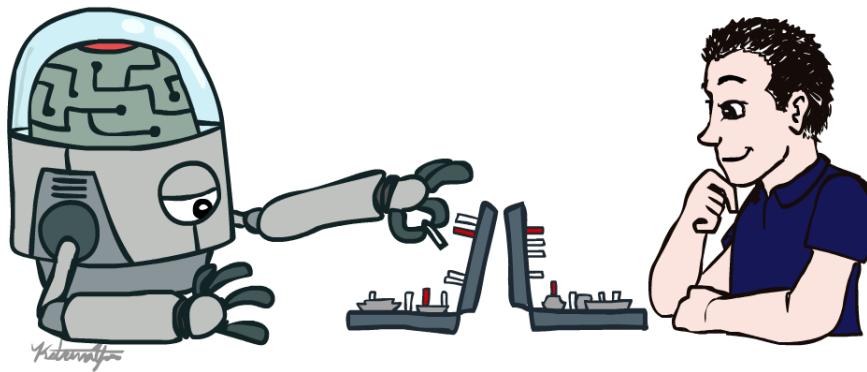


# CSE 3521: Introduction to Artificial Intelligence



[Many slides are adapted from the [UC Berkeley. CS188 Intro to AI](#) at UC Berkeley and previous CSE 3521 course at OSU.]



THE OHIO STATE UNIVERSITY

# Regression (curve fitting): what model?

- Relationship (function):

- Linear

- $f(x) = a x + b$       or       $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \underbrace{w[1]x[1] + \dots + w[D]x[D]}_{\text{D-dimensional input}} + b$

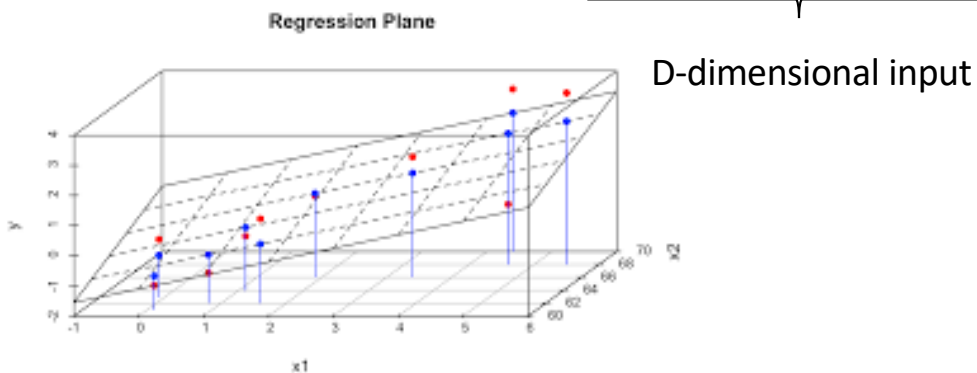
1-dimensional input

- Non-linear

- $f(x) = a x^2 + b x + c$

- Parameter Estimation

- Given data points and desired labels/target values  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$
  - Find values for parameters  $a, b, c, \mathbf{w}$ , etc.
  - Aim to find parameters that make  $f(\mathbf{x}_n)$  as close as possible to  $y_n$



# Regression (curve fitting): close?

---

- Difference or error between predictions and labels
  - Total error: error over all training data instances
- Different types of total error
  - Sum of Absolute Error (SAE)

$$E_1 = \sum_i |y_i - f(\mathbf{x}_i)|$$

- Maximum Error

$$E_\infty = \max_i |y_i - f(\mathbf{x}_i)|$$

- Sum of Squared Error (SSE) or residual sum of squares (RSS)

$$E_2 = \sum_i [y_i - f(\mathbf{x}_i)]^2$$

# Linear least squares

---

- $f$  is a linear function (linear combination of feature variables)
  - Linear in parameters: The function is a sum of terms  $x[d]$ , where parameters are *only* coefficients of those terms
$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = w[1]x[1] + \dots + w[D]x[D] + b$$
  - Thus, only works when we expect relationship is a line/plane/hyperplane
- Minimize sum of squared error
- Why SSE?
  - Short answer: The math is easier!
  - Other answers: You will see more after probabilistic modeling
- Estimation of model parameters derived from
  - Calculus
  - Linear Algebra

## Example: 1-dimensional data

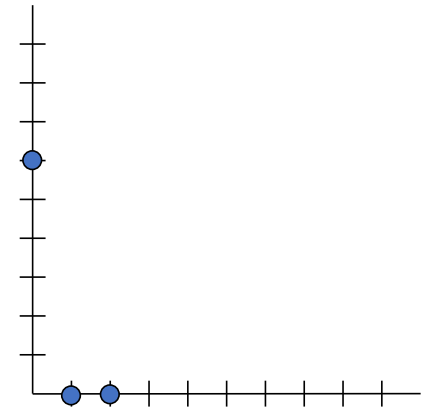
---

- Find a straight line through points (0,6), (1,0), and (2,0)
  - Model equation:  $y = wx + b$
- Need values for parameters  $w$ ,  $b$  that satisfy

$$6 = w \cdot 0 + b$$

$$0 = w \cdot 1 + b$$

$$0 = w \cdot 2 + b$$



# Linear system of equations

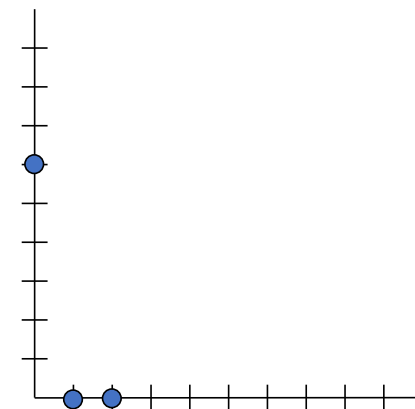
$$y = w \cdot x + b$$

$$6 = w \cdot 0 + b$$

$$0 = w \cdot 1 + b$$

$$0 = w \cdot 2 + b$$

$$\begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix} = w \cdot \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} + b \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \longrightarrow \begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} w \\ b \end{bmatrix}$$



$$\mathbf{y} = \tilde{\mathbf{X}}^T \tilde{\mathbf{w}}$$

Expand x with 1

Combine w with b

$\tilde{\mathbf{w}} = [w, b]$  unsolvable for these actual points! But...

