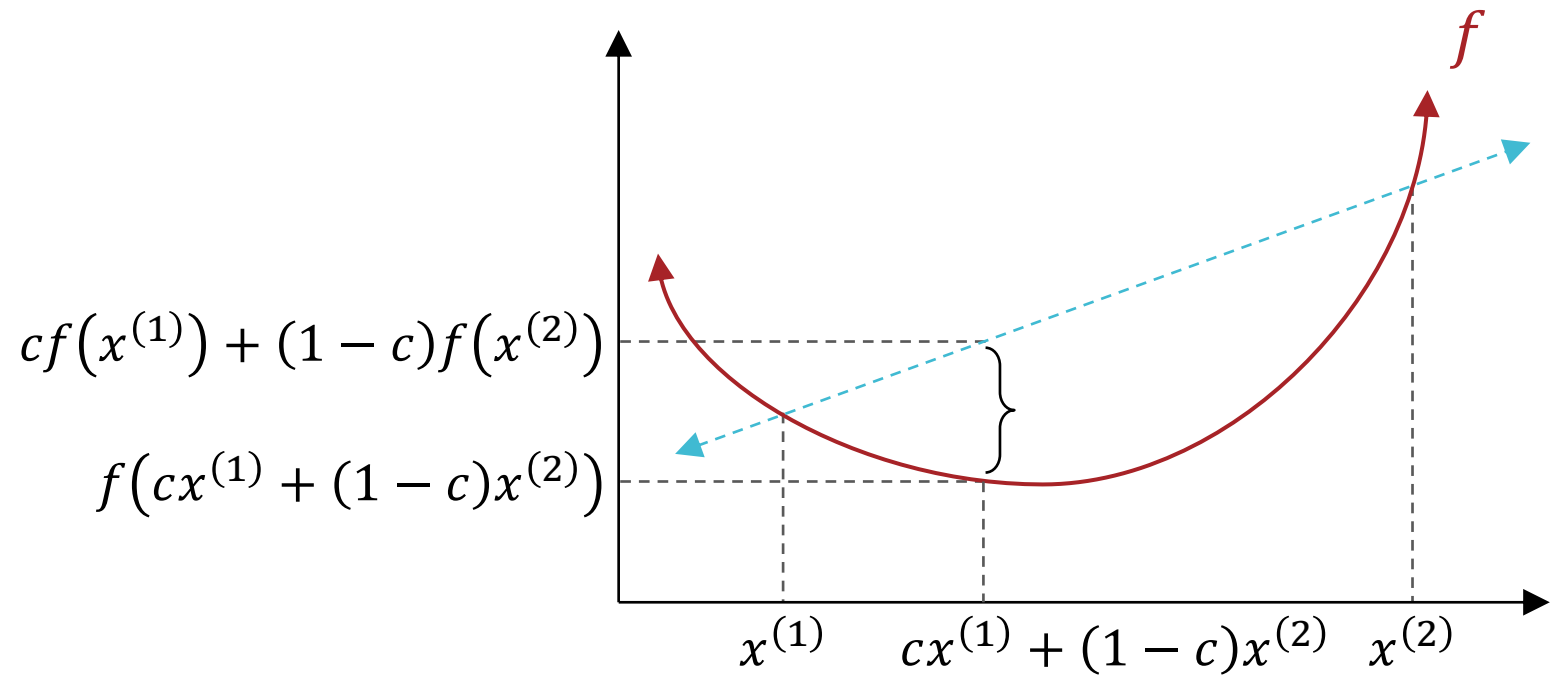
# Convexity

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

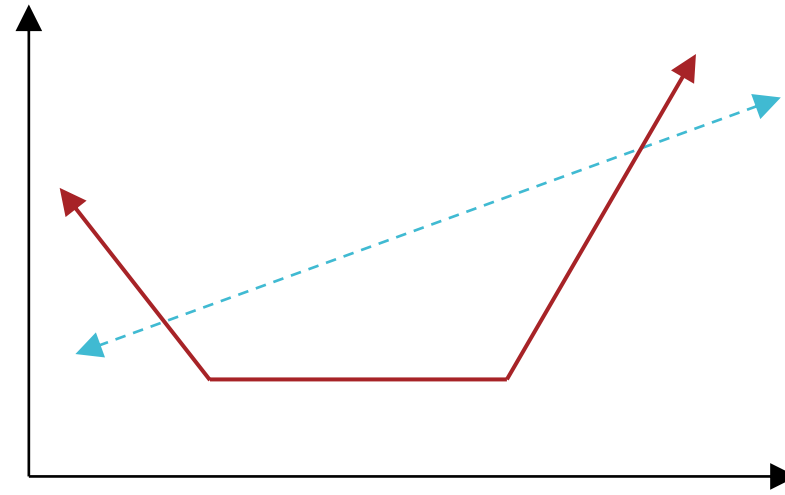


# Convexity

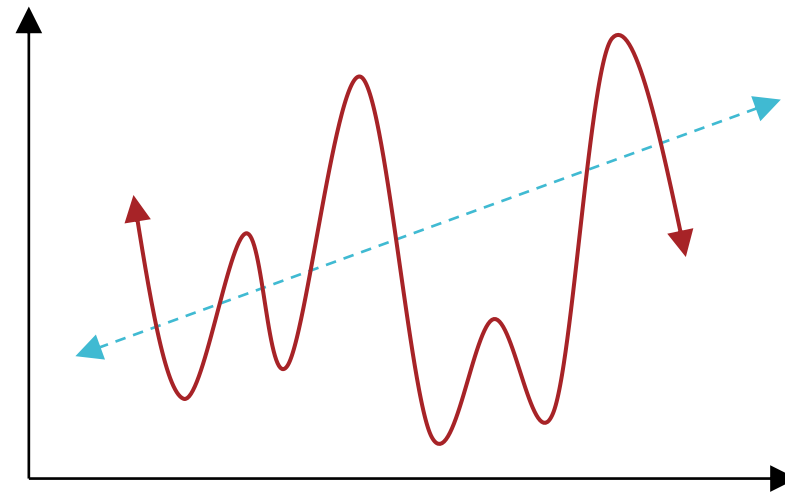
- A function  $f: \mathbb{R}^D \rightarrow \mathbb{R}$  is *strictly convex* if  
 $\forall \mathbf{x}^{(1)} \in \mathbb{R}^D, \mathbf{x}^{(2)} \in \mathbb{R}^D$  and  $0 < c < 1$   
 $f(c\mathbf{x}^{(1)} + (1 - c)\mathbf{x}^{(2)}) < cf(\mathbf{x}^{(1)}) + (1 - c)f(\mathbf{x}^{(2)})$



