

# CSE 446

# Gradient Descent

---

Natasha Jaques



# How do we find optimal weights?

- This is related some questions you might have so far in this course

- Why do we use quadratic loss,  $\sum_{i=1}^n (y_i - w^T x_i)^2$ ?

- Why is Gaussian noise so popular?

$$f(z) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(z-\mu)^2}{\sigma^2}}$$

- Why was Ridge Regression with  $L_2$  regularizer,  $\|w\|_2^2$ , the first to be used?

- When we want sparsity, why do we use  $L_1$  regularizer,  $\|w\|_1$ , and not  $L_{0.5}$  regularizer,  $\|w\|_{0.5}$ ?

convexity

local minima is the global minimum



# Why gradient descent?

- Standard ML paradigm: Define loss, then optimize

$$\hat{w}_{LS} = \arg \min_w \|y - Xw\|_2^2$$

$$\rightarrow 2X^T(Xw - y) = 0$$

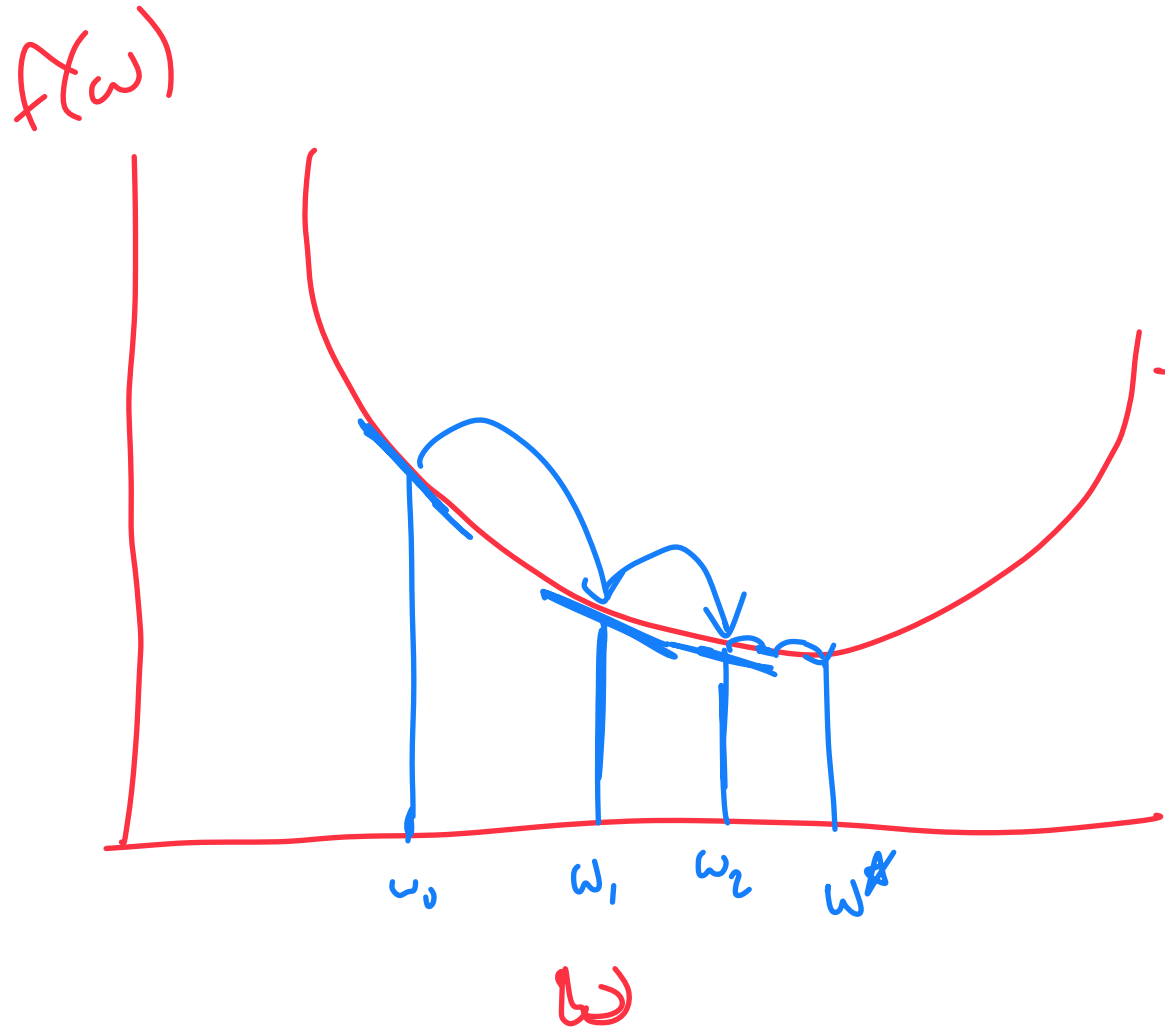
$$\hat{w} = (X^T X)^{-1} X^T y$$

$\rightarrow \lambda \|w\|_1$   
∴ no closed form solution

- But, no closed-form solutions for most losses we use in practice.
- Key idea: Iterative methods
- Used everywhere!

$\rightarrow$  start with a guess for  $w$   
 $\rightarrow$  iteratively refine to minimize loss

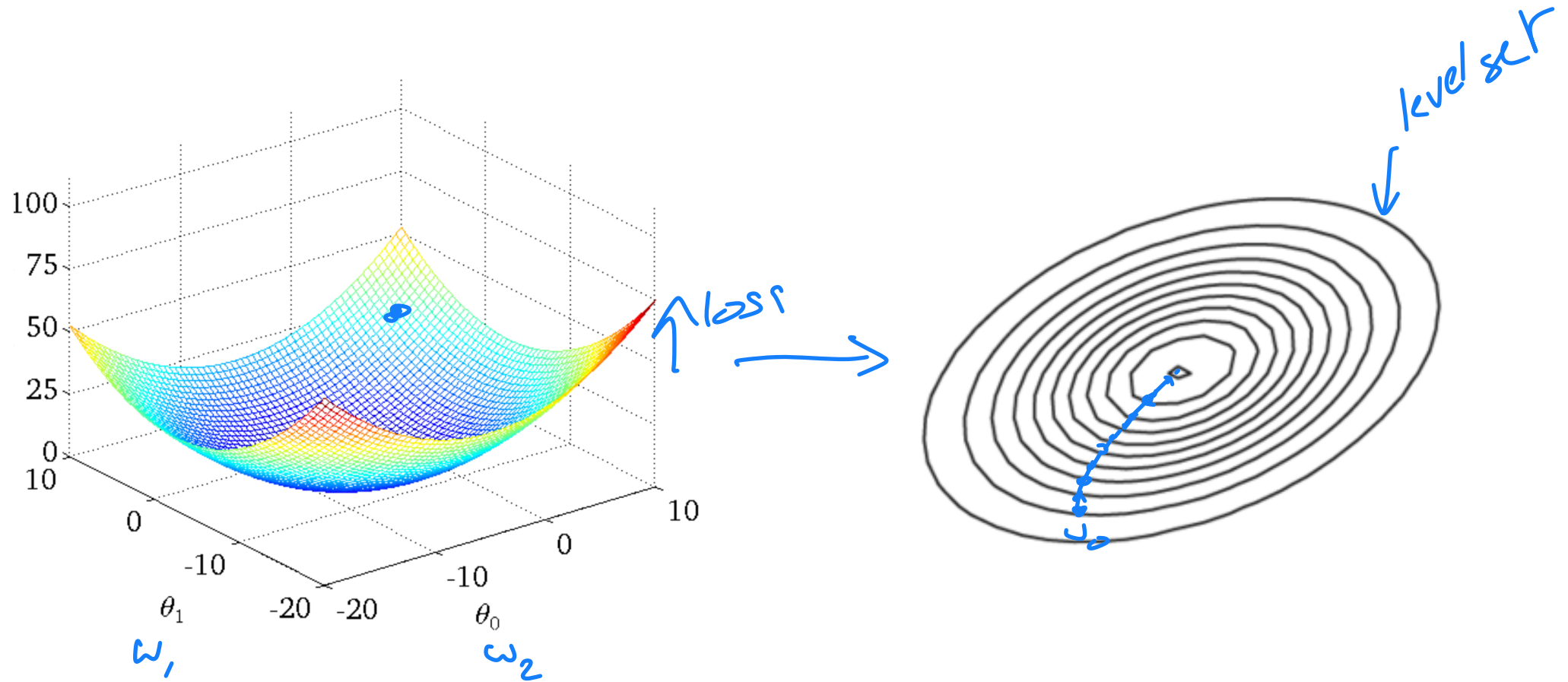
# Gradient descent in one dimension



Step direction:  $-\text{gradient}$   
Step size:  $\propto |\text{gradient}|$

$$\eta \times |\text{gradient}|$$

# Gradient descent in multiple dimensions



# Lecture plan

- Gradient descent algorithm + examples
- Theoretical analysis
  - When does it work?
  - How quickly does it converge?
  - How do we choose a step size?
  - Key idea: Convexity
- Not tested on proof details, but concepts are important & practical

