

# Recap: Linear Regression; Closed Form Solution 2

$$W^* = \underset{W}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - W^T x_i)^2 = \underset{W}{\operatorname{argmin}} \frac{1}{n} \|Y - XW\|_2^2$$

Where  $X = \begin{bmatrix} \leftarrow x_1^T \rightarrow \\ \vdots \\ \leftarrow x_n^T \rightarrow \end{bmatrix}_{n \times (d+1)}$

$$x_i = \begin{bmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{id} \end{bmatrix}_{(d+1) \times 1}$$

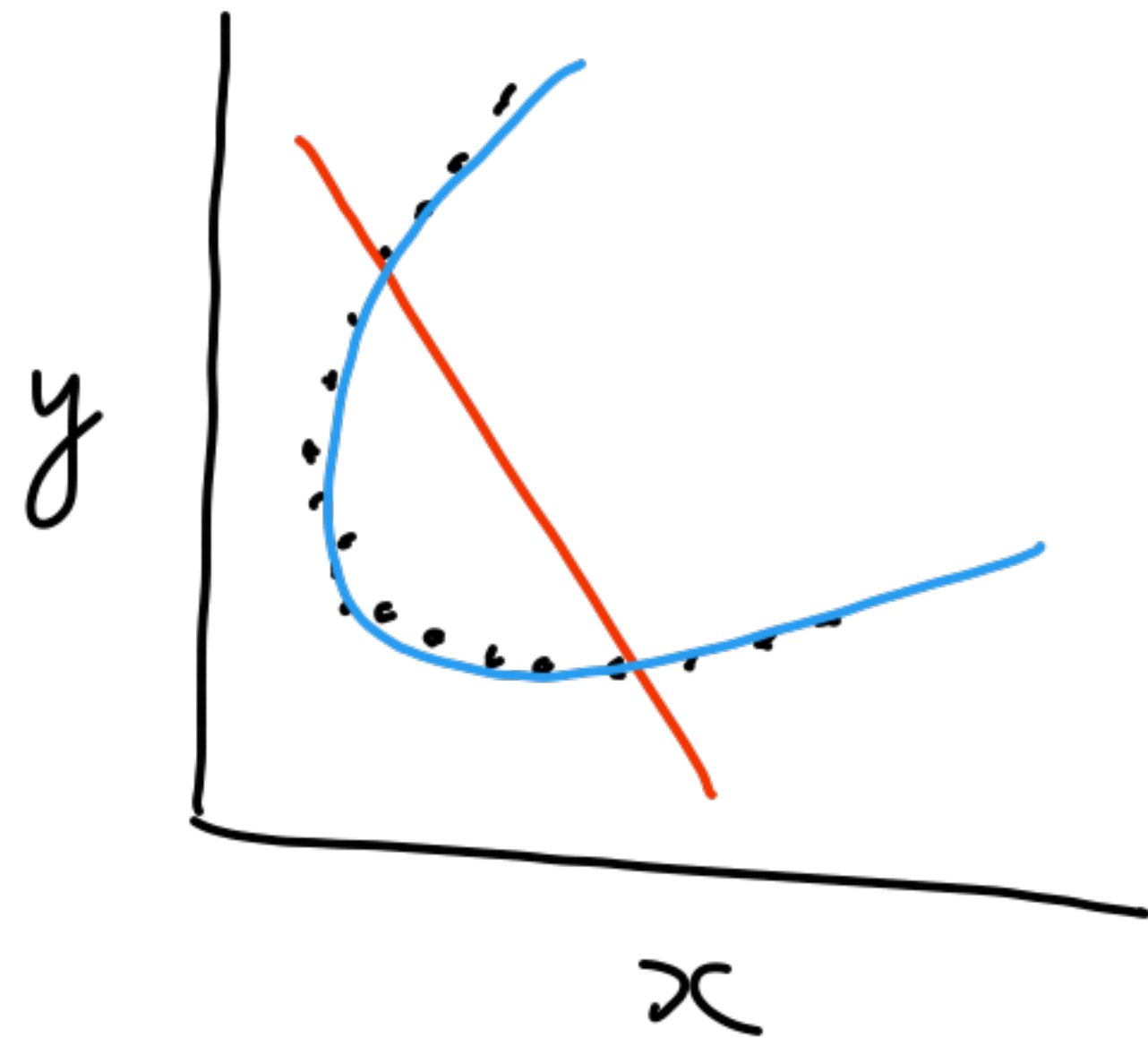
$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}_{(d+1) \times 1}$$