# Convexity

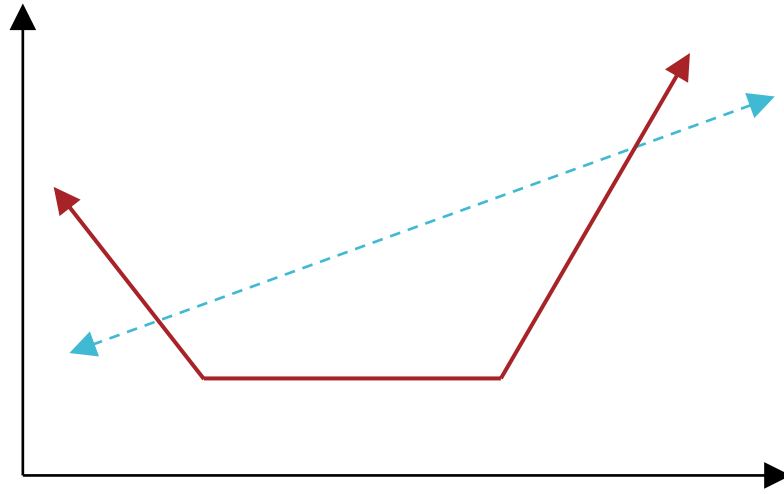


Convex functions



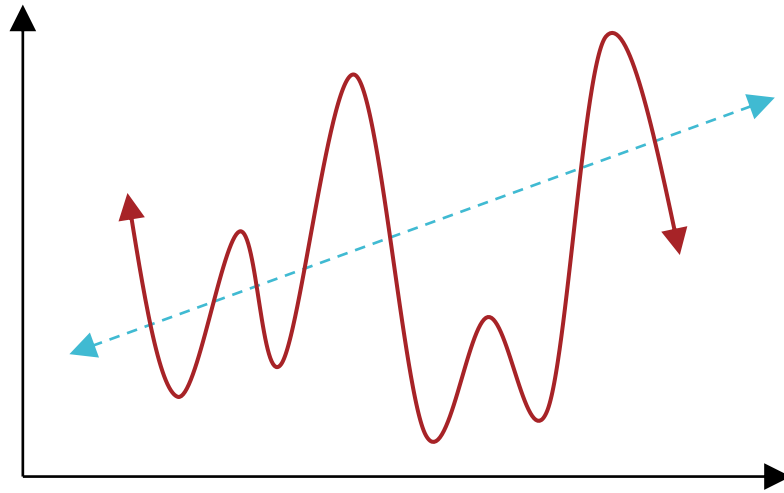
Non-convex functions

# Convexity



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

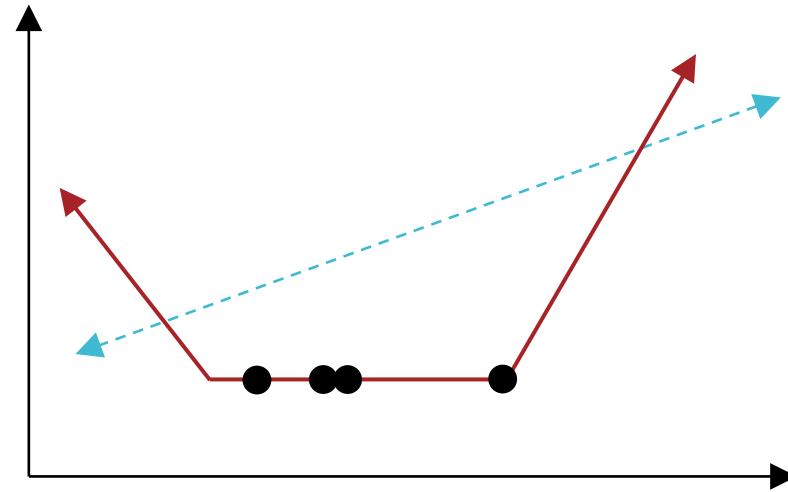
- $\mathbf{x}^*$  is a *global* minimum iff  $f(\mathbf{x}^*) \leq f(\mathbf{x}) \forall \mathbf{x} \in \mathbb{R}^D$



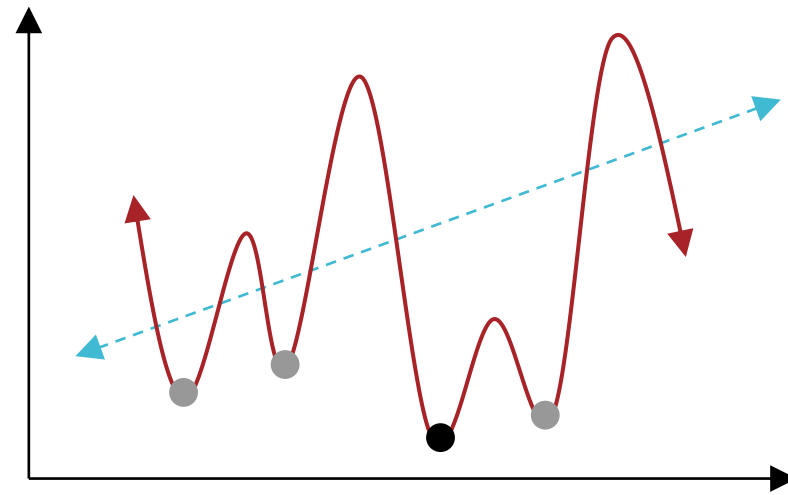
- $\mathbf{x}^*$  is a *local* minimum iff  $\exists \epsilon$  s.t.  $f(\mathbf{x}^*) \leq f(\mathbf{x}) \forall \mathbf{x}$  s.t.  $\|\mathbf{x} - \mathbf{x}^*\|_2 < \epsilon$



# Convexity

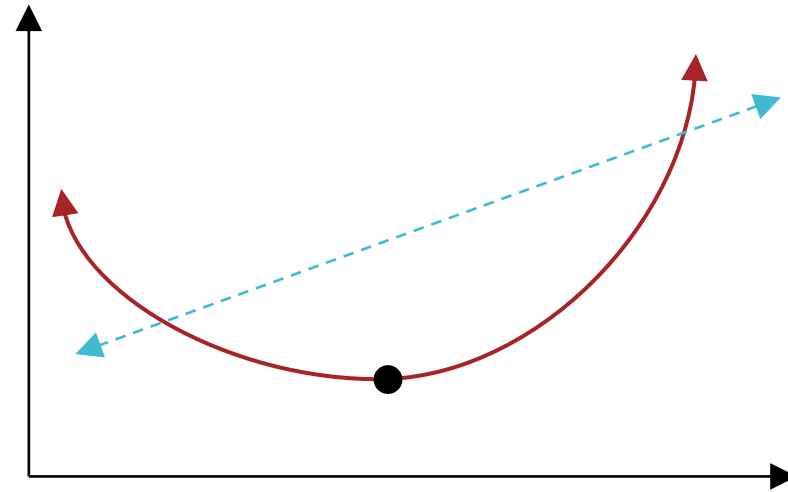


Convex functions:  
Each local minimum is a  
global minimum!

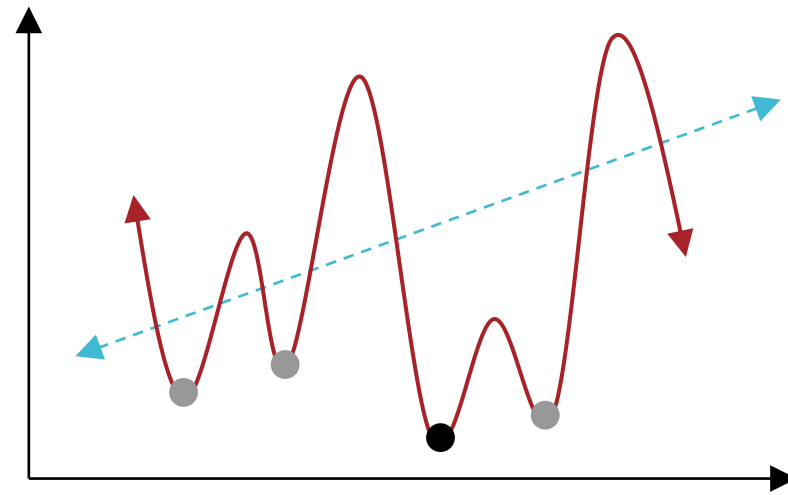


Non-convex functions:  
A local minimum may or may  
not be a global minimum...

# Convexity



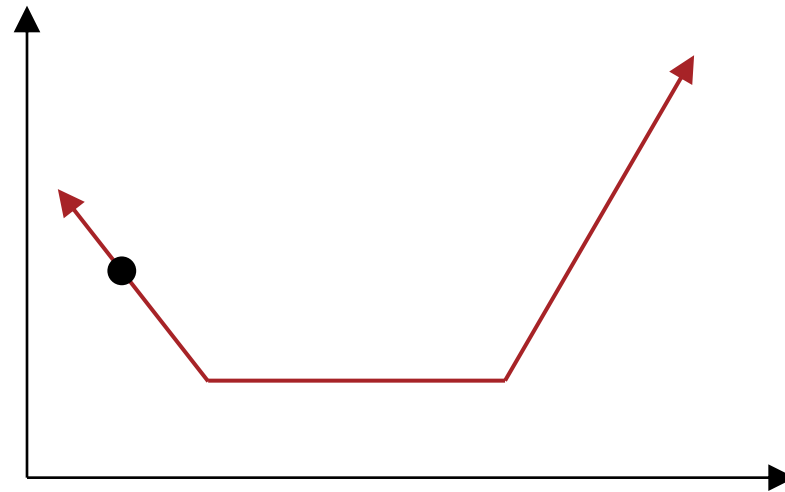
Strictly convex functions:  
There exists a unique global minimum!



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

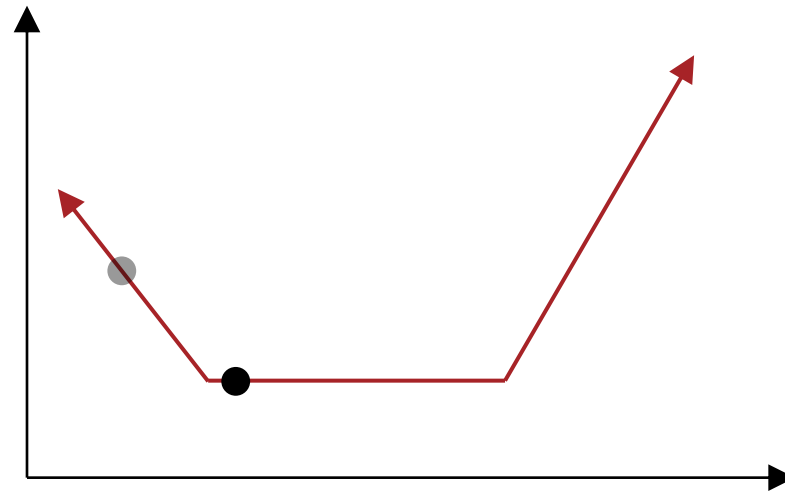
# Gradient Descent & Convexity

- 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!