# Algorithm: Gradient descent

**Algorithm**

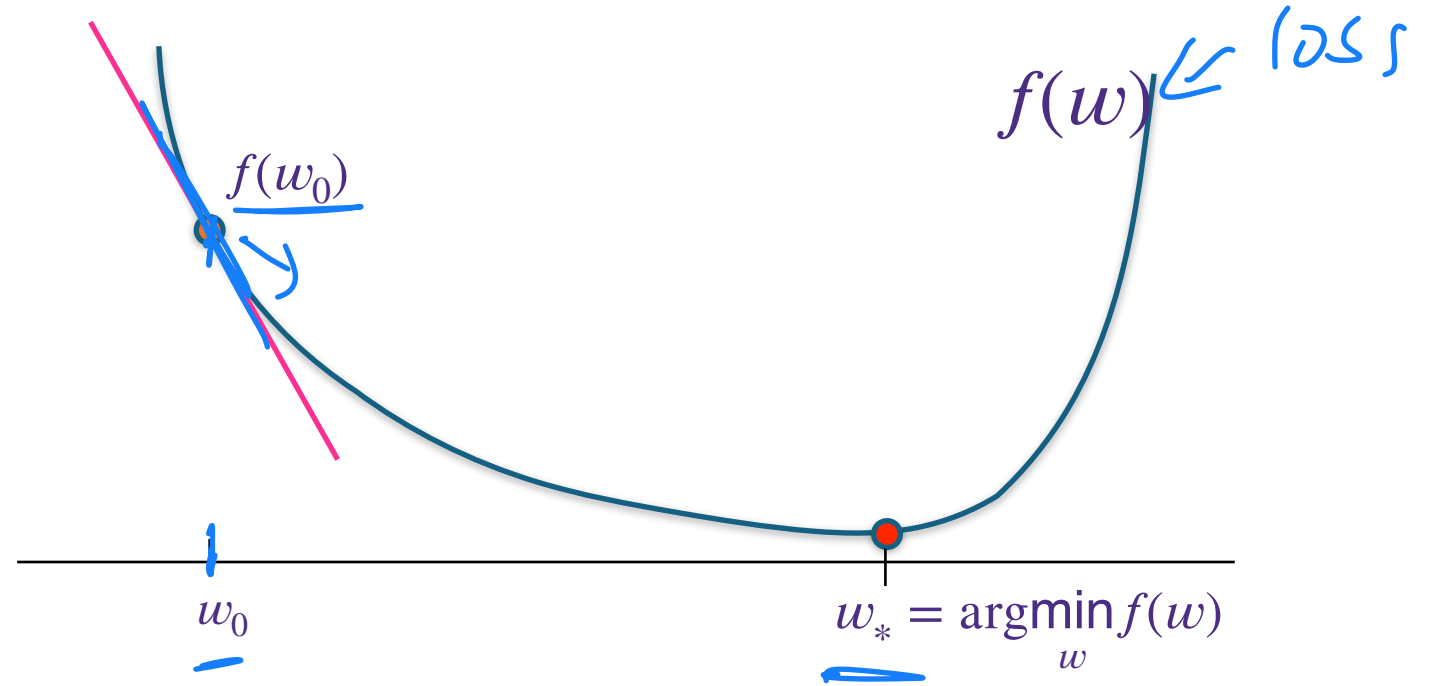
For  $t=0,1,2,3, \dots$

$$w_{t+1} = w_t - \eta \left. \frac{df(w)}{dw} \right|_{w=w_t}$$

Hyperparameters:

- Initial point  $w_0$
- Step size  $\eta$

*Handwritten annotations:* "step" with an arrow pointing to the iteration loop; "gradient" with an arrow pointing to the derivative term; "step size" with an arrow pointing to  $\eta$ .



# Algorithm: Gradient descent

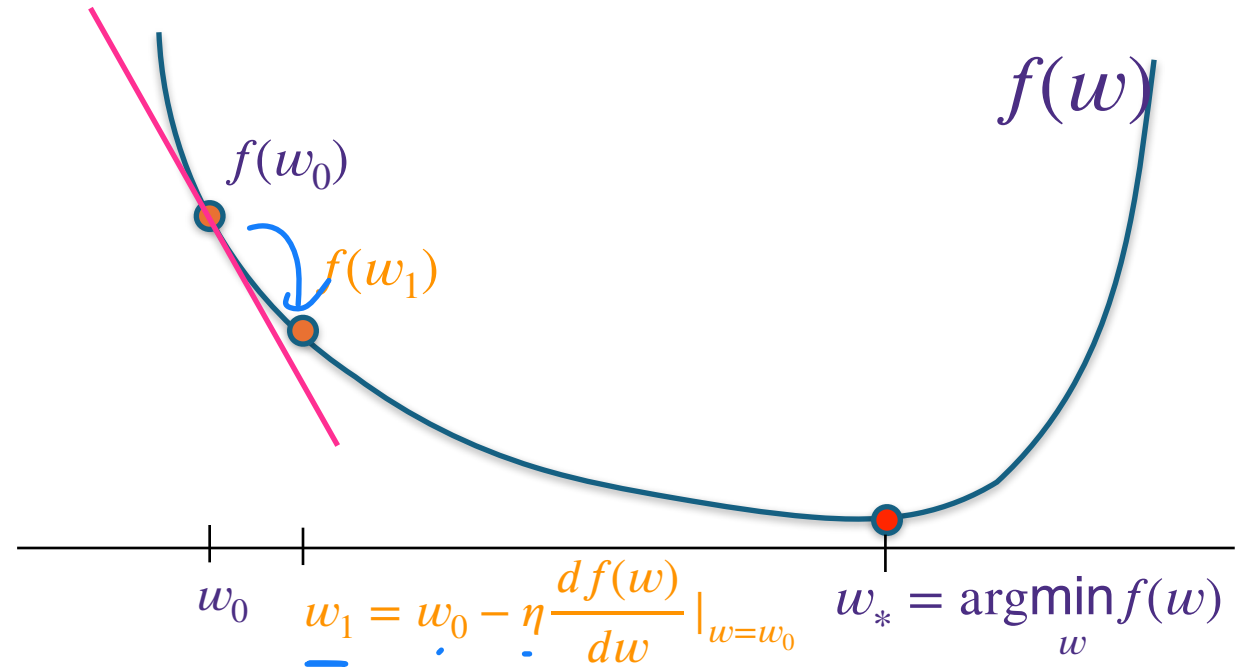
## Algorithm

For  $t=0,1,2,3, \dots$

$$w_{t+1} = w_t - \eta \frac{df(w)}{dw} \Big|_{w=w_t}$$

Hyperparameters:

- Initial point  $w_0$
- Step size  $\eta$



# Algorithm: Gradient descent

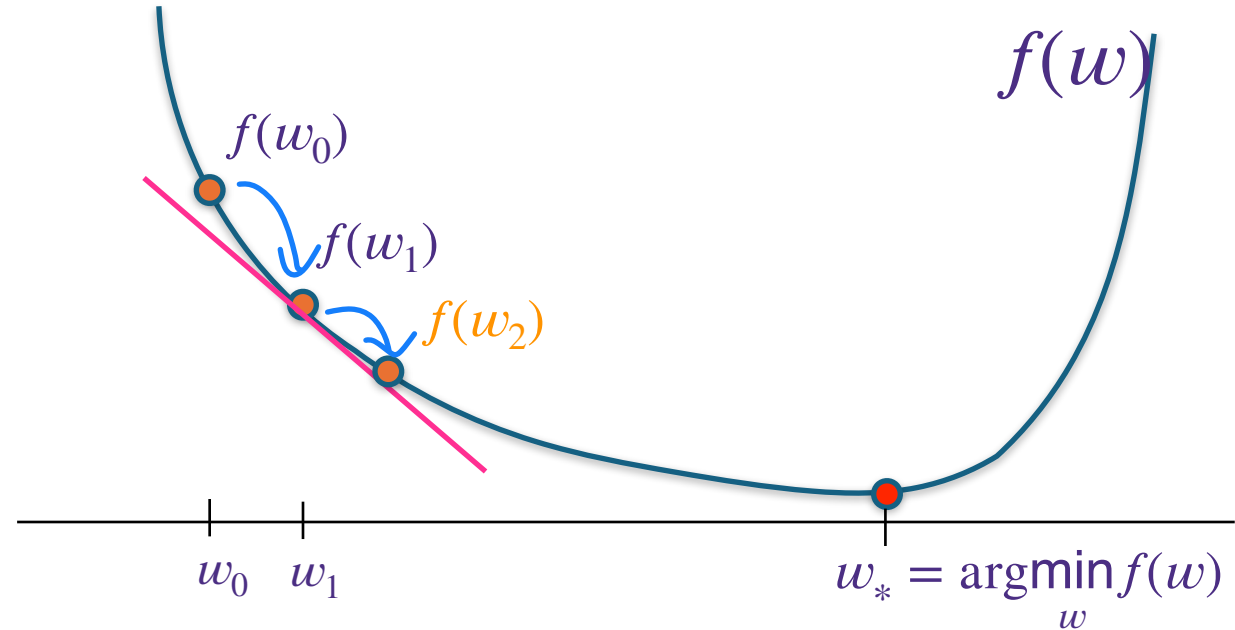
## Algorithm

For  $t=0,1,2,3, \dots$

$$w_{t+1} = w_t - \eta \frac{df(w)}{dw} \Big|_{w=w_t}$$

Hyperparameters:

- Initial point  $w_0$
- Step size  $\eta$



# Algorithm: Gradient descent

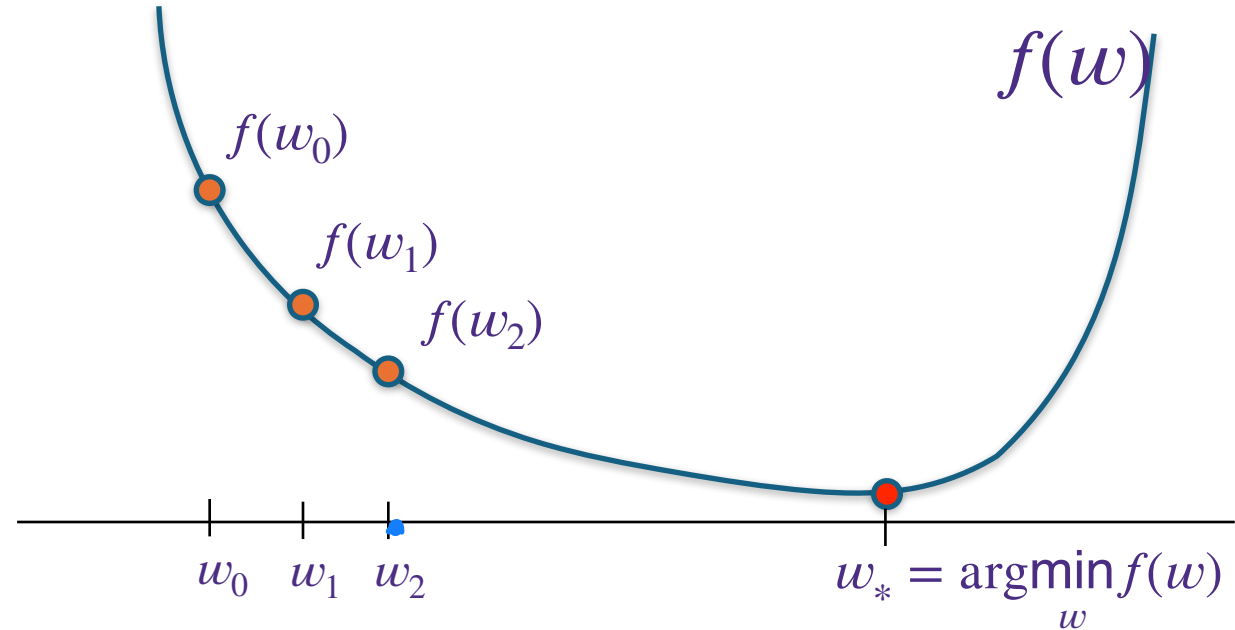
## Algorithm

For  $t=0,1,2,3, \dots$

$$w_{t+1} = w_t - \eta \left. \frac{df(w)}{dw} \right|_{w=w_t}$$

Hyperparameters:

- Initial point  $w_0$
- Step size  $\eta$



# Algorithm: Gradient descent

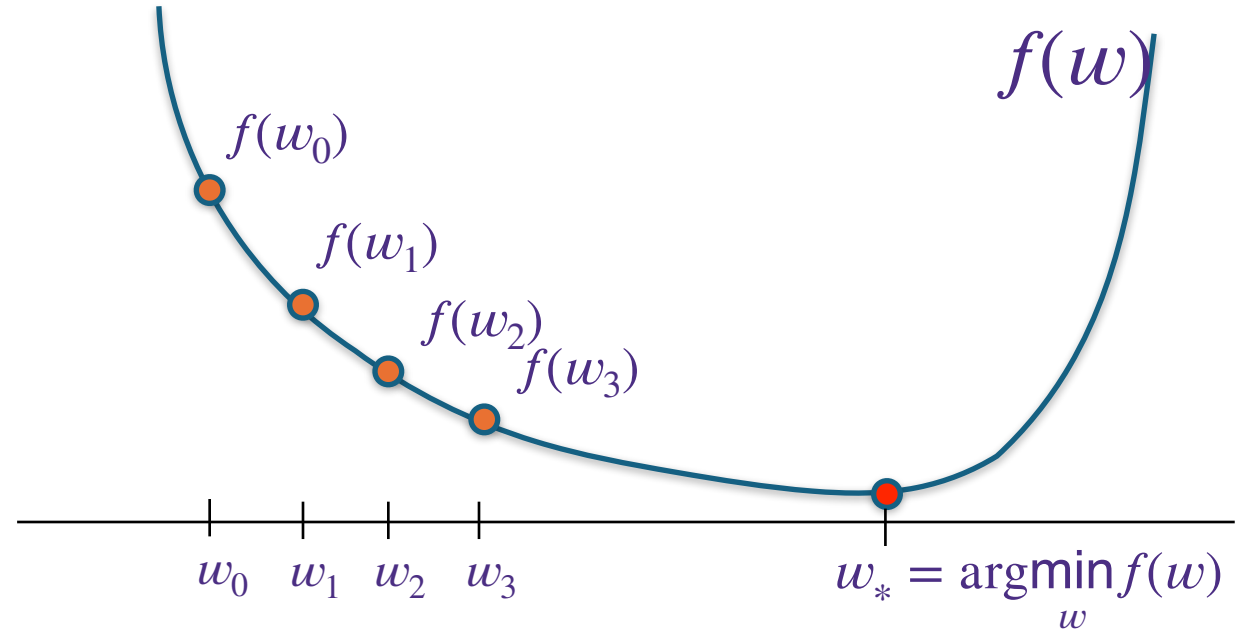
## Algorithm

For  $t=0,1,2,3, \dots$

$$w_{t+1} = w_t - \eta \frac{df(w)}{dw} \Big|_{w=w_t}$$

Hyperparameters:

- Initial point  $w_0$
- Step size  $\eta$



Left off

# Algorithm: Gradient descent

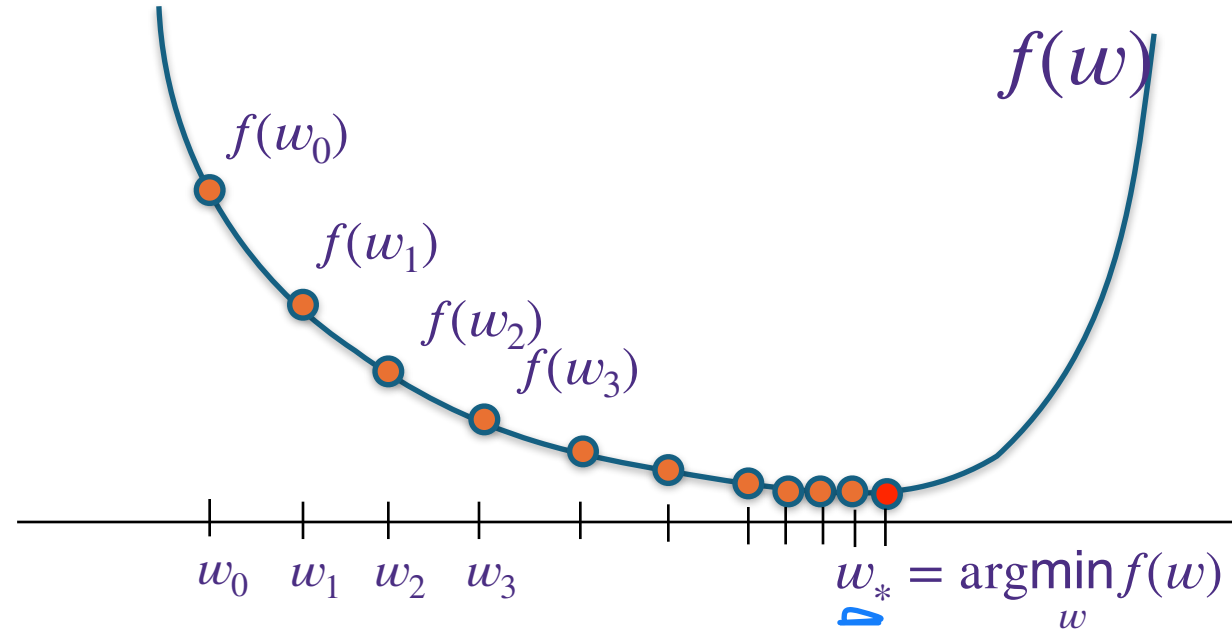
## Algorithm

For  $t=0,1,2,3, \dots$

$$w_{t+1} = w_t - \eta \frac{df(w)}{dw} \Big|_{w=w_t}$$

Hyperparameters:

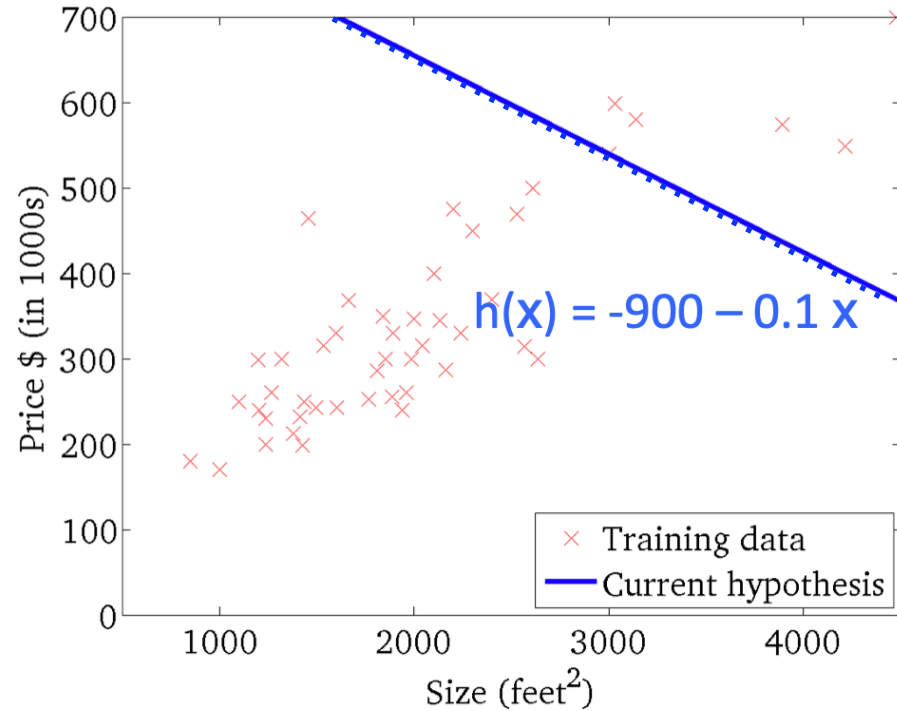
- Initial point  $w_0$
- Step size  $\eta$



Note that as  $t \rightarrow \infty$  we have  $\frac{df(w)}{dw} \Big|_{w=w_t} \rightarrow 0$

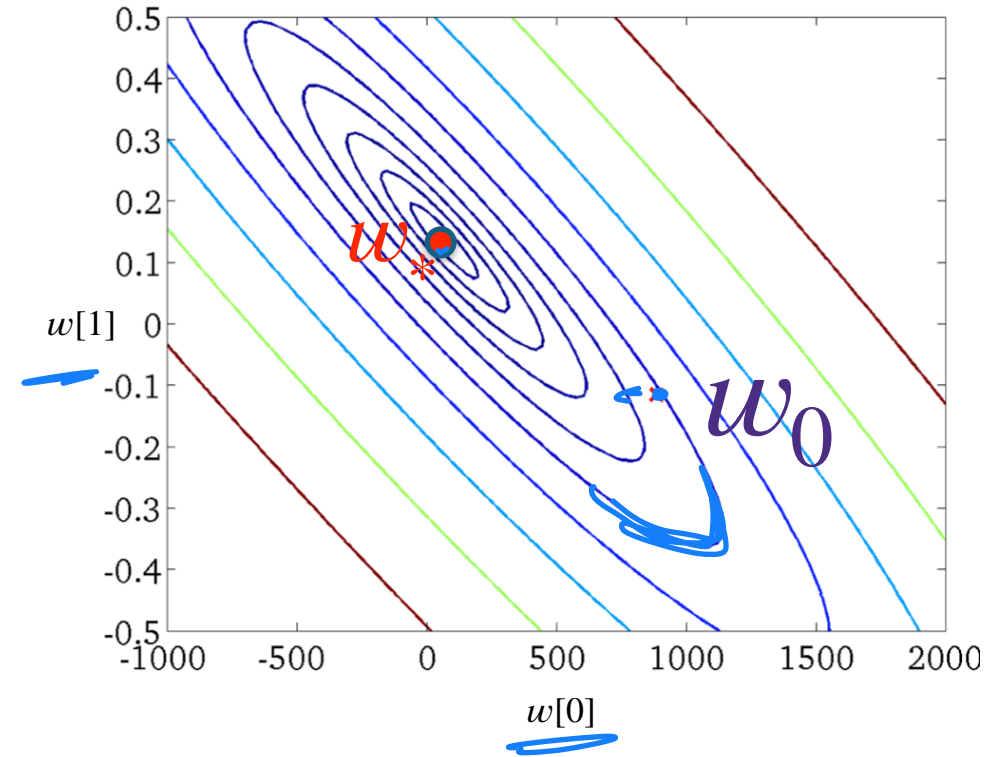
# 1-dimensional linear regression with 2 parameters