# Minimizing the error

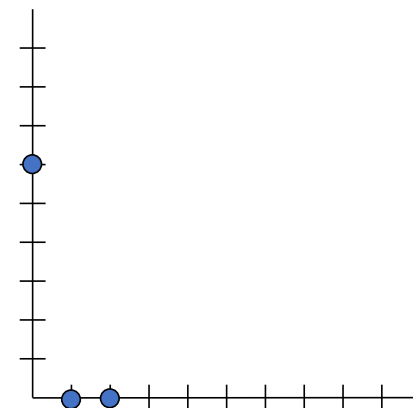
- Find a set of parameters ( $w, b$ ) that “best” fit the data

- Minimize the difference between  $\mathbf{y}$  and  $\tilde{\mathbf{X}}^T \tilde{\mathbf{w}}$

- $\begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} w \\ b \end{bmatrix}$

- “Least squares” approximation

- Look for the best  $\tilde{\mathbf{w}}$  that makes sum of the squared error (SSE) as small as possible



$$E = \sum_i (wx_i + b - y_i)^2 = \sum_i \left( [w, b] \begin{bmatrix} x_i \\ 1 \end{bmatrix} - y_i \right)^2 = \sum_i (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i - y_i)^2 = \sum_i (\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - y_i)^2$$
$$= [(\tilde{\mathbf{x}}_1^T \tilde{\mathbf{w}} - y_1), \dots, (\tilde{\mathbf{x}}_N^T \tilde{\mathbf{w}} - y_N)] \begin{bmatrix} \tilde{\mathbf{x}}_1^T \tilde{\mathbf{w}} - y_1 \\ \vdots \\ \tilde{\mathbf{x}}_N^T \tilde{\mathbf{w}} - y_N \end{bmatrix} = \|\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y}\|_2^2$$

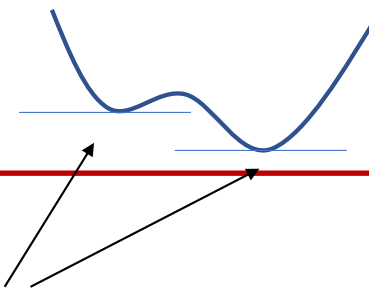
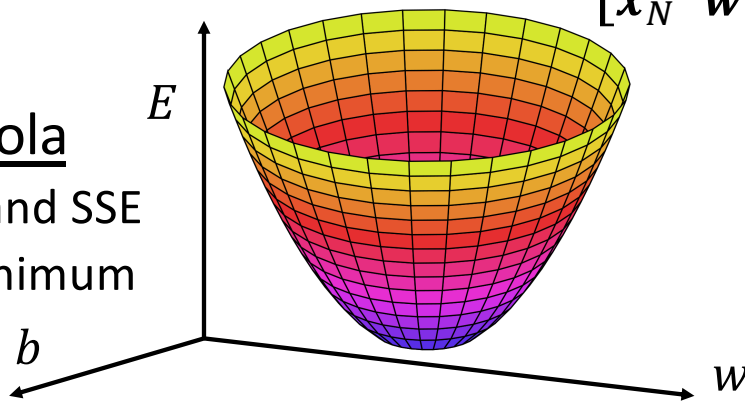
# Minimizing the error

- Many error functions are minimized by calculus
  - Derivatives (gradients) are zero when function at a (local) minimum

- Error function

$$E = \sum_i (wx_i + b - y_i)^2 = \sum_i \left( [w, b] \begin{bmatrix} x_i \\ 1 \end{bmatrix} - y_i \right)^2 = \sum_i (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i - y_i)^2$$
$$= [(\tilde{\mathbf{x}}_1^T \tilde{\mathbf{w}} - y_1), \dots, (\tilde{\mathbf{x}}_N^T \tilde{\mathbf{w}} - y_N)] \begin{bmatrix} \tilde{\mathbf{x}}_1^T \tilde{\mathbf{w}} - y_1 \\ \vdots \\ \tilde{\mathbf{x}}_N^T \tilde{\mathbf{w}} - y_N \end{bmatrix} = \|\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y}\|_2^2$$

- $E$  has the shape of a parabola
  - Due to linear the function and SSE
  - Local minimum = global minimum





# Minimizing the error

---

- We have two unknowns ( $w, b$ ), therefore we have two derivatives
  - Both derivatives are zero at the function minimum
- Compute the partial derivatives and set to 0

$$\begin{aligned}\frac{\partial E}{\partial w} &= \frac{\partial}{\partial w} \sum_i (w \cdot x_i + b - y_i)^2 = 0 \\ &= 2 \sum_i (w \cdot x_i + b - y_i) \cdot x_i = 0\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial b} &= \frac{\partial}{\partial b} \sum_i (w \cdot x_i + b - y_i)^2 = 0 \\ &= 2 \sum_i (w \cdot x_i + b - y_i) = 0\end{aligned}$$

# Minimizing the error

---

- Substitute in data

Recall:  
(0, 6), (1, 0), (2, 0)

$$\begin{aligned}\frac{\partial E}{\partial w} &= 2 \sum_i (w \cdot x_i + b - y_i) \cdot x_i = 0 \\ &= (w \cdot 0 + b - 6) \cdot 0 + (w \cdot 1 + b - 0) \cdot 1 + (w \cdot 2 + b - 0) \cdot 2 = 0 \\ &= 5w + 3b = 0\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial b} &= 2 \sum_i (w \cdot x_i + b - y_i) = 0 \\ &= (w \cdot 0 + b - 6) + (w \cdot 1 + b - 0) + (w \cdot 2 + b - 0) = 0 \\ &= 3w + 3b - 6 = 0 \\ &3w + 3b = 6\end{aligned}$$