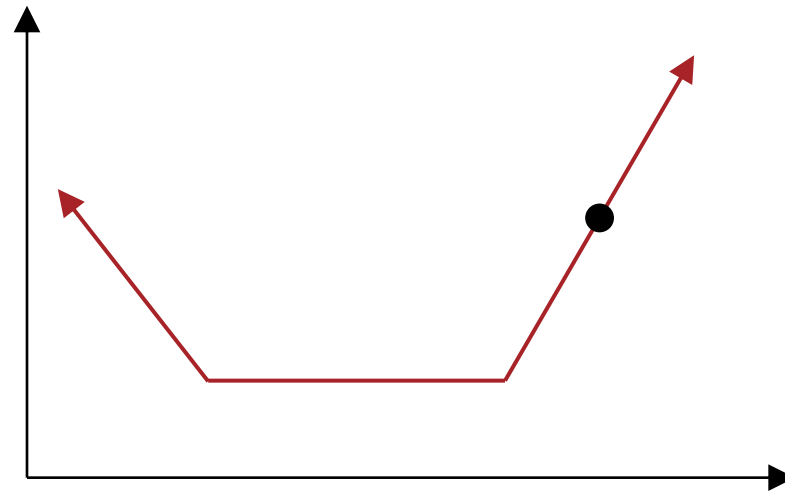
# Gradient Descent & Convexity

- 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!



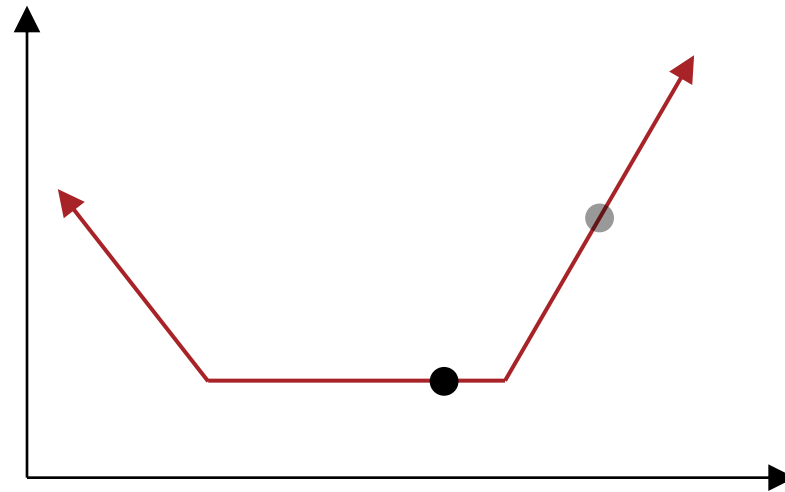
# Gradient Descent & Convexity

- 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!



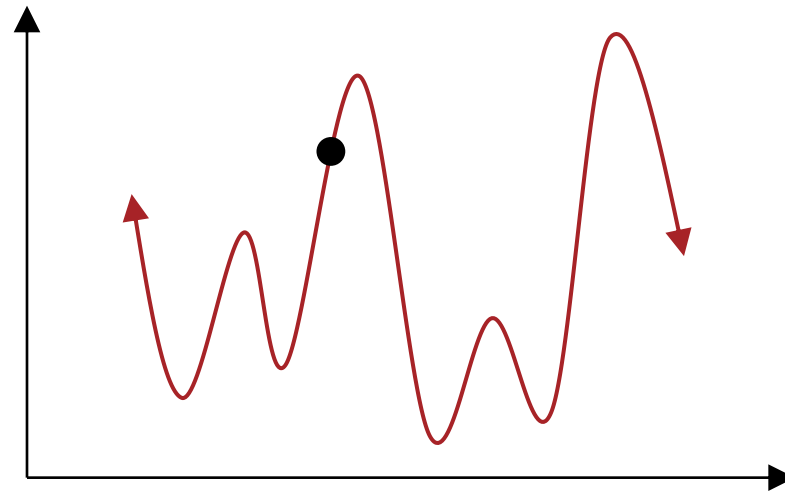
# Gradient Descent & Convexity

- 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!



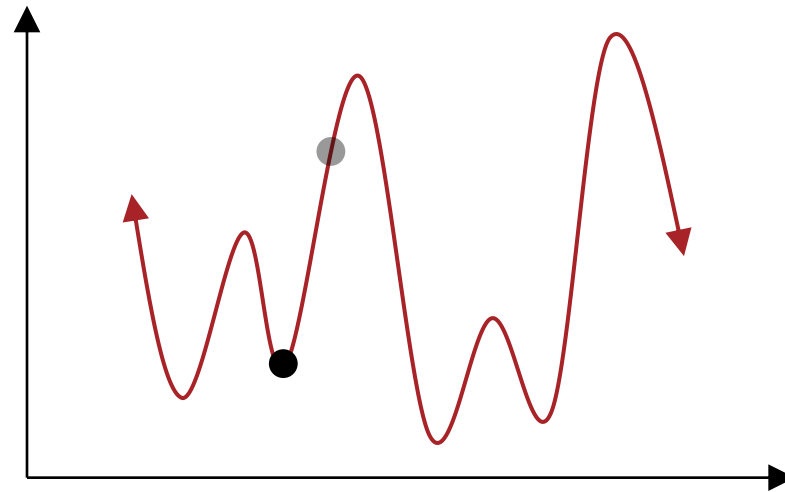
# Gradient Descent & Convexity

- 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

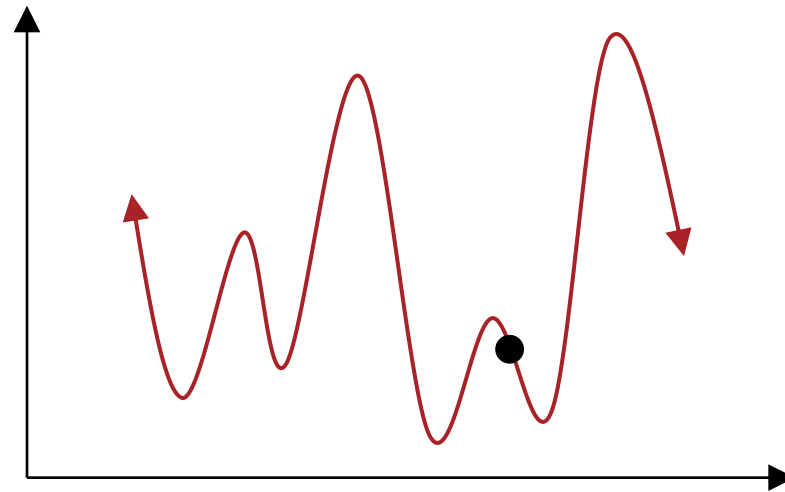
- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
  - Not ideal if the objective function is non-convex...





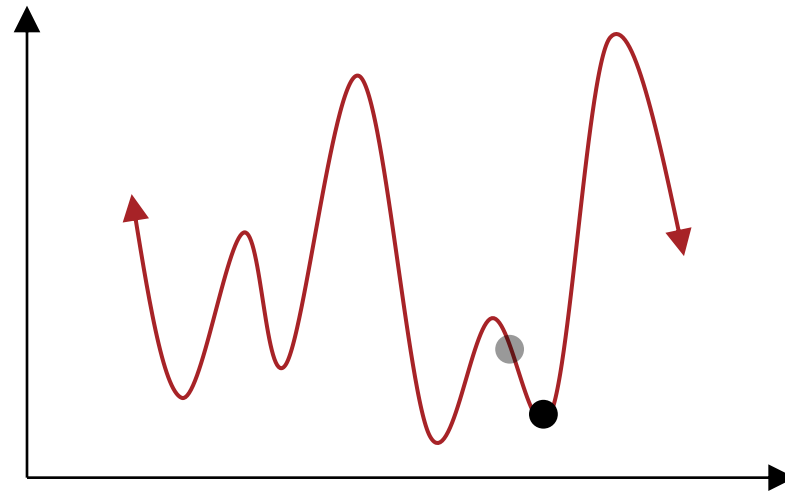
# Gradient Descent & Convexity

- 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...



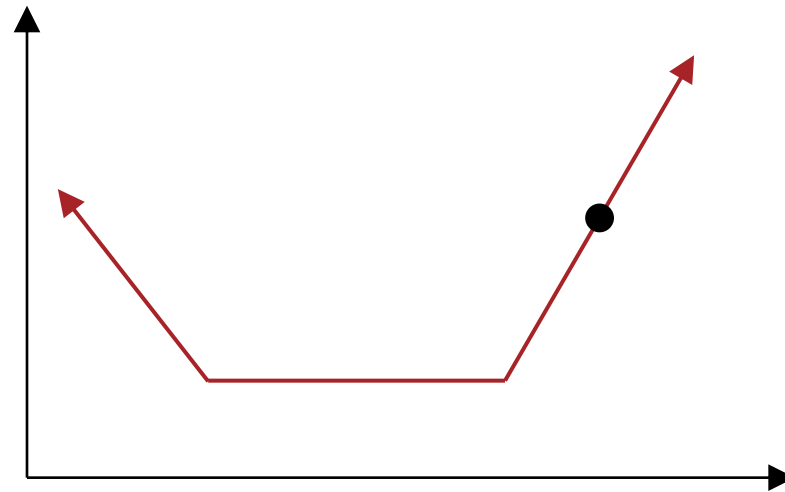
# Gradient Descent & Convexity

- 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 squared error for linear regression is convex (but not strictly convex)!

- 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!