$$\{(x_i, y_i)\}_{i=1}^n$$



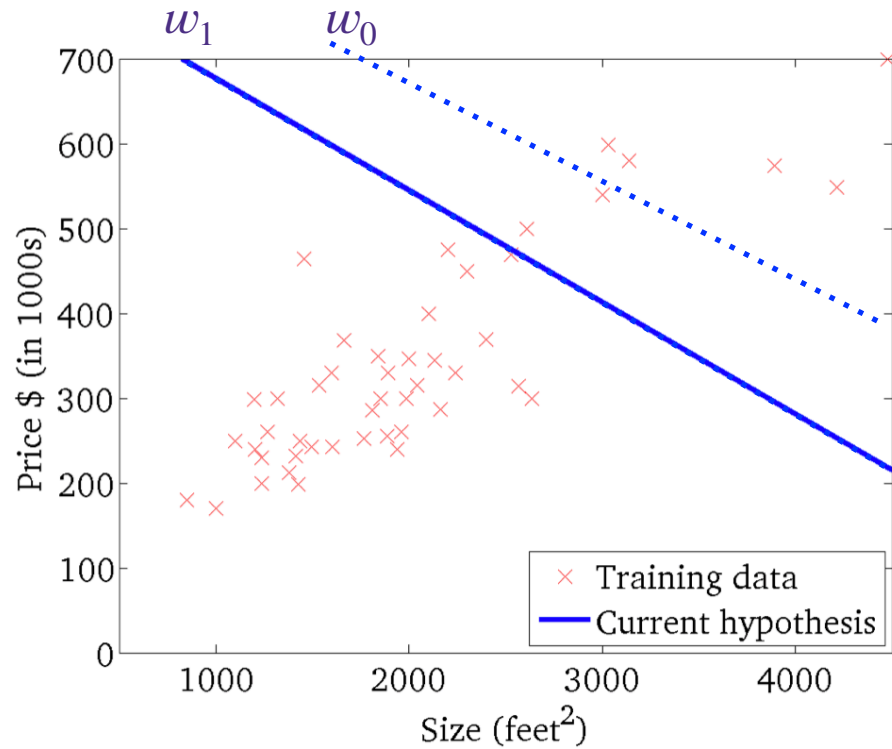
Evolution of the predictor  $y = w[0] + w[1]x$

$$w_{t+1} \leftarrow w_t - \eta \cdot \nabla_w f(w_t)$$

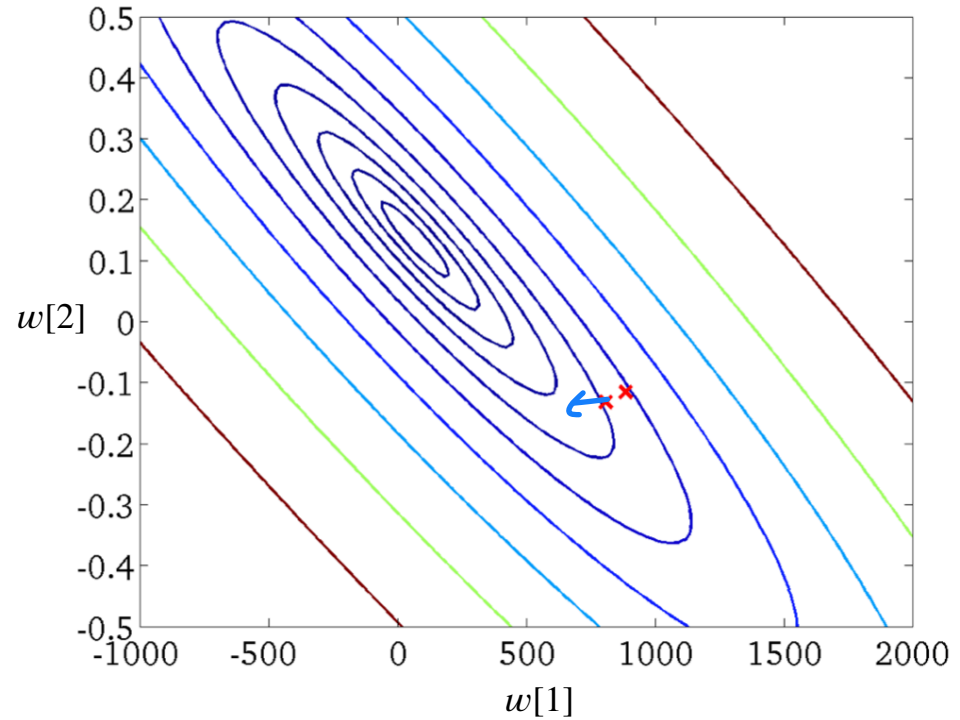


Gradient descent dynamics in the parameter space  $w$   
Ovals show the level set of the objective function

# 1-dimensional linear regression with 2 parameters

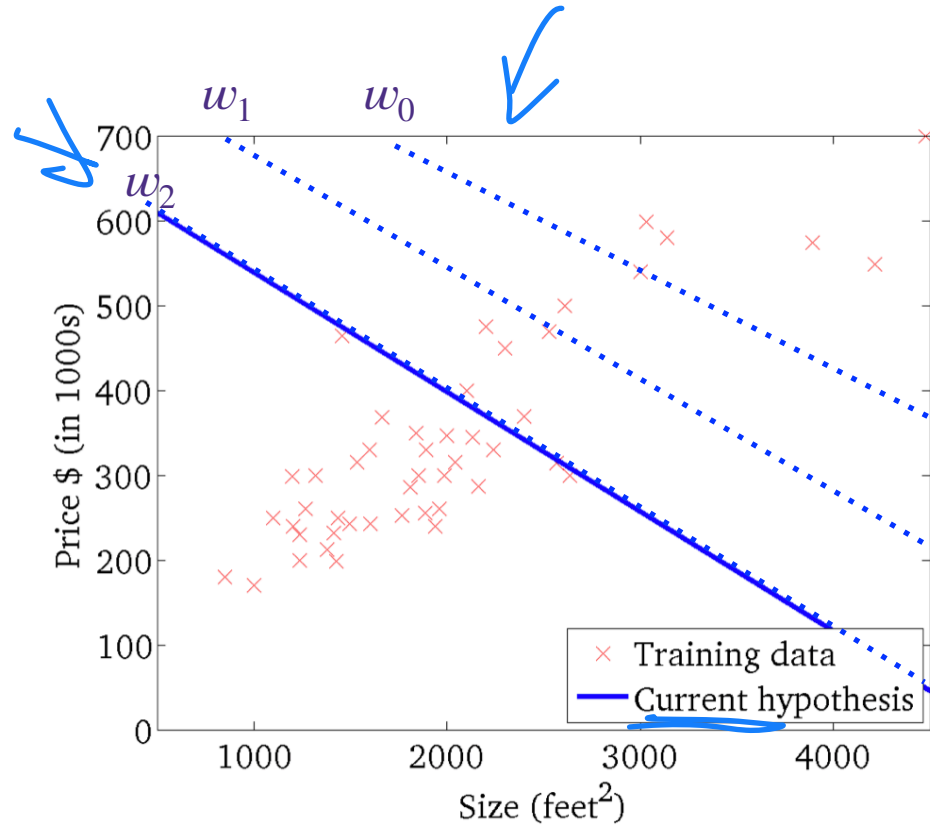


Evolution of the predictor  $y = w[0] + w[1]x$

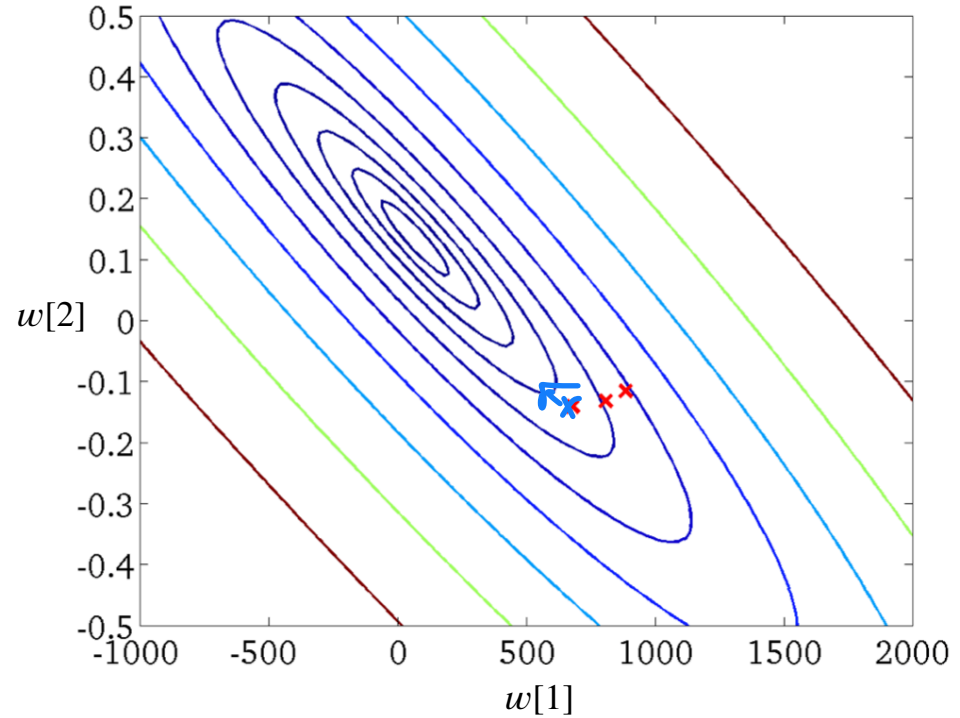


Gradient descent dynamics in the parameter space  $w$   
Ovals show the **level set** of the objective function