# Minimizing the error

---

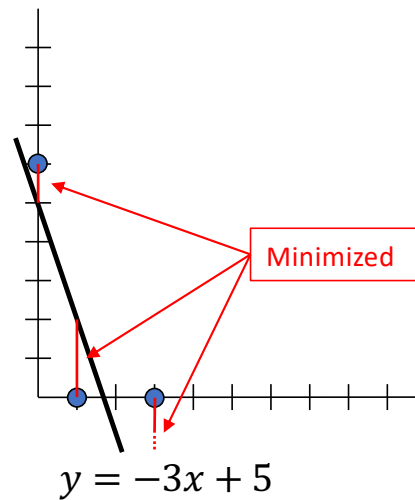
- Solve for model parameters

$$5w + 3b = 0$$

$$3w + 3b = 6$$

$$w = -3$$

$$b = 5$$



# Linear algebra formulation

---

- Solving systems of equations can also be done with linear algebra:

$$\begin{bmatrix} 5 & 3 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} 0 \\ 6 \end{bmatrix}$$

$$\begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} 5 & 3 \\ 3 & 3 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 6 \end{bmatrix} = \begin{bmatrix} -3 \\ 5 \end{bmatrix}$$

- Recall we originally formulated our problem by  $\mathbf{y}$  and  $\tilde{\mathbf{X}}^T \tilde{\mathbf{w}}$ . Is there some relationship with the above?

# Linear algebra formulation

---

- Yes!

$$\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T = \begin{bmatrix} 0 & 1 & 2 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 3 \\ 3 & 3 \end{bmatrix}$$

$$\tilde{\mathbf{X}}\mathbf{y} = \begin{bmatrix} 0 & 1 & 2 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 6 \end{bmatrix} \quad \tilde{\mathbf{w}} = \begin{bmatrix} w \\ b \end{bmatrix}$$

- We can rewrite the previous solution as:

$$\begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} 5 & 3 \\ 3 & 3 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 6 \end{bmatrix} = \begin{bmatrix} -3 \\ 5 \end{bmatrix}$$

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

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

↑  
Optimal solution

# Linear algebra formulation

---

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

This is equivalent to the partial derivative formulation we derived earlier.

In other words, linear least squares is as simple as correctly forming your  $\tilde{\mathbf{X}}$  matrix and  $\mathbf{y}$  vector and then **applying standard linear algebra operations to obtain a** closed-form solution

# Linear algebra formulation

---

- Another way of derivation (for D-dimensional data)

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = w[1]x[1] + \dots + w[D]x[D] + b = [\mathbf{w}^T, b] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

$$\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_N] = \begin{bmatrix} | & & | \\ \tilde{\mathbf{x}}_1 & \dots & \tilde{\mathbf{x}}_N \\ | & & | \end{bmatrix}$$

$$\begin{aligned} E &= \|\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y}\|_2^2 = (\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y})^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y}) = (\tilde{\mathbf{w}}^T \tilde{\mathbf{X}} - \mathbf{y}^T)(\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y}) \\ &= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2\mathbf{y}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} + \mathbf{y}^T \mathbf{y} = \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2\mathbf{y}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} + \text{Const.} \end{aligned}$$

$$\nabla_{\tilde{\mathbf{w}}} E = \begin{bmatrix} \frac{\partial E}{\partial \tilde{\mathbf{w}}[1]} \\ \vdots \\ \frac{\partial E}{\partial \tilde{\mathbf{w}}[D]} \end{bmatrix} = 2\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2\tilde{\mathbf{X}} \mathbf{y} = 0$$

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

# Linear algebra formulation

---

- Works for any model with linear combination of parameters
- Need enough examples for unknowns (i.e., parameters)
  - $D+1$  unknowns  $\rightarrow D+1$  equations/points (unique)
    - Often need many more than  $D$  to get a good result
      - Noise
      - Overfitting
- Another explanation:
  - $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$  is  $(D+1)$ -by- $(D+1)$
  - $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}$  exists only if  $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$  is of full rank:  $N \geq D + 1$



# Extension: non-linear models (1-dimension)

---

- Consider polynomials...

$$y = g + cx + bx^2 + ax^3 = g \cdot 1 + \dots$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad \tilde{\mathbf{X}}^T = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{bmatrix} \quad \tilde{\mathbf{w}} = \begin{bmatrix} g \\ c \\ b \\ a \end{bmatrix}$$

$N \geq 4$

Though  $g$  is a constant term, you can think of it as being the coefficient for 1 (equivalently  $x^0$  in polynomials)

Note that the order in  $\tilde{\mathbf{w}}$  doesn't matter, as long as the rows of  $\tilde{\mathbf{w}}$  correspond with the columns of  $\tilde{\mathbf{X}}^T$

## Extension: non-linear models (2-dimension)

---

- Consider 2<sup>nd</sup>-order polynomials...

$$y = b + w[1]x[1] + w[2]x[2] + w[3](x[1]x[2]) + w[4](x[1]^2) + w[5](x[2]^2)$$

- How many training data are needed?

# Extension: non-linear models

---

- Or non-polynomial terms...
  - But parameters are still linear in terms of these terms

$$y = a \frac{1}{x} + b \log x + c e^x$$

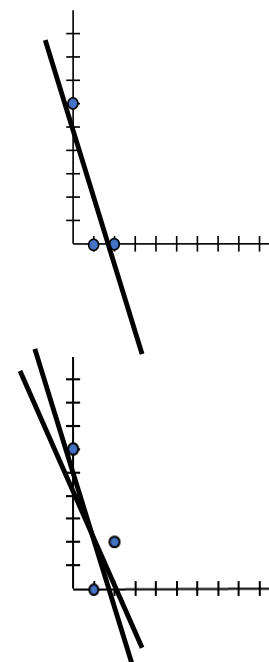
$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad \tilde{\mathbf{X}}^T = \begin{bmatrix} 1/x_1 & \log x_1 & e^{x_1} \\ 1/x_2 & \log x_2 & e^{x_2} \\ \vdots & \vdots & \vdots \\ 1/x_N & \log x_N & e^{x_N} \end{bmatrix} \quad \tilde{\mathbf{w}} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

$N \geq 3$

# Problems

---

- Requires relationship to be *linear*
- Matrix inverse is expensive for large numbers of parameters (complex models)
  - In general,  $O((D + 1)^3)$
  - Sometime, matrix inverse might even not exist
- Outliers (noise in the “label” or “input features”)
  - Data point(s) which don't match the pattern
  - Squared error is large, large influence on SSE
  - Thus, outliers can significantly affect results



# Notes

---

- In different textbooks or papers

- $b$  is sometimes denoted as  $w[0]$ ;
- $\mathbf{w}$  is called weights or parameters or parameter vector, and  $b$  as bias;
- sometimes,  $\tilde{\mathbf{w}}$  is called parameters, too;
- sometimes, people write  $\mathbf{w}$ , but they actually mean  $\tilde{\mathbf{w}}$

○ sometimes, people write  $\tilde{\mathbf{X}}$  (or  $\mathbf{X}$ ) to denote  $\begin{bmatrix} -\tilde{\mathbf{x}}_1^T & - \\ \vdots & \\ -\tilde{\mathbf{x}}_N^T & - \end{bmatrix}$  (or  $\begin{bmatrix} -\mathbf{x}_1^T & - \\ \vdots & \\ -\mathbf{x}_N^T & - \end{bmatrix}$ )

- So please pay attention to contexts when you read papers, textbooks, or assigned reading material.