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

$$H_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}) = 2X^T X \text{ which is positive } \textit{semi-definite}$$

# 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?

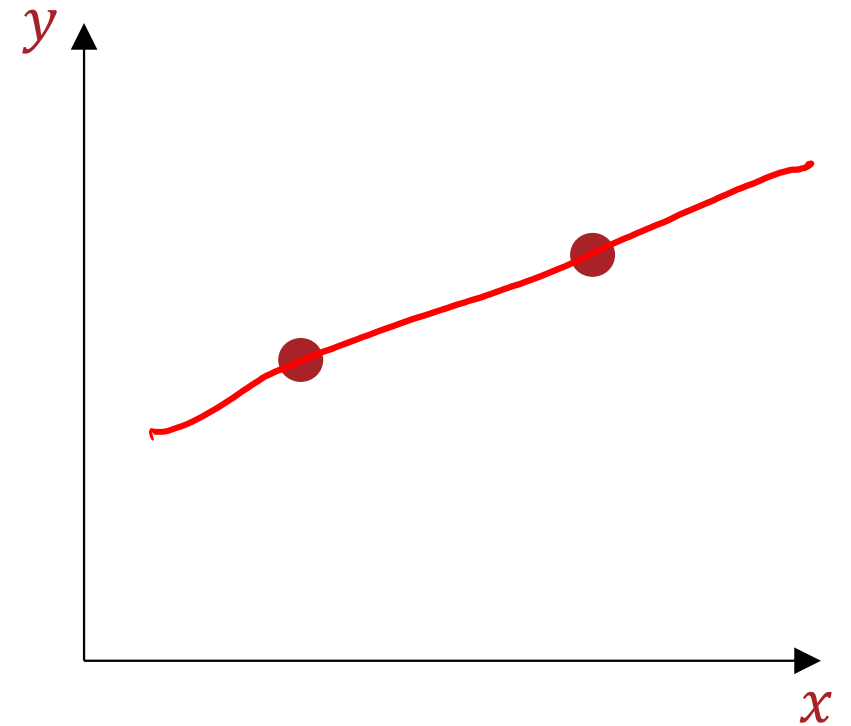
Almost always (as long as features are linearly independent)

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!