# 1-dimensional linear regression with 2 parameters

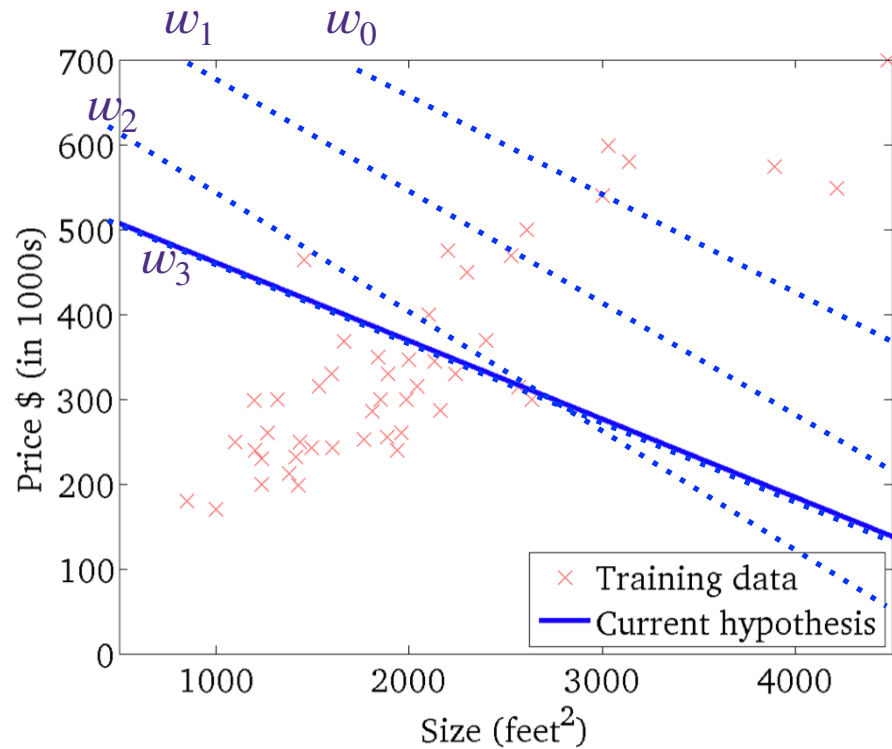


Evolution of the predictor  $y = w[0] + w[1]x$

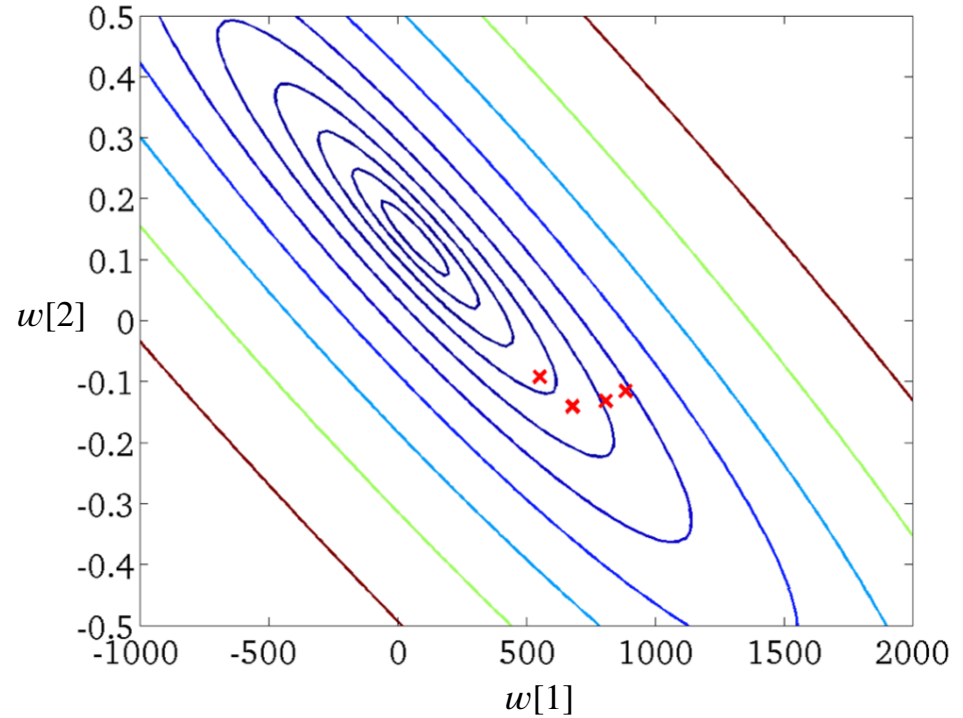


Gradient descent dynamics in the parameter space  $w$   
Ovals show the **level set** of the objective function

# 1-dimensional linear regression with 2 parameters

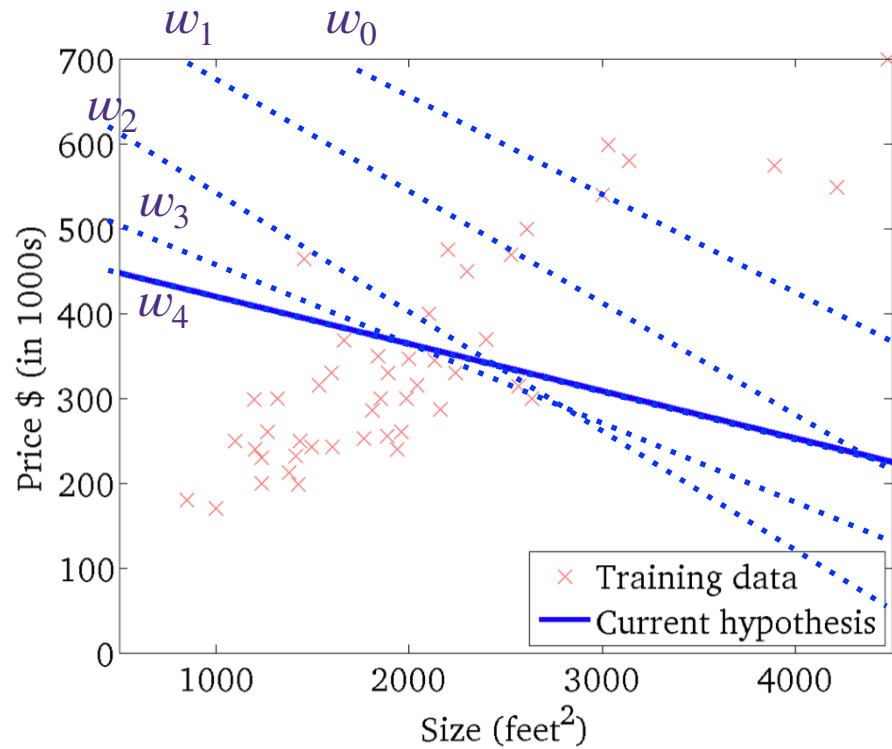


Evolution of the predictor  $y = w[0] + w[1]x$

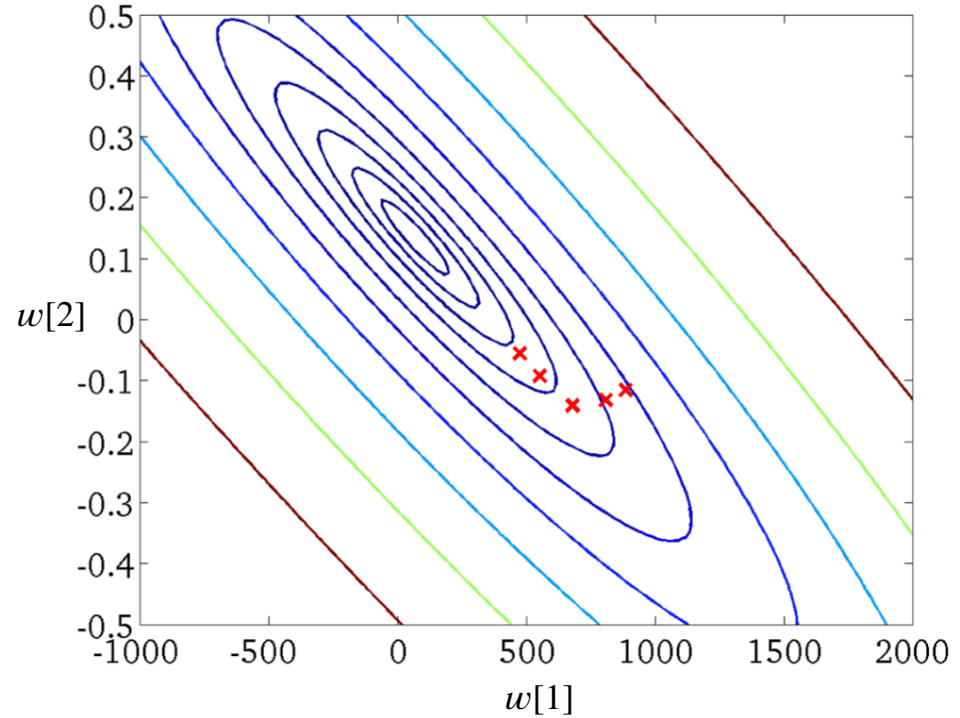


Gradient descent dynamics in the parameter space  $w$   
Ovals show the **level set** of the objective function