# Summary

- Linear regression (LR)

- Fit a linear line to labeled data:

- For one-dimensional cases,  $f(x) = wx + b$
    - For multi-dimensional cases,  $f(x) = \sum_{d=1}^D w[d]x[d] + b$   
 $= f(x) = \mathbf{w}^T \mathbf{x} + b$

- Loss/error: absolute, square, etc.

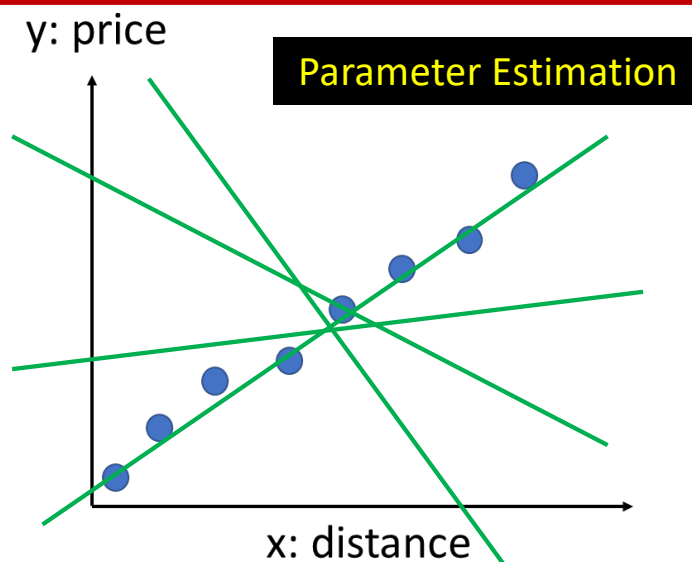
- Parameter estimation for LR

- For the sum of square error:  $E = \sum_i [y_i - f(x_i)]^2$
  - The closed-form solution that minimizes  $E$  (if we have the inverse)

- $\tilde{\mathbf{w}}^* = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\tilde{\mathbf{X}}\mathbf{y}$ ,

$$\tilde{\mathbf{w}}^* = \begin{bmatrix} \mathbf{w}^* \\ b^* \end{bmatrix}, \tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_N] = \begin{bmatrix} \mathbf{x}_1 & & \mathbf{x}_N \\ 1 & \dots & 1 \end{bmatrix}$$

- Careful about how the data is represented



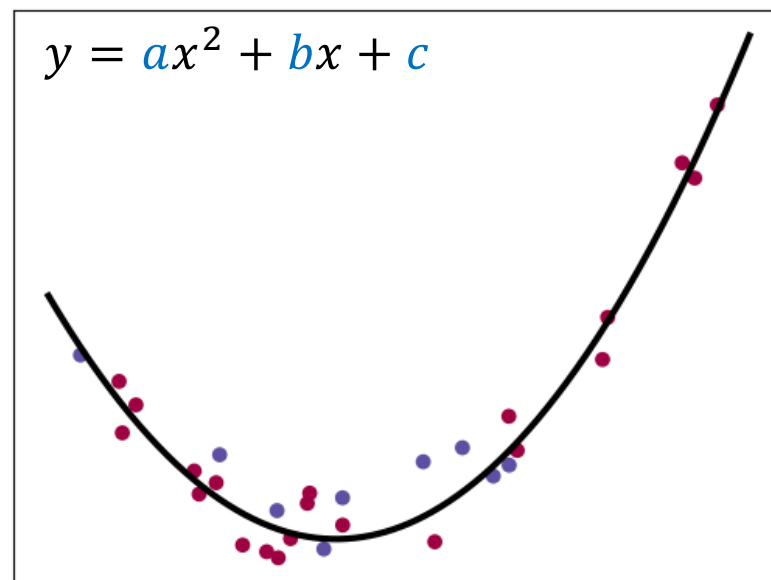
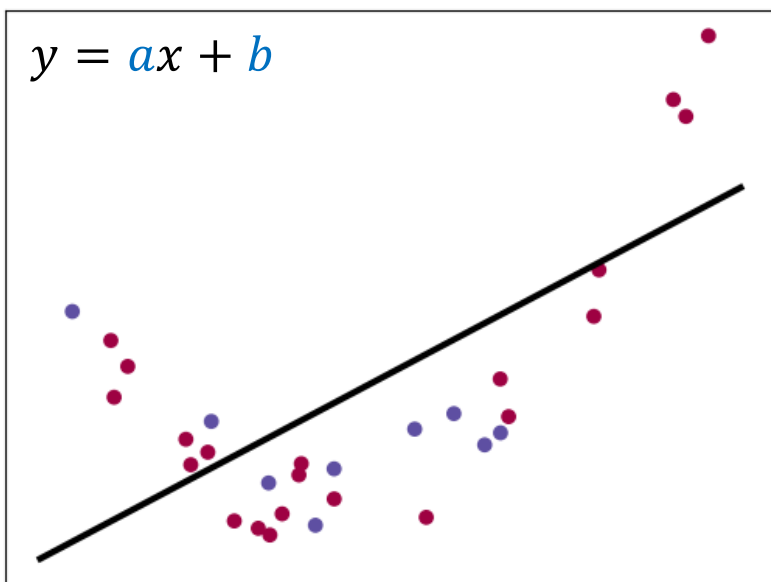
Concatenate  $b$  at the end of the column vector  $\mathbf{w}$

Concatenate 1 at the end of each column vector

# Recap: 1-dimensional case

---

- The closed-form solution is applicable to **some** non-linear regression
  - Polynomial, or when  $f(x_i)$  is the linear combination of terms derived by  $x$ .
  - $a$ ,  $b$ ,  $c$  just are just to indicate what to estimate. You can change the notations.