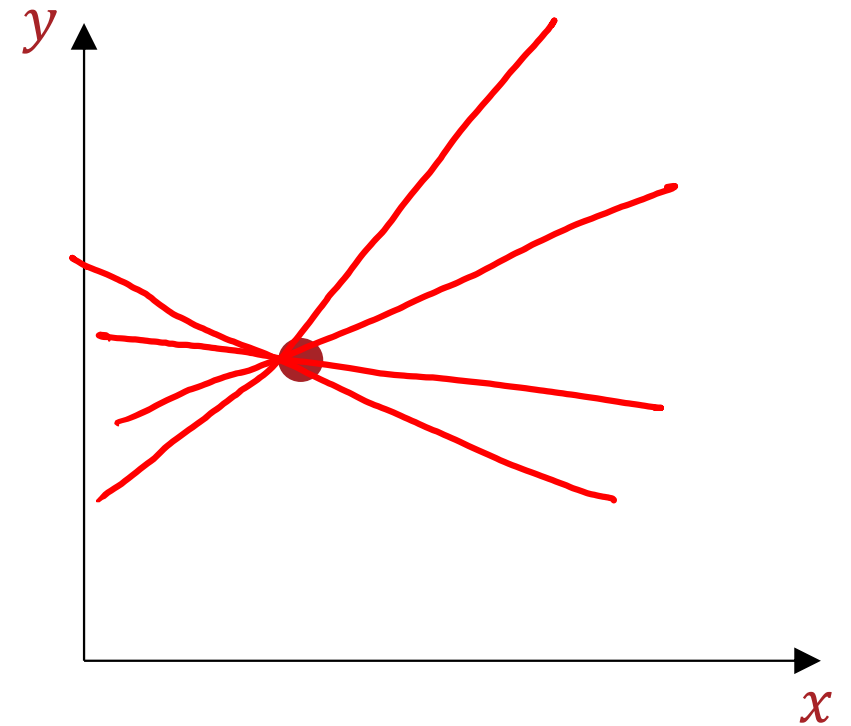
# Linear Regression: Uniqueness

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



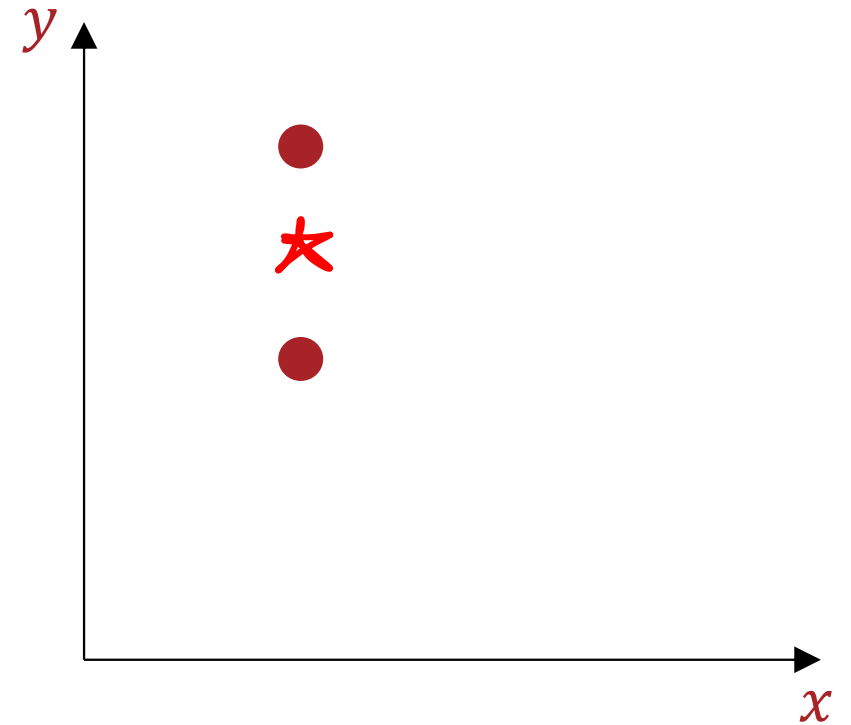
# Linear Regression: Uniqueness

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



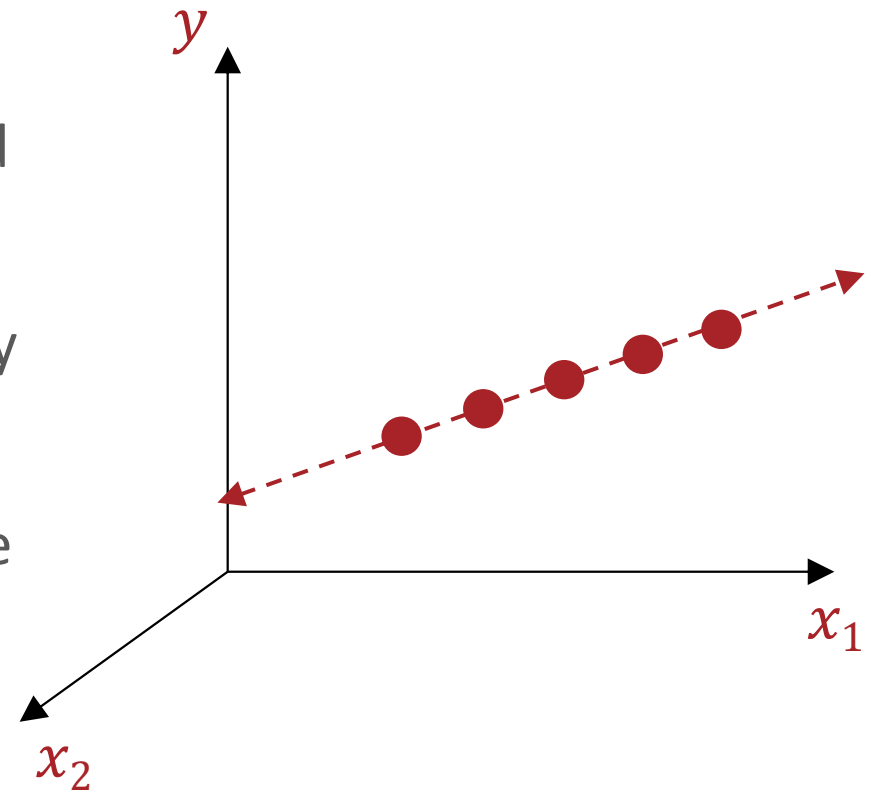
# Linear Regression: Uniqueness

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



# Linear Regression: Uniqueness

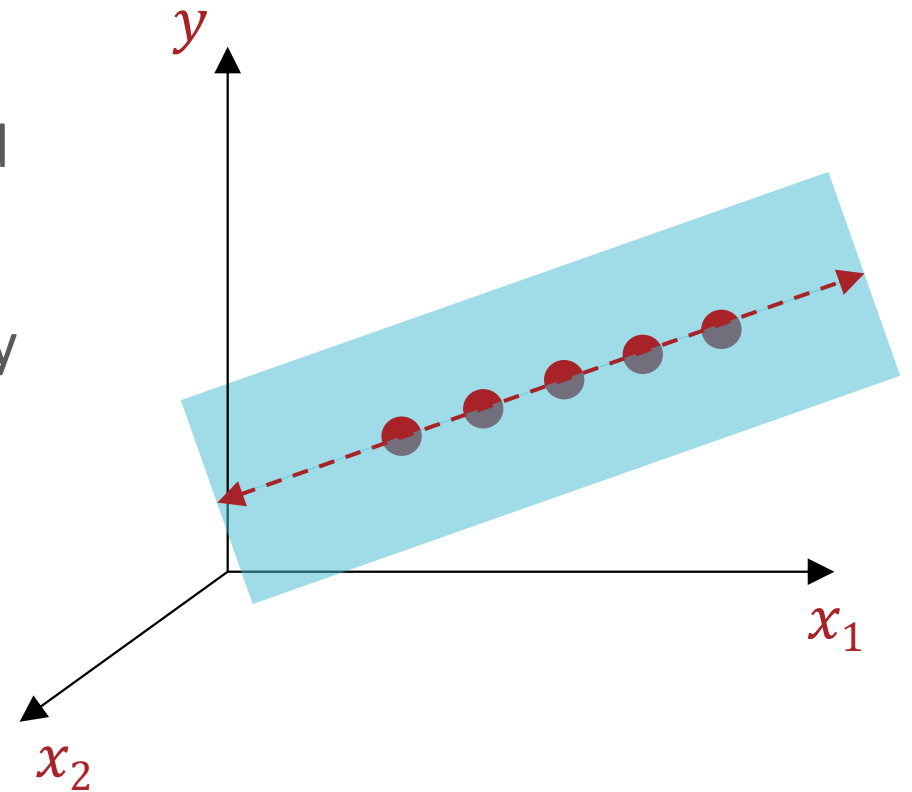
- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters  $\theta$ ) are there for the given dataset?





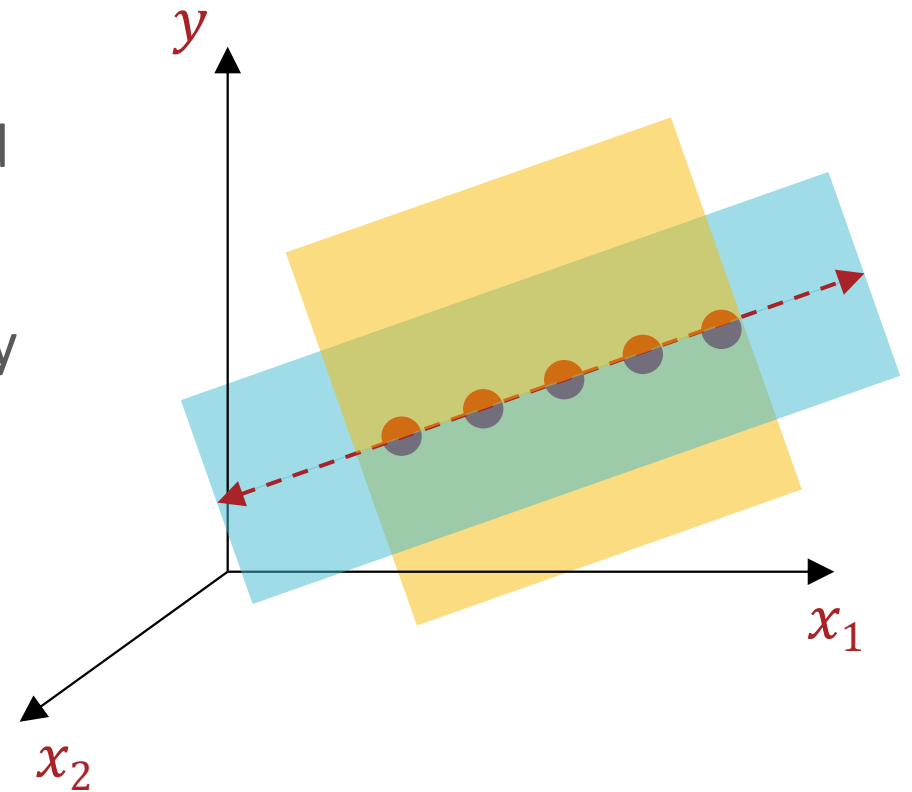
# Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of weights  $\mathbf{w}$ ) 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 weights  $\mathbf{w}$ ) are there for the given dataset?



# Key Takeaways

- Closed form solution for linear regression
  - Setting the gradient equal to 0 and solving for critical points
  - Potential issues: invertibility and computational costs
- Gradient descent
  - Effect of step size
  - Termination criteria
- Convexity vs. non-convexity
  - Strong vs. weak convexity
  - Implications for local, global and unique optima

# Bias-Variance Tradeoff

- Suppose you have a regression task and your goal is to minimize the *true* squared error:

$$err(h) = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} \left[ (h(\mathbf{x}) - f(\mathbf{x}))^2 \right]$$

where  $f$  is the target function and

$\mathcal{P}$  is some distribution of interest over all possible inputs

- Let  $h_{\mathcal{D}}$  be the hypothesis returned when the input training dataset is  $\mathcal{D}$
- Assume each data point in  $\mathcal{D}$  is drawn independently from  $\mathcal{P}$

# Bias-Variance Tradeoff

- $err(h_{\mathcal{D}}) = \mathbb{E}_{x \sim \mathcal{P}} [(h_{\mathcal{D}}(x) - f(x))^2]$

- $\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})]$

$$= \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{x \sim \mathcal{P}} [(h_{\mathcal{D}}(x) - f(x))^2]]$$

$$= \mathbb{E}_{x \sim \mathcal{P}} [\mathbb{E}_{\mathcal{D}} [(h_{\mathcal{D}}(x) - f(x))^2]]_{f(x)}$$

$$= \mathbb{E}_{x \sim \mathcal{P}} [\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(x)^2 - 2h_{\mathcal{D}}(x)f(x) + f(x)^2]]$$

$$= \mathbb{E}_{x \sim \mathcal{P}} [\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(x)^2] - 2\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(x)]f(x) + f(x)^2]$$

$$\bar{h}(x) = \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(x)] \approx \frac{1}{C} \sum_{i=1}^C h_{\mathcal{D}_i}(x)$$

# Bias-Variance Tradeoff

- $\mathbb{E}_{\mathcal{D}}[\text{err}(h_{\mathcal{D}})]$   
$$= \mathbb{E}_{x \sim p} [\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(x)^2] - 2\bar{h}(x)f(x) + f(x)^2]$$
$$= \mathbb{E}_{x \sim p} [\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(x)^2] - \bar{h}(x)^2 + \bar{h}(x)^2 - 2\bar{h}(x)f(x) + f(x)^2]$$
$$= \mathbb{E}_{x \sim p} [\underbrace{\mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(x)^2 - \bar{h}(x)^2]}_{\text{variance of } h_{\mathcal{D}}(x)} + \underbrace{(\bar{h}(x) - f(x))^2}_{(\text{bias of } \bar{h}(x))^2}]$$

# Bias-Variance Tradeoff

How variable is  $h_{\mathcal{D}}$ ?

$$\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})] = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} \left[ \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})^2 - \bar{h}(\mathbf{x})^2] + (\bar{h}(\mathbf{x}) - f(\mathbf{x}))^2 \right]$$

How well, on average, does  $h_{\mathcal{D}}$  approximate  $f$ ?

# Bias-Variance Tradeoff

How well could  $h_{\mathcal{D}}$  approximate anything?

$$\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})] = \mathbb{E}_{x \sim \mathcal{P}} \left[ \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)^2 - \bar{h}(x)^2] + (\bar{h}(x) - f(x))^2 \right]$$

How well, on average, does  $h_{\mathcal{D}}$  approximate  $f$ ?



# Bias-Variance Tradeoff

How well could  $h_{\mathcal{D}}$  approximate random noise?

$$\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})] = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} \left[ \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})^2 - \bar{h}(\mathbf{x})^2] + (\bar{h}(\mathbf{x}) - f(\mathbf{x}))^2 \right]$$

How well, on average, does  $h_{\mathcal{D}}$  approximate  $f$ ?