# 1-dimensional linear regression with 2 parameters

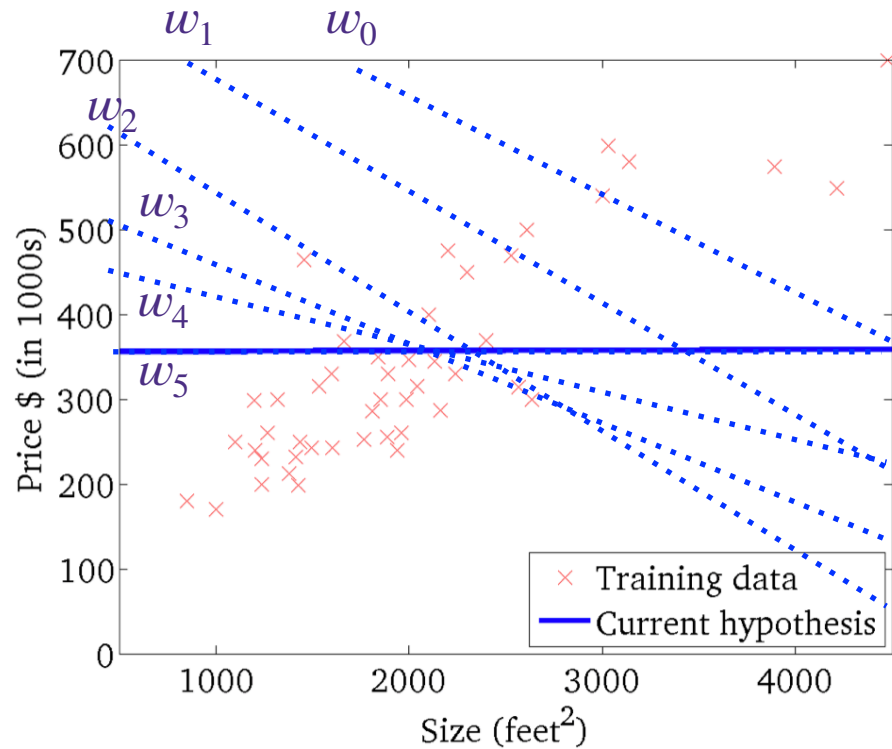


Evolution of the predictor  $y = w[0] + w[1]x$

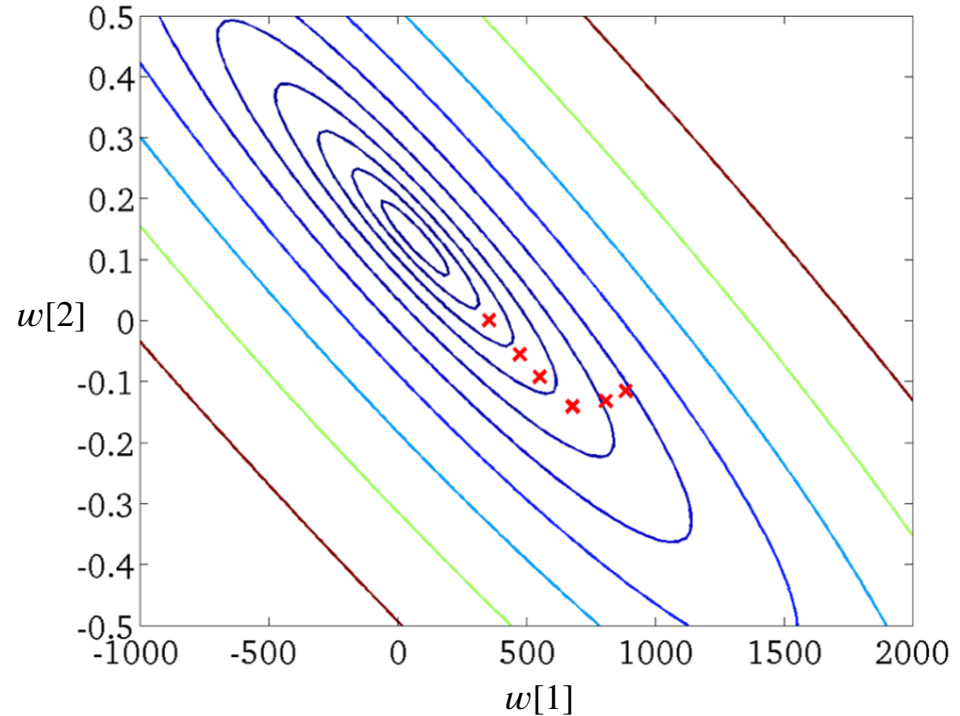


Gradient descent dynamics in the parameter space  $w$   
Ovals show the **level set** of the objective function

# 1-dimensional linear regression with 2 parameters

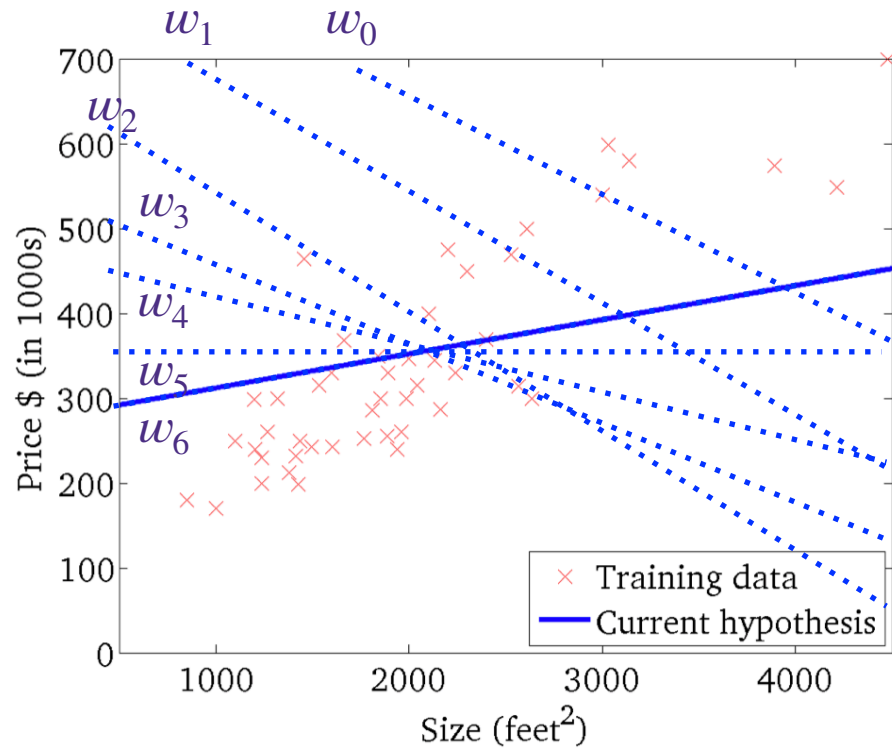


Evolution of the predictor  $y = w[0] + w[1]x$

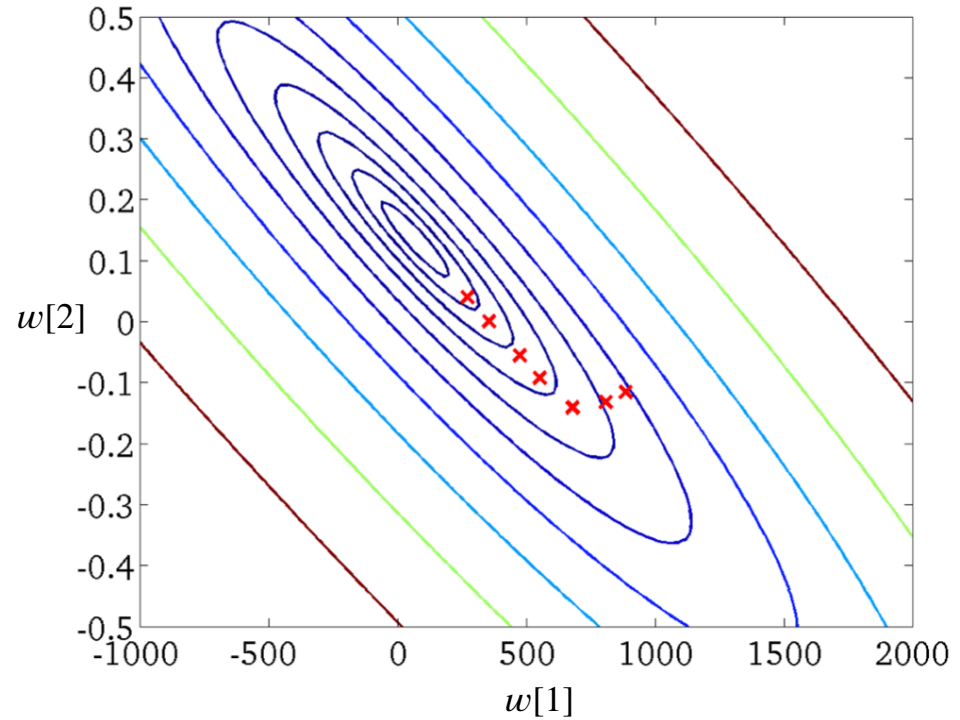


Gradient descent dynamics in the parameter space  $w$   
Ovals show the **level set** of the objective function

# 1-dimensional linear regression with 2 parameters

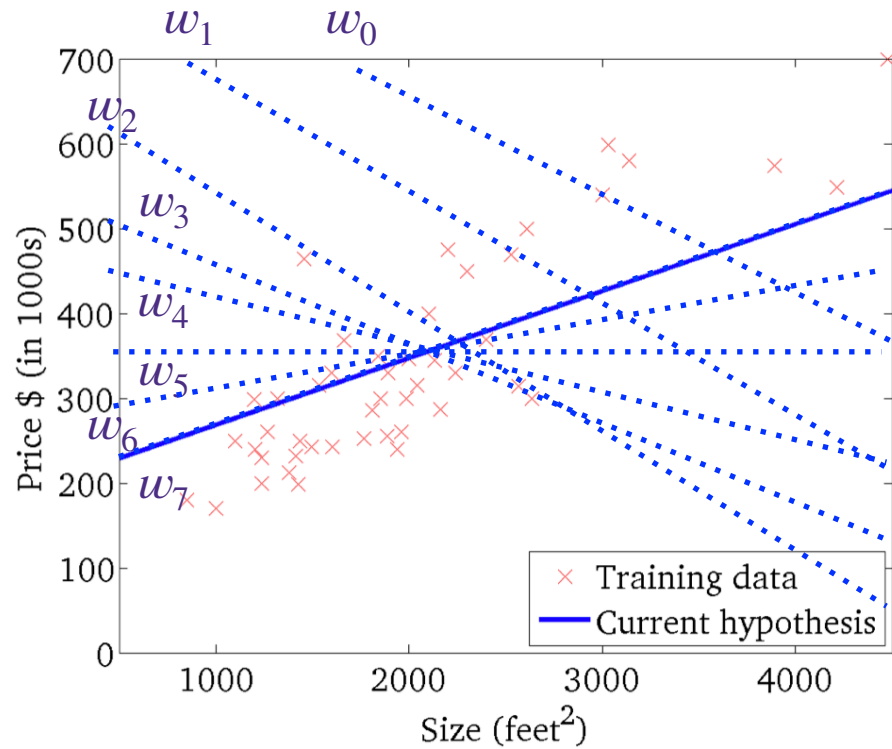


Evolution of the predictor  $y = w[0] + w[1]x$

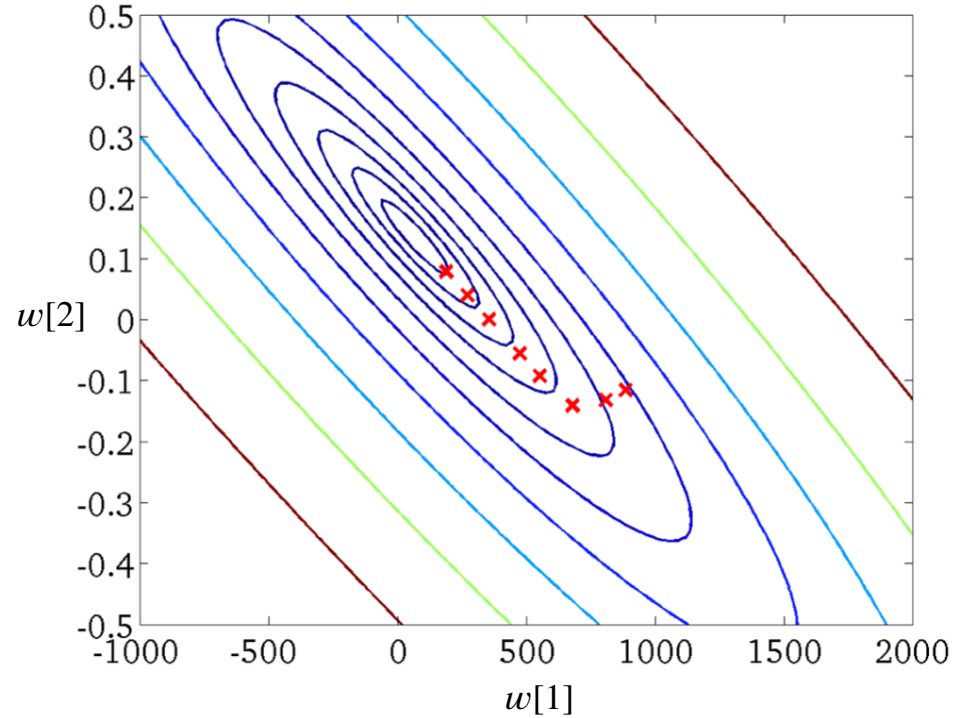


Gradient descent dynamics in the parameter space  $w$   
Ovals show the **level set** of the objective function