## Recap: 1-dimensional case

---

$$y = g + cx + bx^2 + ax^3$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad \tilde{\mathbf{X}} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_N \\ x_1^2 & x_2^2 & \cdots & x_N^2 \\ x_1^3 & x_2^3 & \cdots & x_N^3 \end{bmatrix} \quad \tilde{\mathbf{X}}^T = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{bmatrix} \quad \tilde{\mathbf{w}} = \begin{bmatrix} g \\ c \\ b \\ a \end{bmatrix}$$

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

$x_i$ : the i-th data instance

$x[d]$ : the d-th dimension of a data instance

$x_i[d]$ : the d-th dimension of the i-th data instance



# Topics

---

- Linear w.r.t.  $x$  or  $w$
- Non-linear least squares
  - Estimate parameters for non-linear functions w/o closed form solutions
- Gradient descent
  - Alternate solution methods
  - Not limited to SSE error functions

**We will focus on 1-dimensional cases first!**

# In finding the solution

---

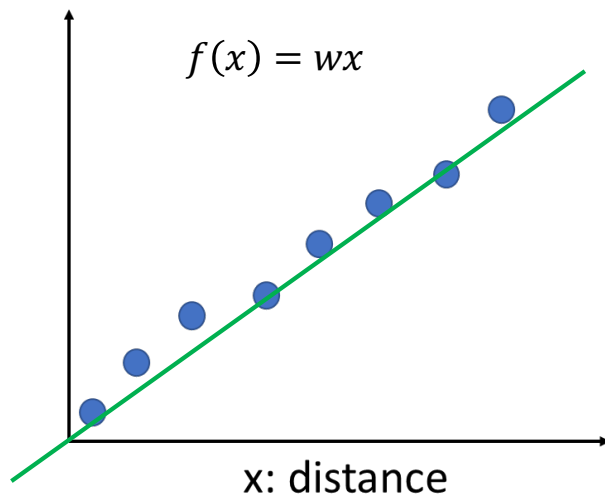
- The training data instances  $(x_i, y_i)$  are fixed
- The variables are the parameters to estimate
- A machine learning algorithm finds the parameters to minimize  $E$
- For the users, we care about given a future  $x$ , what the prediction will be



$$E = \sum_i [y_i - f(x_i)]^2$$

# Different figures

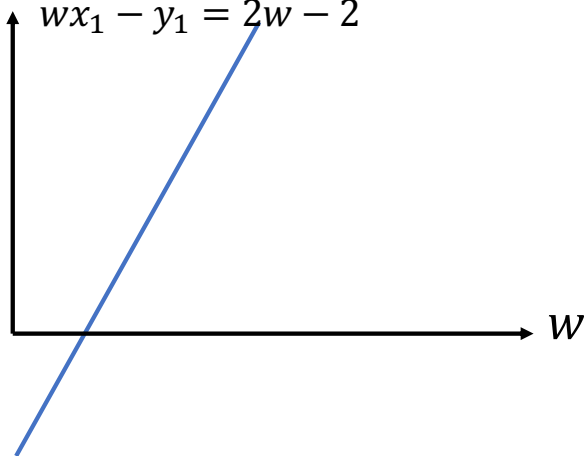
y: price



For simplicity, only estimate  $w$

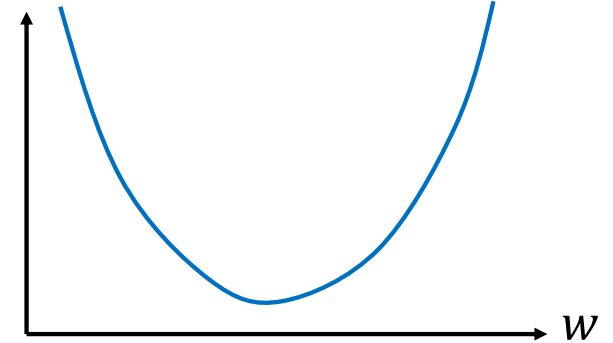
$f(x) - y$

$$wx_1 - y_1 = 2w - 2$$



Let  $(x_1, y_1) = (2, 2)$

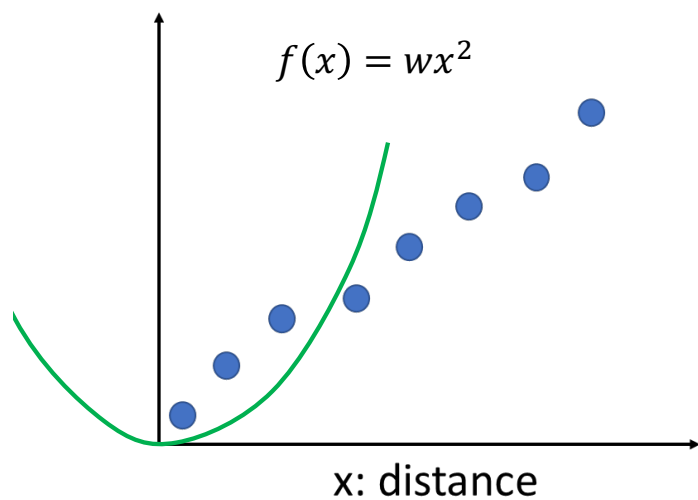
$$E(w) = \sum_i (wx_i - y_i)^2$$



Parabolic!

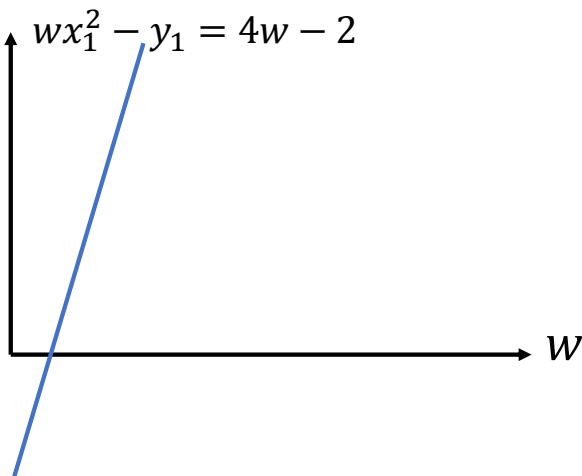
# Different figures

y: price



Nonlinear regression as  $f(x) = wx^2$  is not a straight line w.r.t.  $x$ .