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$\left[ \begin{array}{c} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_d} \end{array} \right]$$

$$\begin{aligned} \nabla_W L = 0 &\Rightarrow -\frac{2}{n} \sum_i (y_i - W^T x_i) x_i = 0 \Rightarrow -2 \sum_i y_i x_i + 2 \sum_i (W^T x_i) x_i = 0 \\ &\Rightarrow 2X^T Y = 2X^T X W \Rightarrow W = (X^T X)^{-1} X^T Y \end{aligned}$$

Consider the following dataset



$$h_w(x) = w_0 + w_1 x$$

$$h_w(x) = w_0 + w_1 x + w_2 x^2$$

More generally, can we transform the datapoints to a different dimensionality so as to enable a better regression fit

Do this with "BASIS FUNCTIONS"

# BASIS FUNCTION

Define  $\phi: \mathcal{X} \rightarrow \mathbb{R}^m$ ,  $\phi(x) = \begin{bmatrix} 1 \\ \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_m(x) \end{bmatrix}$   $\xrightarrow{\text{Scalars}} (m+1) \times 1$

$x \in \mathbb{R}^d$   
 $= \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$

Typically  $m \gg d$

but  $d > m$  as well if there are correlated/spurious features  
that we want to eliminate in  $\mathbb{R}^d$

## Common types of basis functions

- ① Polynomial basis; for 1-D points,  $\phi_j(x) = x^j$
- ② Radial basis function (RBF); for 1D points,  
(Gaussian basis)  
$$\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{\sigma_j^2}\right\}$$

Where  $\mu_j$  and  $\sigma_j$  are predefined



③ Piecewise linear basis

④ Periodic basis ( $\sin x, \cos x, \text{etc.}$ )

⑤ Fourier basis

⋮



What is the loss function with basis functions?  
& the corresponding optimization problem

$$w^* = \underset{w}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \left( y_i - w^T \phi(x_i) \right)^2 = \underset{w}{\operatorname{argmin}} \left\| y - \phi w \right\|_2^2$$

Where  $\phi = \begin{bmatrix} \leftarrow \phi(x_1)^T \rightarrow \\ \vdots \\ \leftarrow \phi(x_n)^T \rightarrow \end{bmatrix}_{n \times (m+1)}$

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

# GRADIENT DESCENT<sup>(GD)</sup> FOR LINEAR REGRESSION

---

GD is a first-order iterative algorithm used to find local optima of a differentiable function

GD is derived from gradient descent

Gradient:  $\nabla_w L = \begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix}$ , gradient is the direction of fastest increase in  $L$



descent : We are interested in minimizing a loss fn

So, update  $W$  in the reverse direction of the gradient

Apply GD in order to estimate  $W$  which parameterizes a loss function  $L$

General Template of a GD-style algorithm:

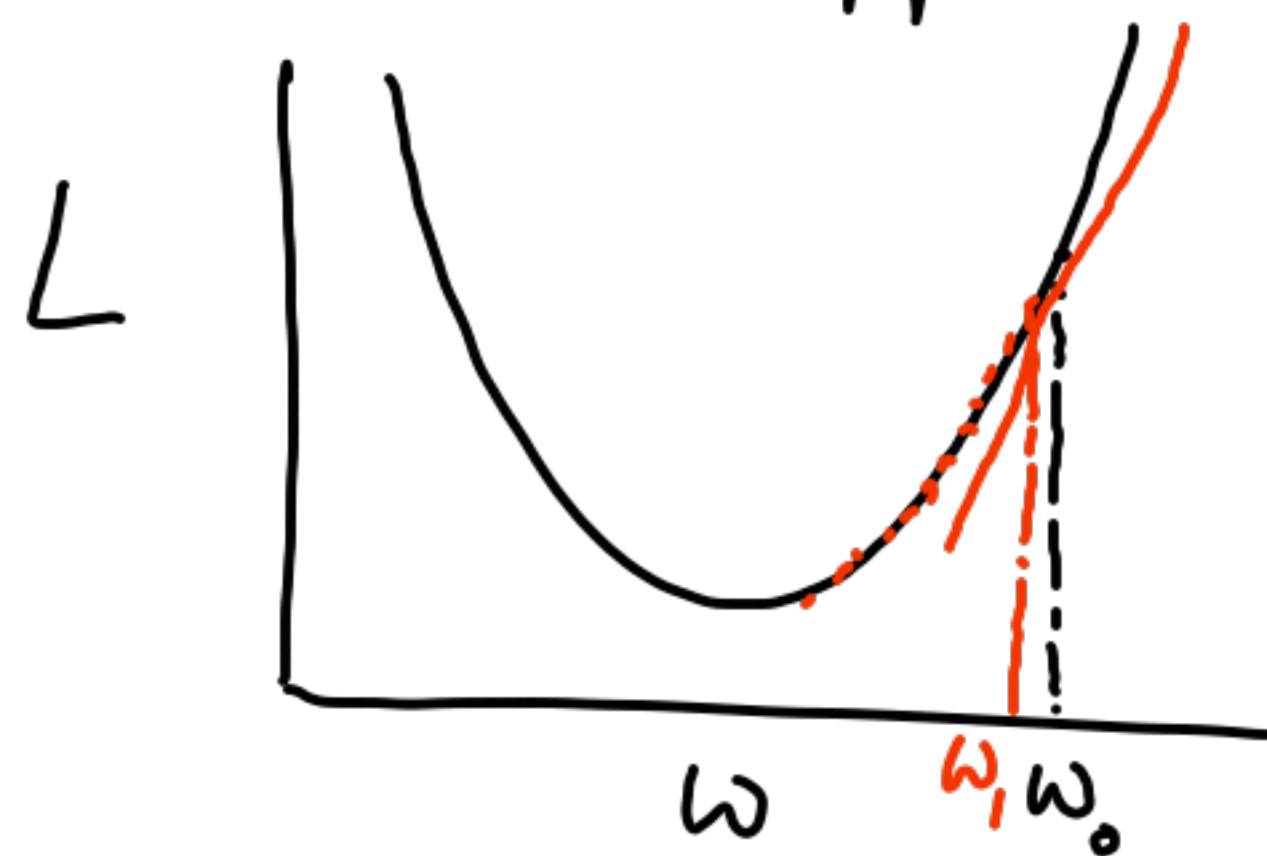
- Initialize  $W$
- Repeat
  - Choose a descent direction
  - Choose a step size
  - Update  $W$  using an update rule
- Exit repeat if some stopping criterion is met

### GRADIENT DESCENT

- $W \leftarrow W_0$  ( $W_0$  can be all 0's, random uniform, random Gaussian, etc.)
- Repeat
  - Descent direction:  $-\nabla_W L$
  - Step size is determined by  $\eta > 0$  ("LEARNING RATE")
  - $W_{t+1} \leftarrow W_t - \eta \nabla_W L$
- Exit when  $\|W_{t+1} - W_t\|^2 < \epsilon$