# Bias-Variance Tradeoff

Increases as the model becomes more complex

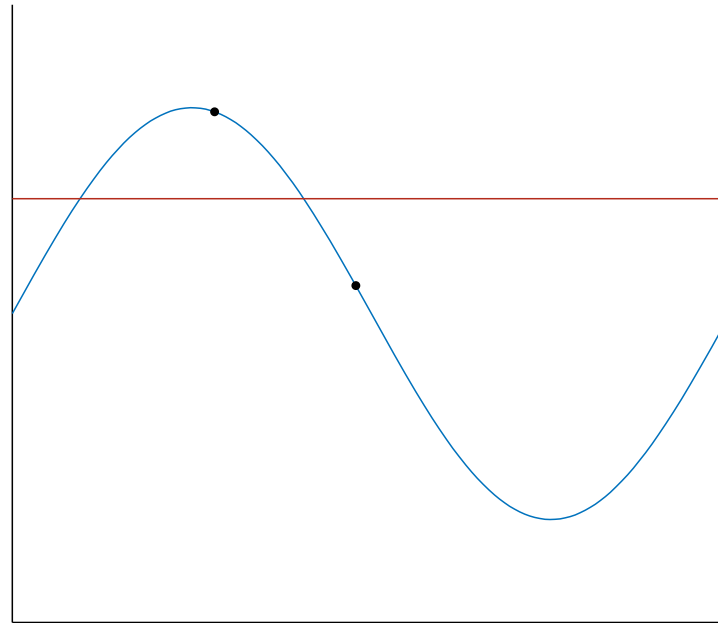
$$\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})] = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} \left[ \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})^2 - \bar{h}(\mathbf{x})^2] + (\bar{h}(\mathbf{x}) - f(\mathbf{x}))^2 \right]$$

Decreases as the model becomes more complex

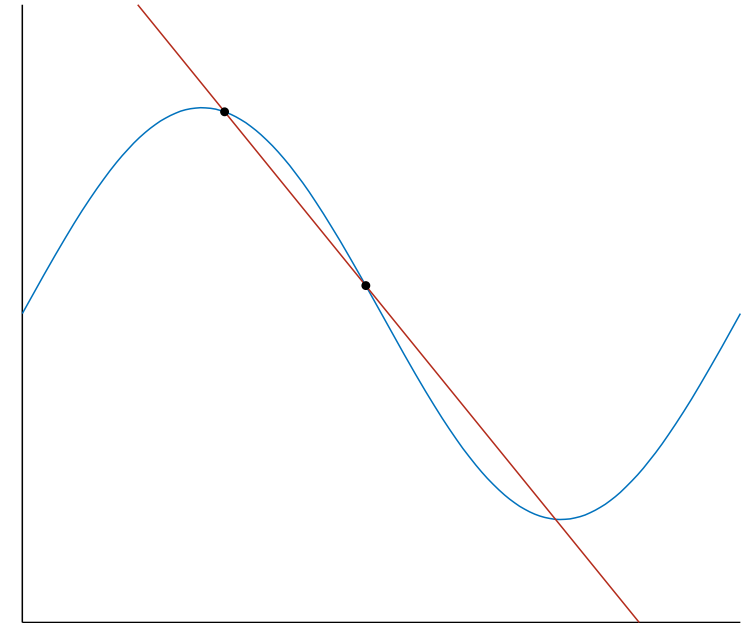
# Bias-Variance Tradeoff (Example)

- $\mathcal{X} = \mathbb{R}$  and  $\mathcal{P} = \text{Uniform}(0, 2\pi)$
- $f(x) = \sin(x)$
- $N = 2 \rightarrow \mathcal{D} = \{(x_1, \sin(x_1)), (x_2, \sin(x_2))\}$
- Consider two models:
  - The “constant” model -  $\mathcal{H}_0 = \{h : h(x) = b\}$
  - Linear regression -  $\mathcal{H}_1 = \{h : h(x) = ax + b\}$

# Bias-Variance Tradeoff (Example)

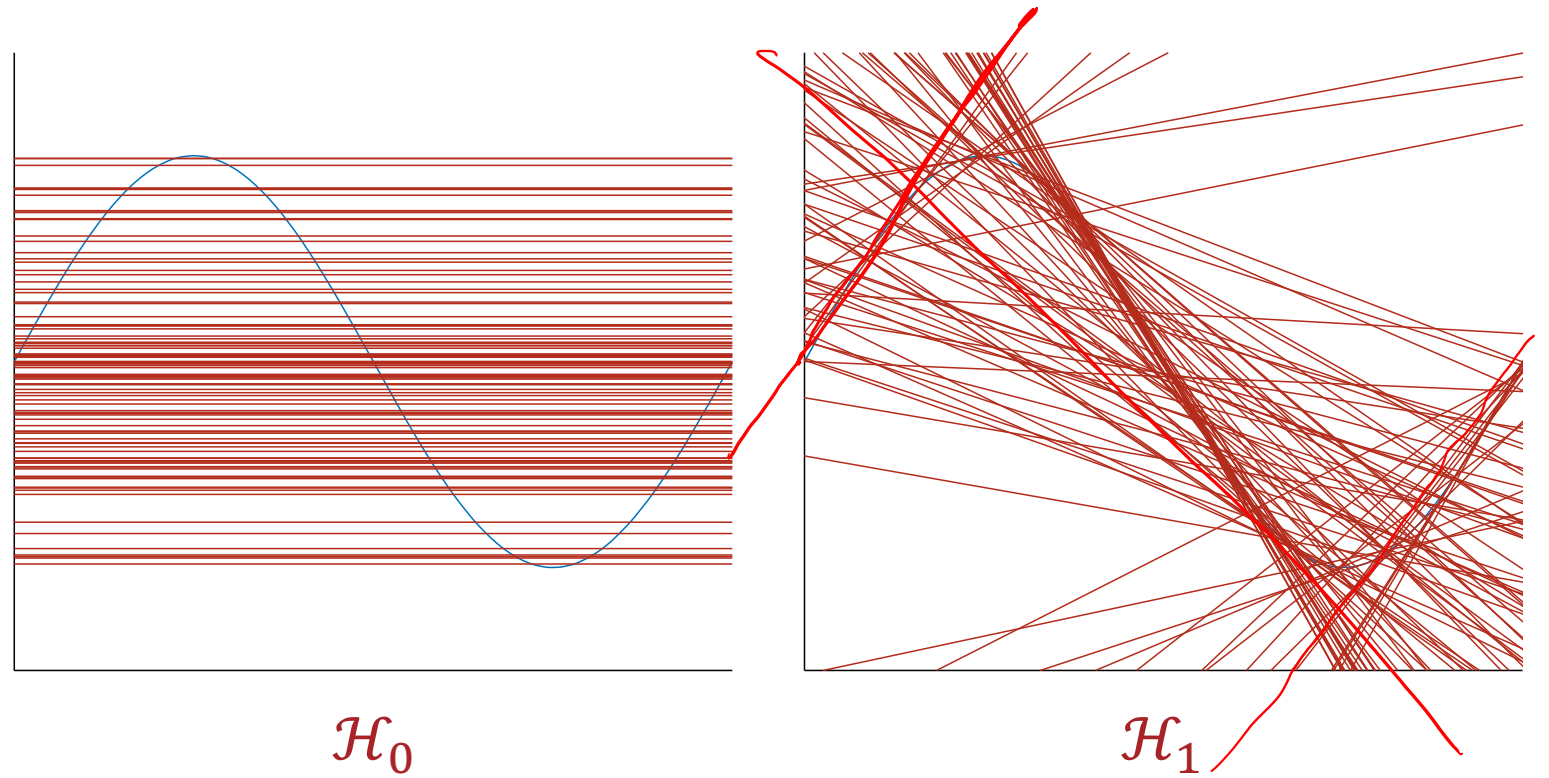


$\mathcal{H}_0$

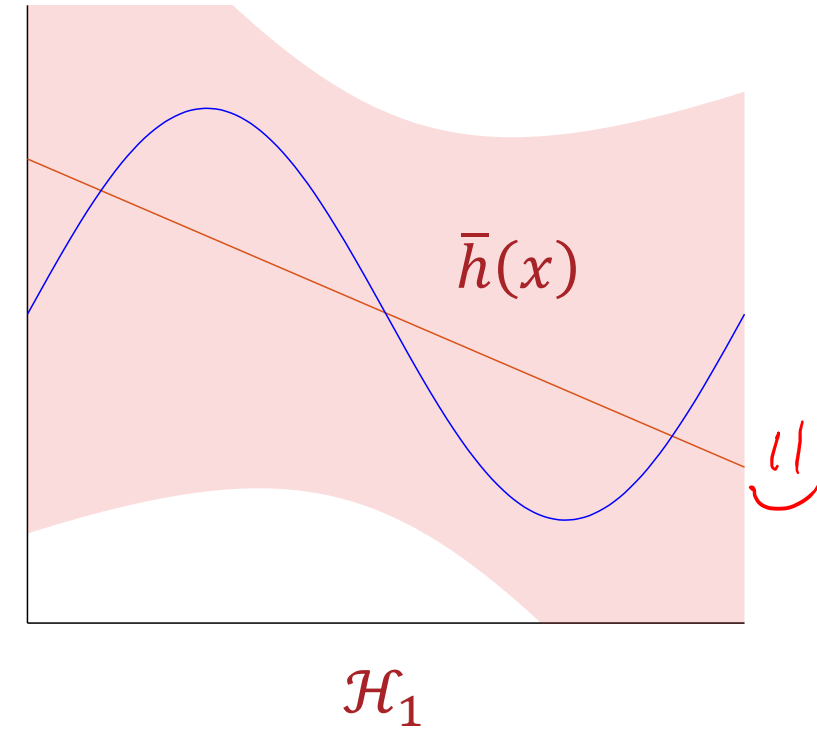
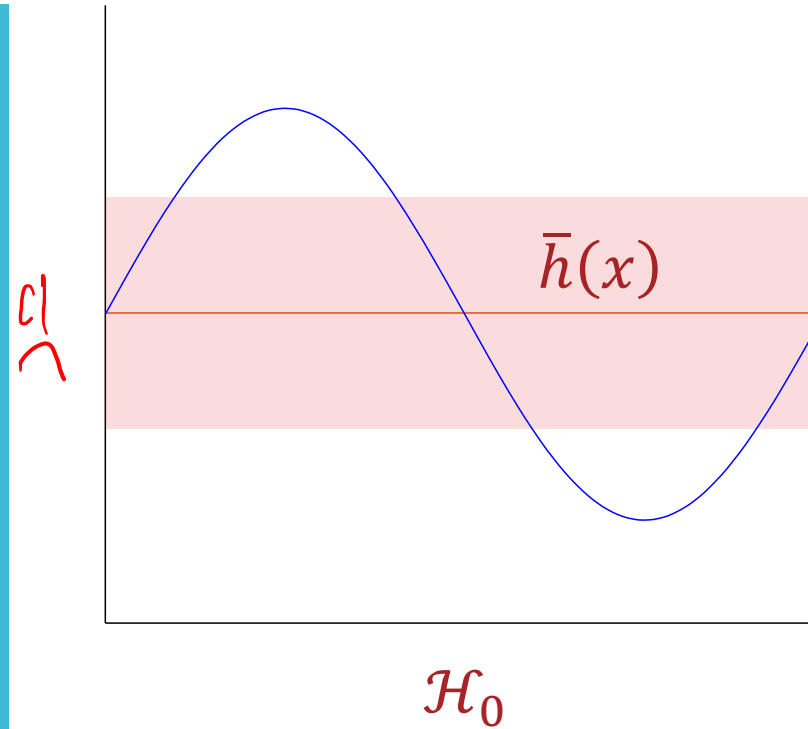