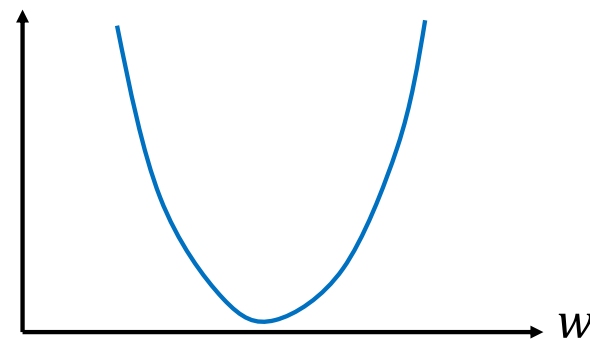
$$f(x) - y$$
$$wx_1^2 - y_1 = 4w - 2$$



Let  $(x_1, y_1) = (2, 2)$

Still a straight line w.r.t.  $w$ !

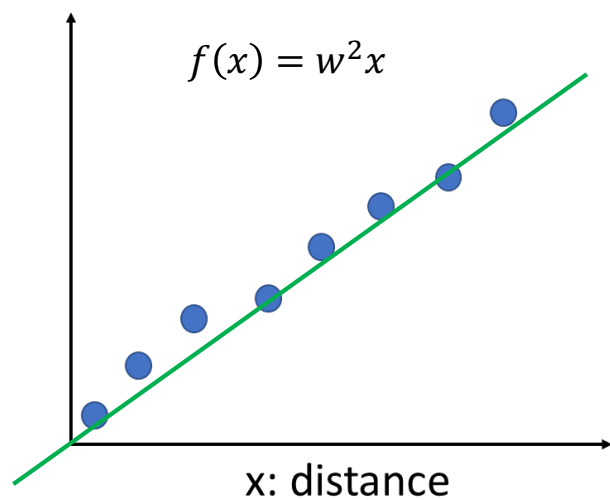
$$E(w) = \sum_i (wx_i^2 - y_i)^2$$



Parabolic!

# Different figures

y: price



Linear regression as  $f(x) = w^2x$  is a straight line w.r.t.  $x$ .

$f(x) - y$

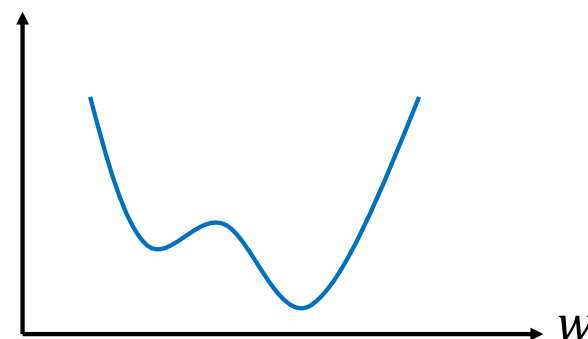
$$w^2x_1 - y_1 = 2w^2 - 2$$



Let  $(x_1, y_1) = (2, 2)$

Not a straight line w.r.t.  $w$ !

$$E(w) = \sum_i (w^2x_i - y_i)^2$$



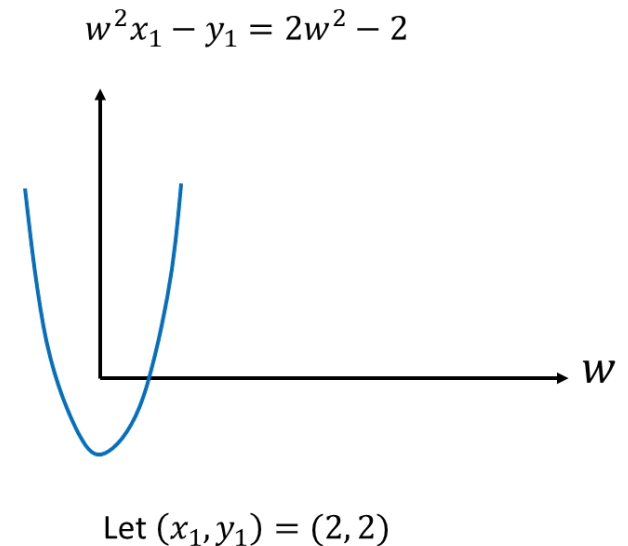
Not parabolic!

# When ...

---

To indicate that  $w$  is the parameter!

- Given a random  $(x, y)$ , when the prediction function  $f(x; w)$  or  $f(x; w) - y$  is not a straight line (i.e., nonlinear) **w.r.t.  $w$** 
  - there may be no closed form solution
  - $E$  may not be parabolic anymore
- The closed form solution  $\tilde{\mathbf{w}}^* = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\tilde{\mathbf{X}}\mathbf{y}$  exists
  - for the SSE loss  $E(w) = \sum_i (f(x; w) - y_i)^2$
  - for  $f(x; w)$  linear w.r.t.  $w$
- **We say linear or nonlinear regression, w.r.t.  $x$**