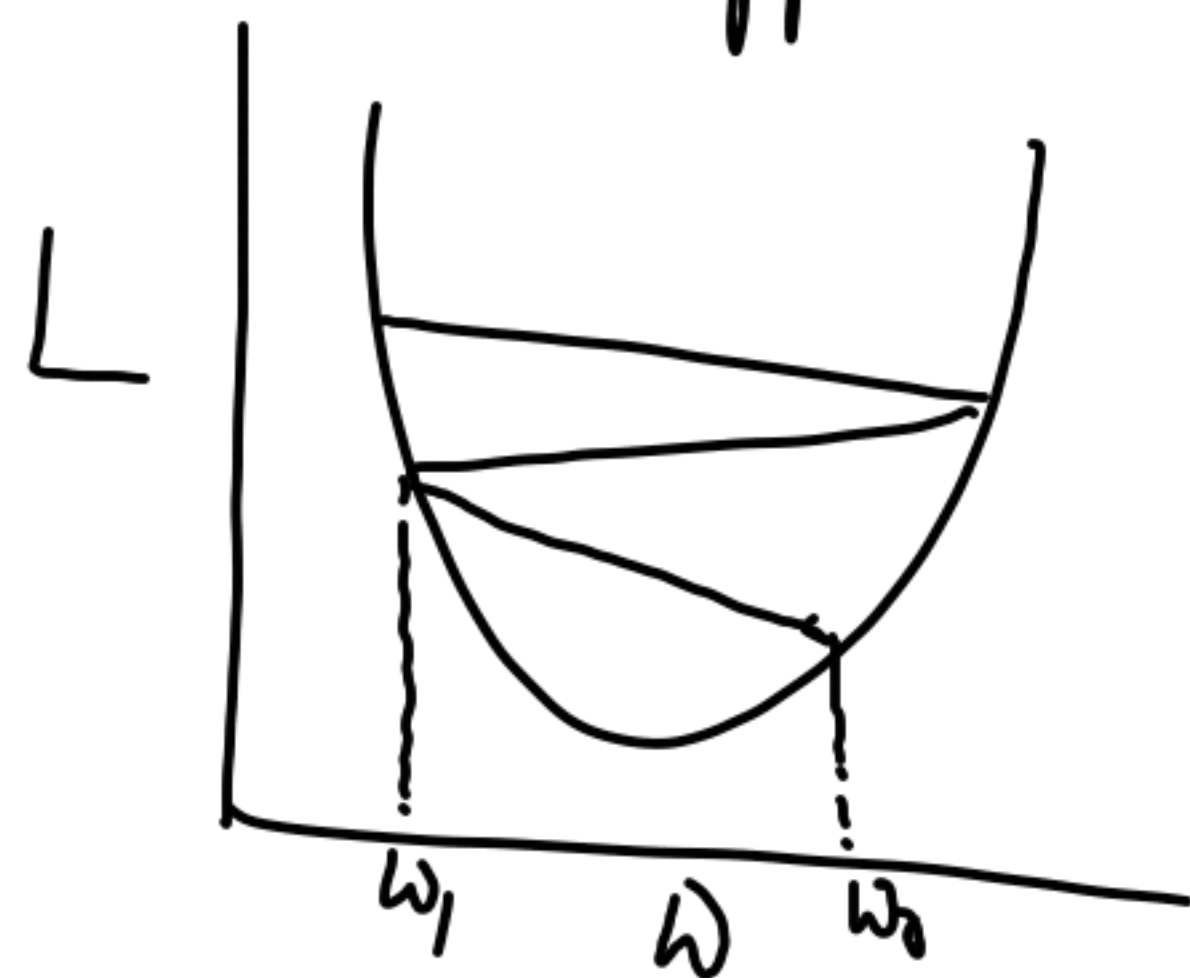


① What happens with GD if the step size is too small?



Going to take many iterations to converge

② What happens when the step size is too large?



Could delay arriving at the final solution  
At worst, the loss could diverge