$\mathcal{H}_1$

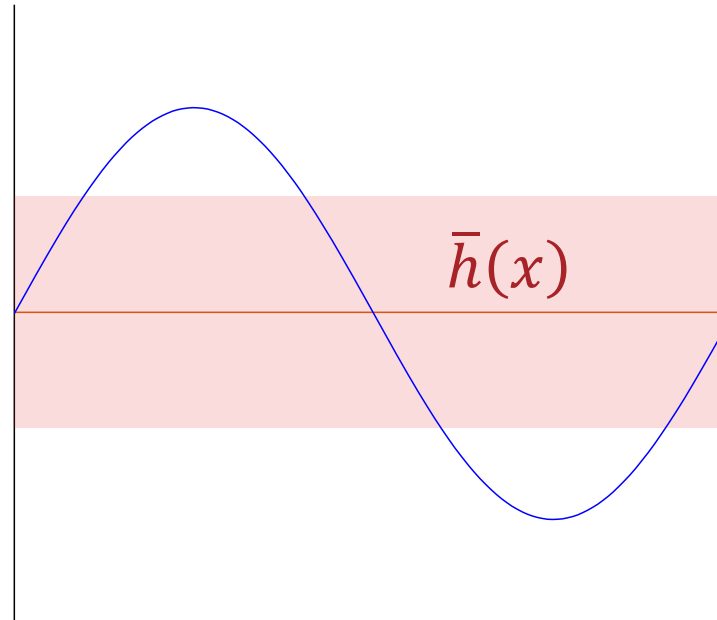
# Bias-Variance Tradeoff (Example)



# Bias-Variance Tradeoff (Example)

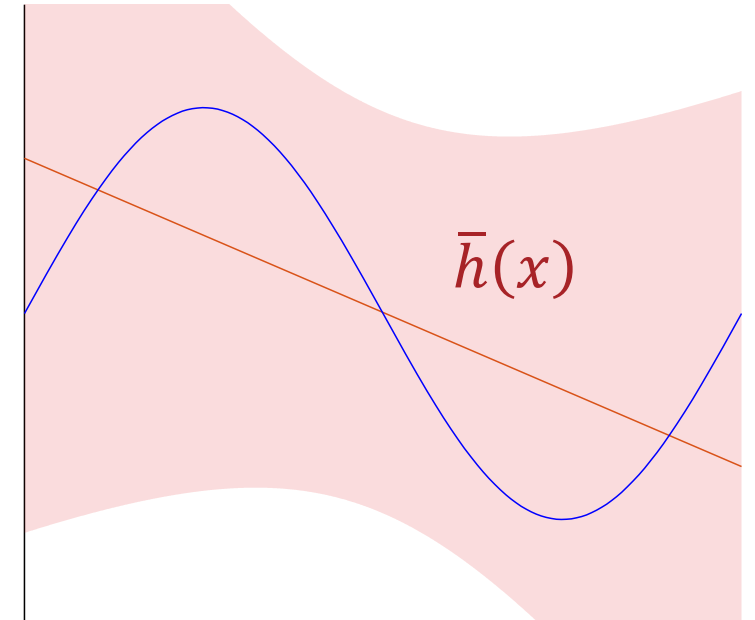


# Bias-Variance Tradeoff ( $N = 2$ )



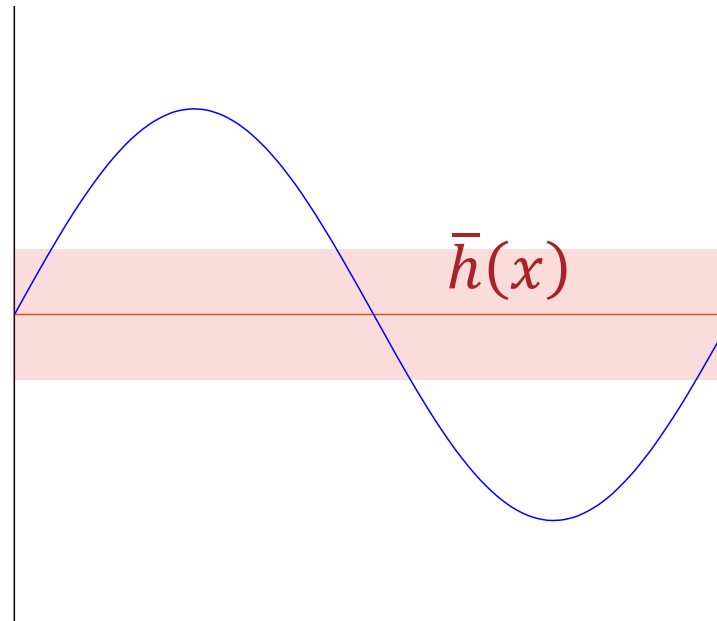
Bias of  $\bar{h}(x) \approx 0.50$   
Variance of  $h_{\mathcal{D}}(x) \approx 0.25$   
 $\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})] \approx 0.75$

$\uparrow$

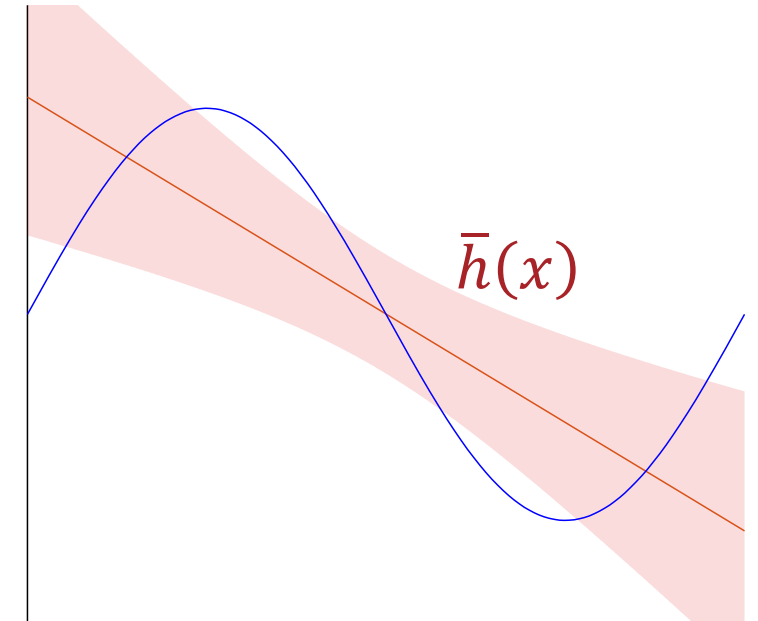


Bias of  $\bar{h}(x) \approx 0.21$   
Variance of  $h_{\mathcal{D}}(x) \approx 1.74$   
 $\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})] \approx 1.95$