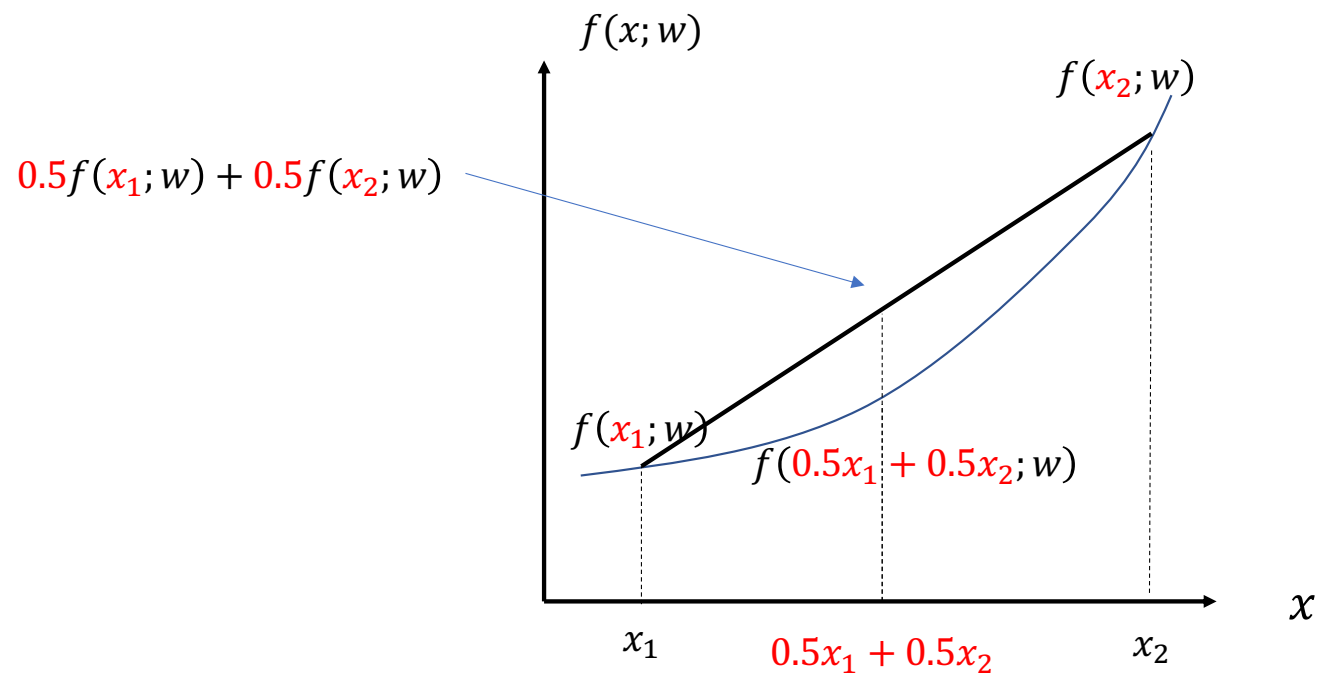
# How to check if $f(x; w)$ linear w.r.t. $x$ or $w$ ?

---

- If you want to check **w.r.t.  $x$** 
  - Ex.  $f(x; w) = wx$
- Find 2 random  $x$ :  $x_1$  and  $x_2$ 
  - Ex.  $x_1 = 1, x_2 = 2$
- Find 2 random **scalars** (even for  $D$ -dimensional cases):  $v_1$  and  $v_2$ 
  - Ex.  $v_1 = 1, v_2 = 2$
- If  $f(v_1x_1 + v_2x_2; w) = v_1f(x_1; w) + v_2f(x_2; w)$ 
  - $f(x; w)$  is linear w.r.t.  $x$ ; otherwise, not
  - Need to hold for all possible  $x_1, x_2, v_1, v_2$
- Practice:
  - Is  $f(x; w)$  linear w.r.t.  $x$ ?

# Interpretation 1

---





# Interpretation 2

$x$

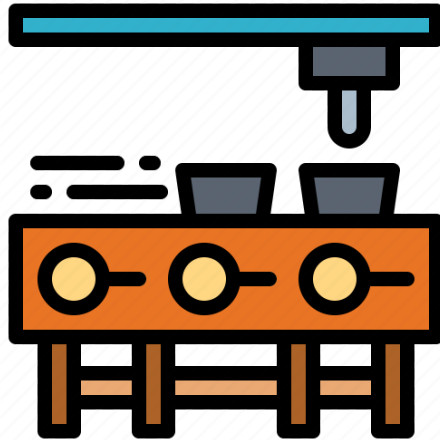
$f$

$f(x; w)$

$x_1$



$x_2$



$f(x_1; w)$



$f(x_2; w)$



$x_1 + x_2$



$f(x_1 + x_2; w)$



# How to check if $f(x; w)$ linear w.r.t. $x$ or $w$ ?

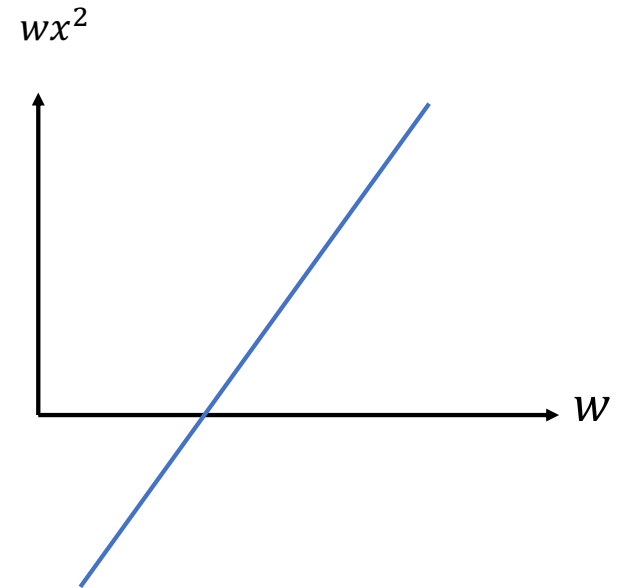
---

- If you want to check **w.r.t.  $w$** 
  - Ex.  $f(x; w) = wx^2$
- Find 2 random  $w$ :  $w_1$  and  $w_2$ 
  - Ex.  $w_1 = 1, w_2 = 2$
- Find 2 random **scalars** (even for  $D$ -dimensional cases):  $v_1$  and  $v_2$ 
  - Ex.  $v_1 = 1, v_2 = 2$
- If  $f(x; v_1w_1 + v_2w_2) = v_1f(x; w_1) + v_2f(x; w_2)$ 
  - $f(x; w)$  is linear w.r.t.  $w$ ; otherwise, not
  - Need to hold for all possible  $w_1, w_2, v_1, v_2$
- Practice:
  - Is  $f(x; w)$  linear w.r.t.  $w$ ?