# 1-dimensional linear regression with 2 parameters



Evolution of the predictor  $y = w[0] + w[1]x$

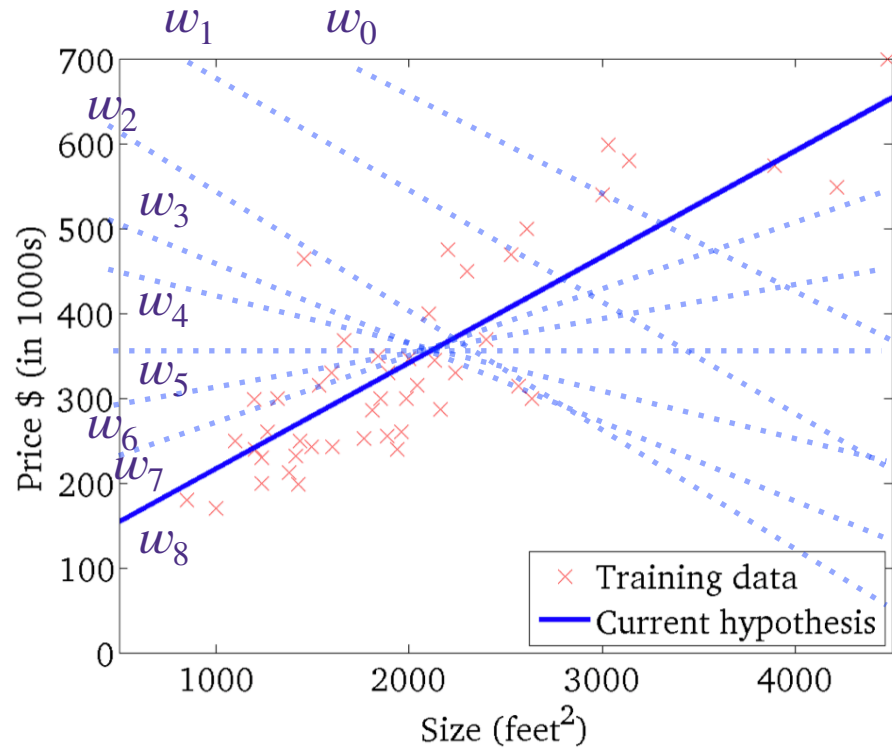


Gradient descent dynamics in the parameter space  $w$   
Ovals show the **level set** of the objective function

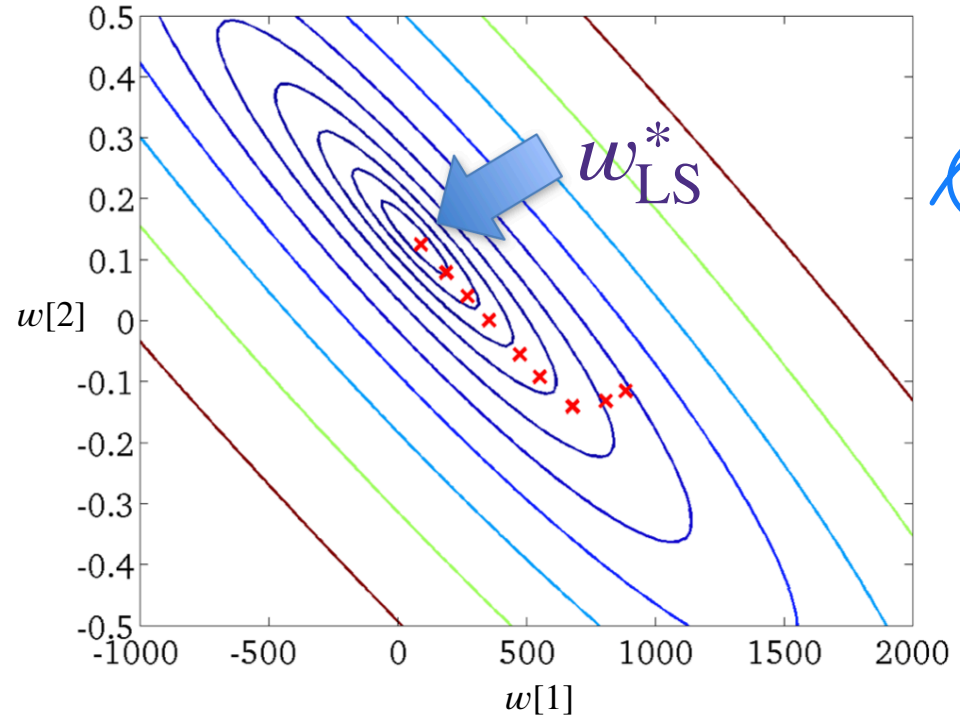
# 1-dimensional linear regression with 2 parameters

$$w_0 \sim N(0, I_{d \times d} \sigma^2)$$

$\downarrow$   
 $\downarrow$   
 Hypoparams  
 $\mathcal{N}$   
 Gaussian



Evolution of the predictor  $y = w[0] + w[1]x$



Gradient descent dynamics in the parameter space  $w$   
 Ovals show the **level set** of the objective function

# Warmup: Quadratic functions

$\rightarrow f(w)$  (loss)

$$\hat{w} = \operatorname{argmin}_w \underline{aw^2 + bw + c}$$

$$w_0 \sim \mathcal{N}(0, I\sigma^2)$$

$$\left. \frac{df(w)}{dw} \right|_{w=w_0} = 2aw_0 + b$$

$$w_1 = w_0 - \eta(2aw_0 + b)$$

recall!

$$\boxed{w_{t+1} = w_t - \eta \left. \frac{df(w)}{dw} \right|_{w=w_t}}$$

# Example: Linear regression

$f(w)$

$$\hat{w} = \underset{w}{\operatorname{argmin}} \frac{1}{2} \sum_i (y_i - x_i^\top w)^2$$

$$\nabla_w f(w_0) = X^\top (Xw - Y)$$

Diagram illustrating the dimensions of the matrices in the gradient formula:

- $X$ :  $d \times n$  matrix
- $w$ :  $n \times 1$  vector
- $Y$ :  $d \times 1$  vector

$$w_{t+1} \leftarrow w_t - \eta X^\top (Xw_0 - Y)$$

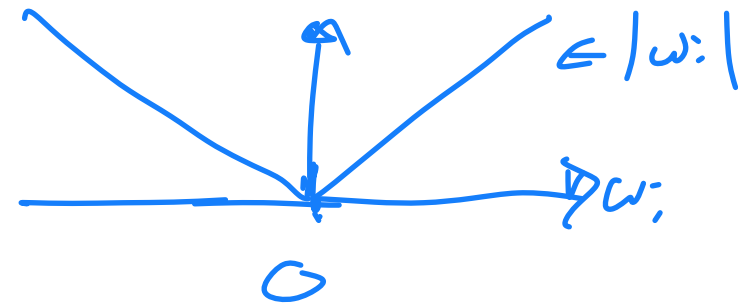
# Example: Lasso

$$\hat{w} = \underset{w}{\operatorname{argmin}} \frac{1}{2} \|y - Xw\|_2^2 + \lambda \|w\|_1$$

? Double check

$$\nabla_w f = X^T(Xw - y) + \lambda \left( \sum_{i=1}^n \right) \operatorname{sign}(w_i)$$

$$w_{t+1} \leftarrow w_t - \eta \cdot \nabla_w f|_{w=w_t}$$



convex  $\rightarrow$  local minima = global minimum

$$\frac{d|w_i|}{dw_i} = \begin{cases} +1 & w_i > 0 \\ [1, 1) & w_i = 0 \\ -1 & w_i < 0 \end{cases}$$

undefined  
 $\hookrightarrow$  so define subgradient

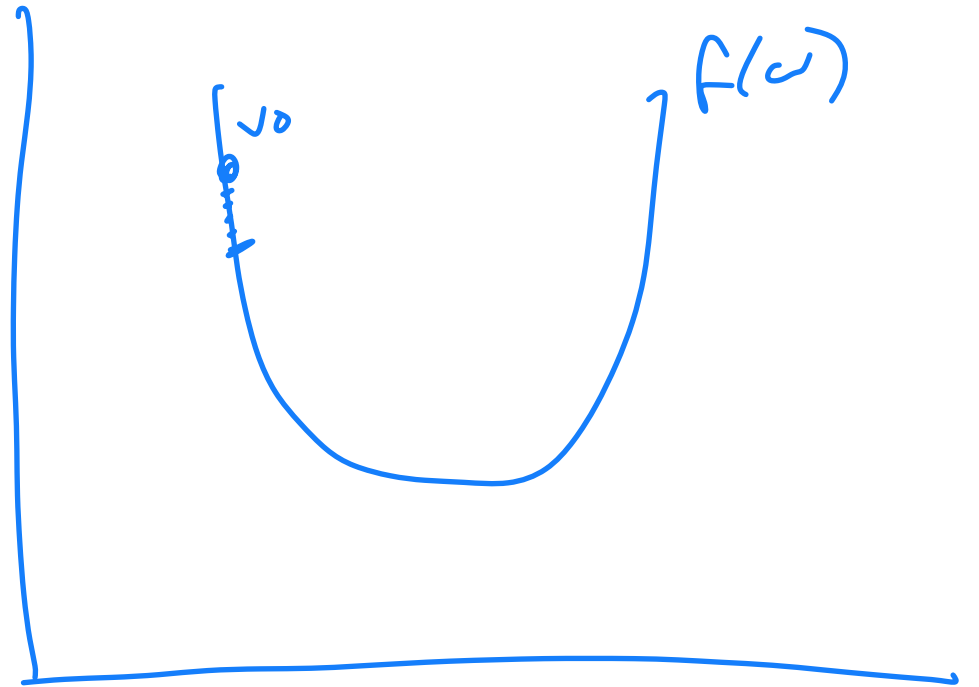
$$f(w) = aw^2 + bw + c + |w|$$

$$\frac{df}{dw} = 2aw + b + \text{sign}(w) = 0 \quad \text{// can find if } w \text{ is 1-d}$$