# Bias-Variance Tradeoff ( $N = 5$ )

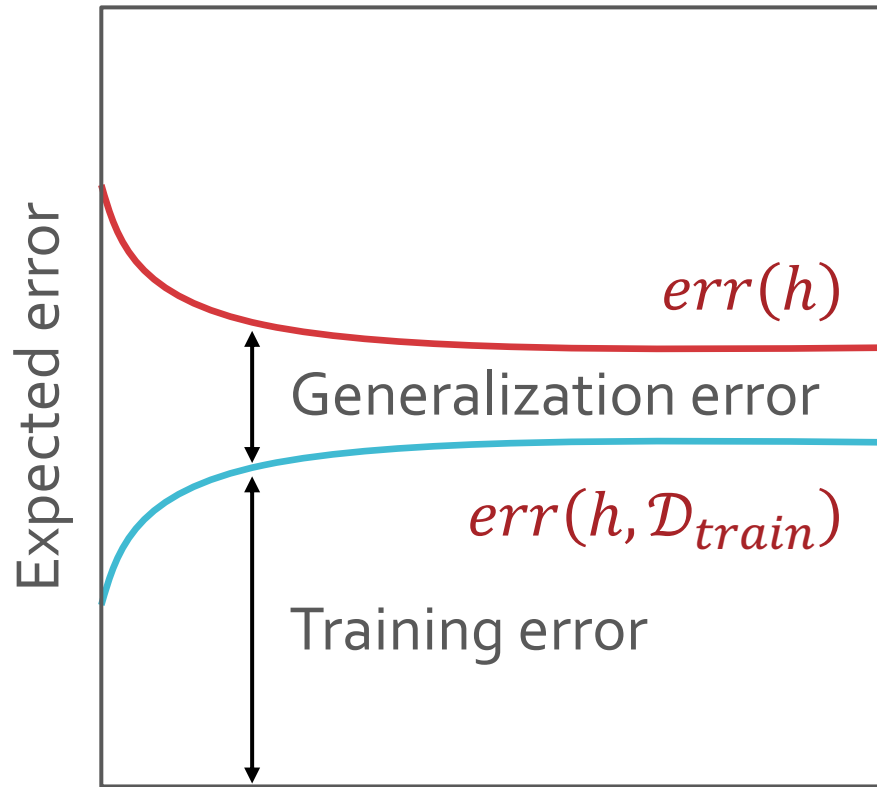


Bias of  $\bar{h}(x) \approx 0.50$   
Variance of  $h_{\mathcal{D}}(x) \approx 0.10$   
 $\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})] \approx 0.60$



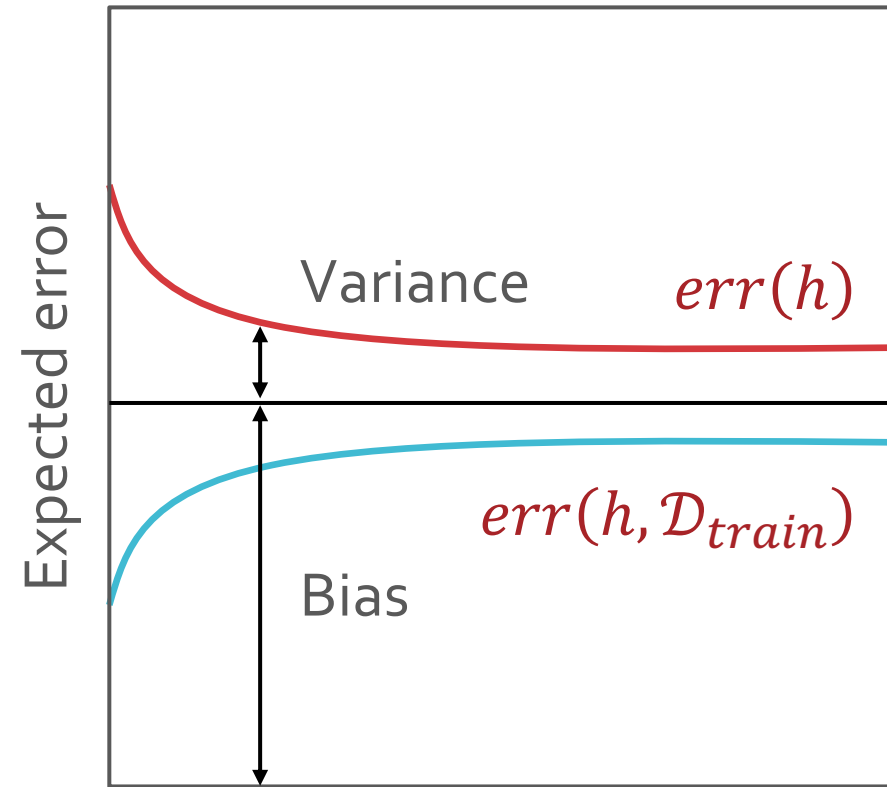
Bias of  $\bar{h}(x) \approx 0.21$   
Variance of  $h_{\mathcal{D}}(x) \approx 0.21$   
 $\mathbb{E}_{\mathcal{D}}[err(h_{\mathcal{D}})] \approx 0.42$





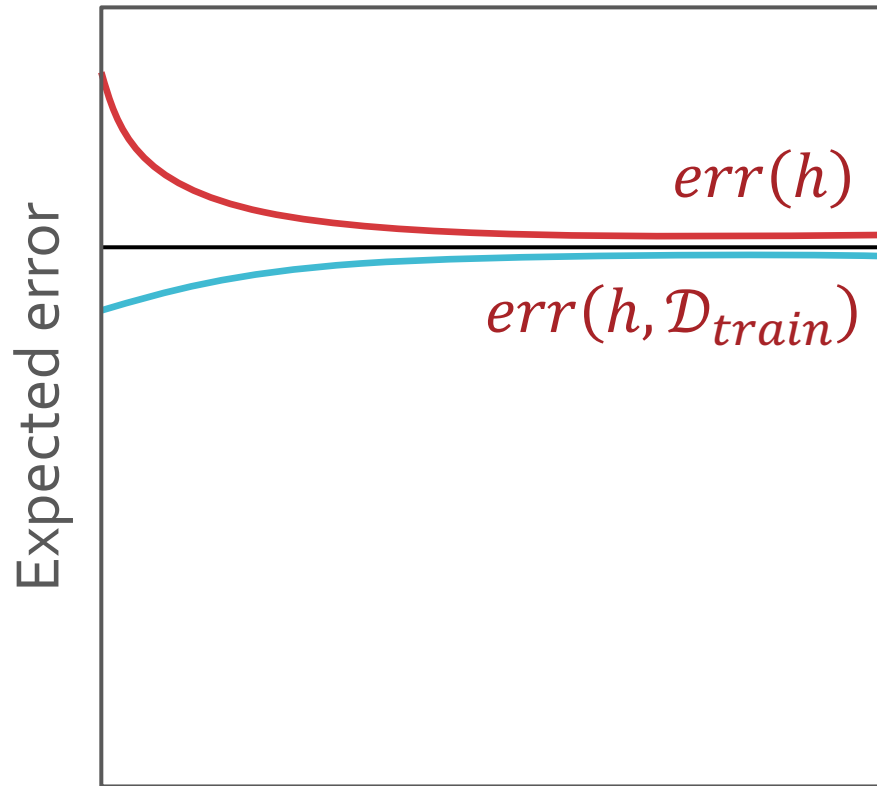
Number of training points,  $N$

Generalization



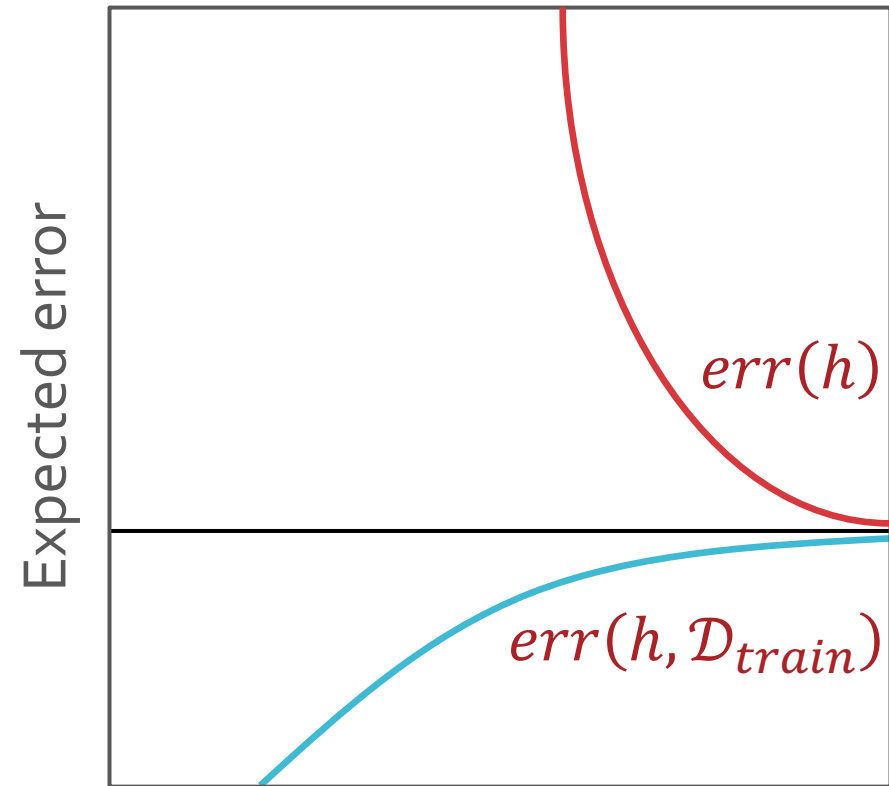
Number of training points,  $N$

Bias-Variance analysis



Number of training points,  $N$

Simple model



Number of training points,  $N$

Complex model