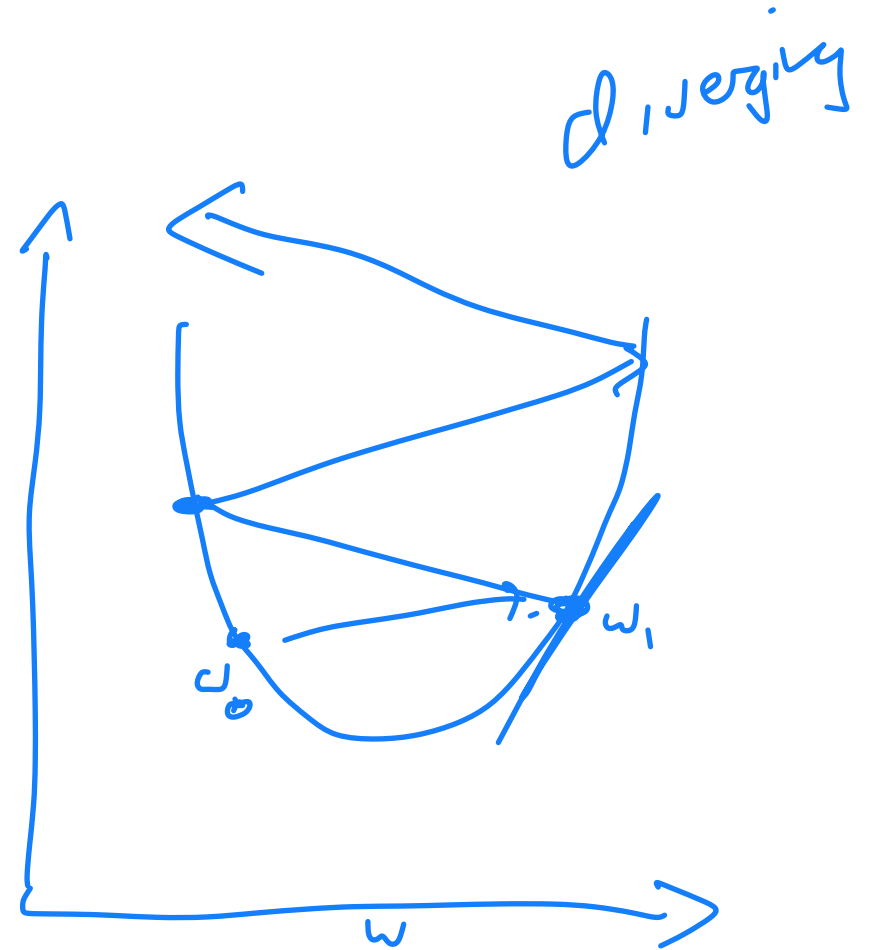
$$w \in \mathbb{R}^d$$

$\text{sign}(w) \rightarrow 2^d$  possibilities

# How do you choose a step size?



step size  $\eta$  is too small  
slow convergence

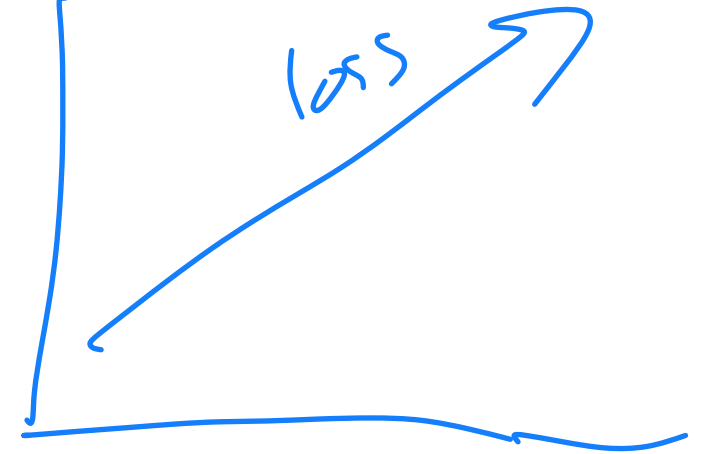
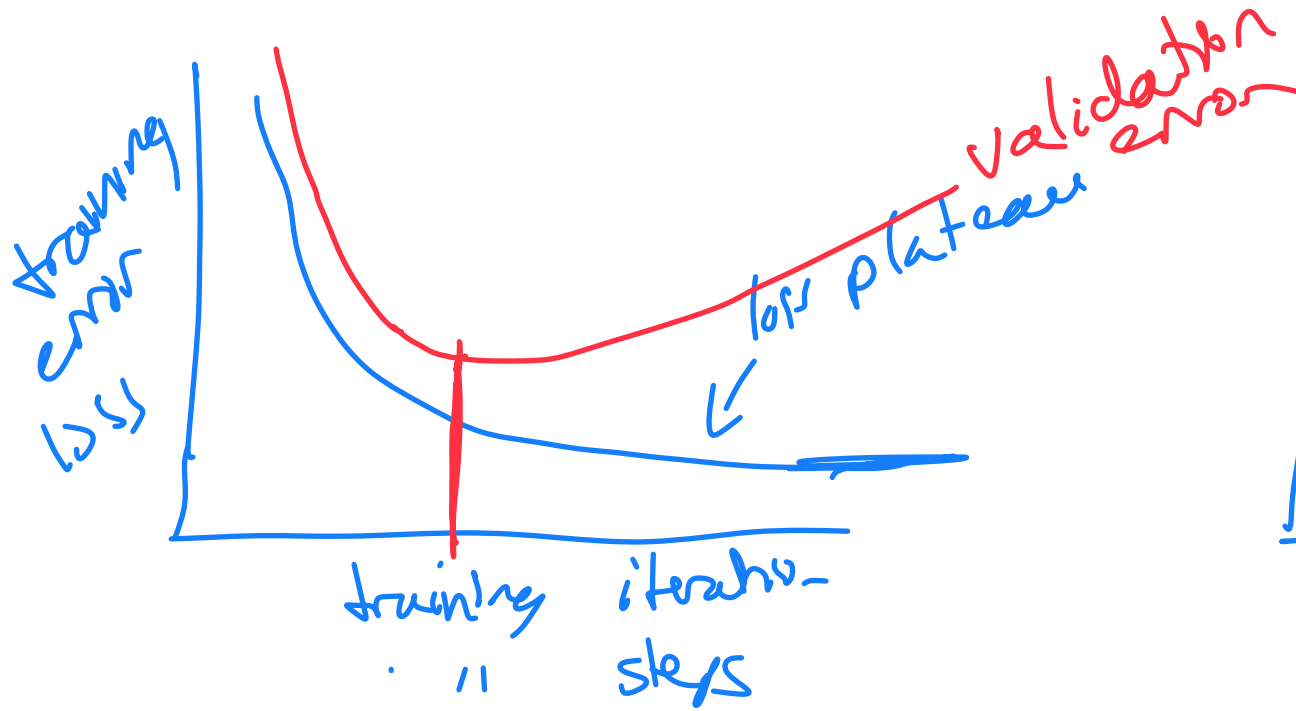


$\eta$  is too big

- If  $\eta$  too small, converges very, very slowly.
- If  $\eta$  too big, does not converge!
- In practice: guess and check

# How do you choose a step size?

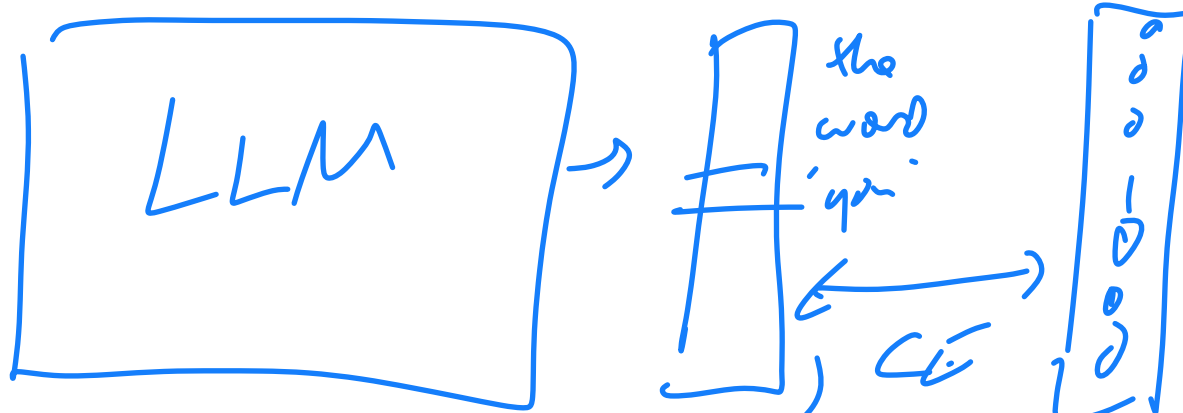
→ visualize loss over training!



train loss effectively minimized,  
but use early stopping to  
prevent over-fitting

loss increasing  
is a sign step  
size might be too  
big.

X  
"hi: how are"

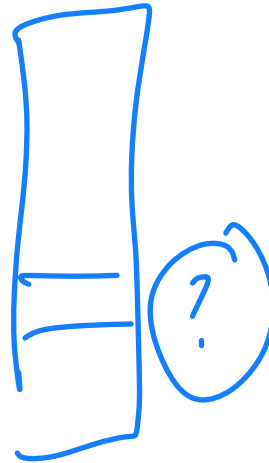
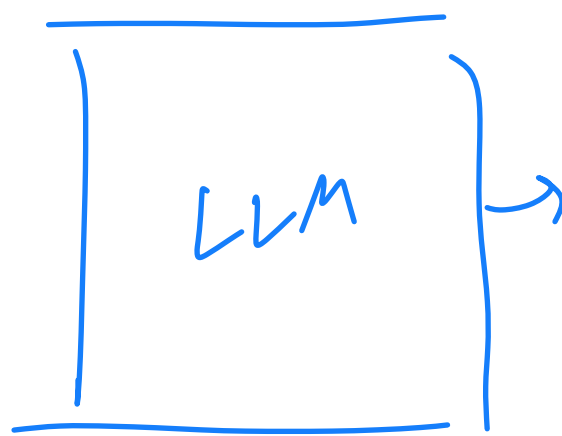


NOT ON EXAM



X  
"hi: how are you"

X



y.  
"hi: how are you?"  
"hey I'm good"