# Practice

---

- $w x^2$ 
  - Is it linear w.r.t.  $w$ ?
  - Let  $w_1 = 1, w_2 = 2, v_1 = 1, v_2 = 2$
  - $(v_1 w_1 + v_2 w_2) x^2 = 5x^2$
  - $v_1(w_1 x^2) + v_2(w_2 x^2) = 5x^2$



# Practice

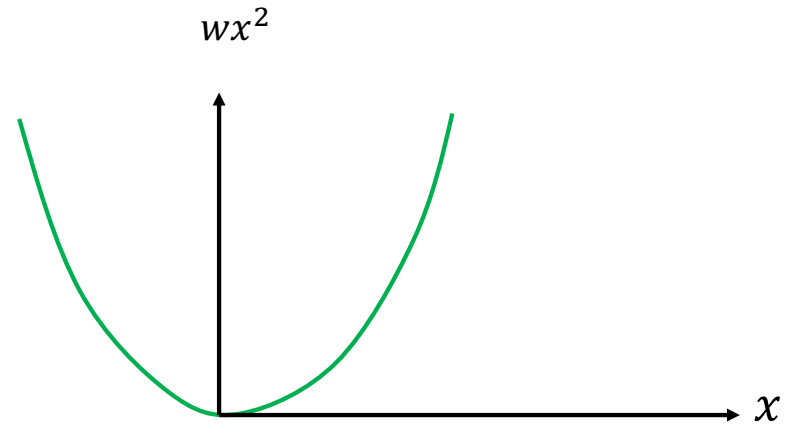
---

- $wx^2$ 
  - Is it linear w.r.t.  $x$ ?

# Practice

---

- $wx^2$ 
  - Is it linear w.r.t.  $x$ ?
  - Let  $x_1 = 1, x_2 = 2, v_1 = 1, v_2 = 2$
  - $w(v_1x_1 + v_2x_2)^2 = 25w$
  - $v_1(wx_1^2) + v_2(wx_2^2) = 9w$



# Topics

---

- Linear w.r.t.  $x$  or  $w$
- Non-linear least squares
  - Estimate parameters for non-linear functions
- Gradient descent
  - Alternate solution methods
  - Not limited to SSE error functions

# Non-linear regression

---

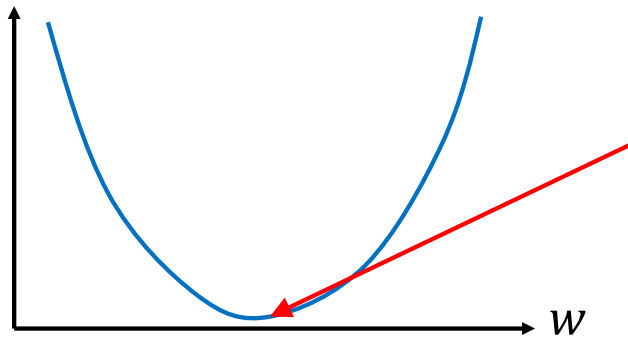
- For many cases when  $f(x; w)$  is not linear w.r.t.  $x$ ,  $f(x; w)$  is not linear w.r.t.  $w$
- Examples:
  - $f(x; w) = e^{wx}$
  - $f(x; w) = \sin(wx)$
- For non-linear regression,  $x$  and  $w$  very likely have different dimensionality
- Example:
  - $f(x; w) = e^{w[1]x} + \sin(w[2]x)$
- We will say  $x$  is D-dimensional and  $w$  is K-dimensional

# When $f(x; w)$ is nonlinear w.r.t. $w$

---

- Squared-error function is no longer parabolic w.r.t. the parameters
  - Can not assume local minimums of error functions are always global minimums
- Often no closed-form solution
  - Need numerical approximation

$$E(w) = \sum_i (wx_i - y_i)^2$$



- ❖  $\frac{\partial E}{\partial w} = 0$  means global minimums
- ❖ Closed-forms to find  $w^*$  for  $\frac{\partial E}{\partial w} = 0$

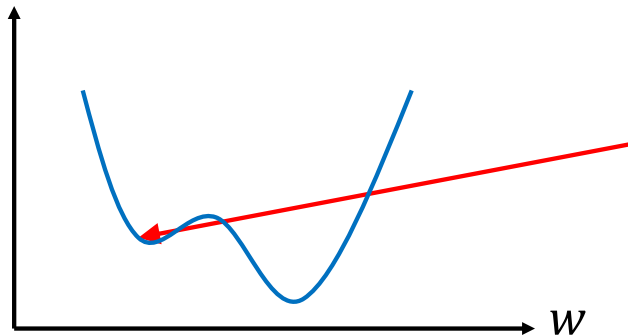


# When $f(x; w)$ is nonlinear w.r.t. $w$

---

- Squared-error function is no longer parabolic w.r.t. the parameters
  - Can not assume local minimums of error functions are always global minimums
- Often no closed-form solution
  - Need numerical approximation

$$E(w) = \sum_i (f(x_i; w) - y_i)^2$$



- ❖  $\frac{\partial E}{\partial w} = 0$  does not mean global minimums
- ❖ No closed-forms to find  $w^*$  for  $\frac{\partial E}{\partial w} = 0$

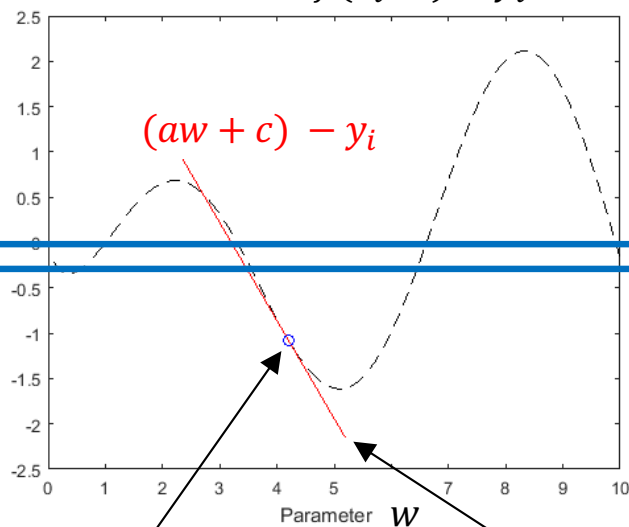
# How to solve? Approximation via linearization

---

- The problem is the nonlinear  $f(\mathbf{x}; \mathbf{w})$  w.r.t. parameters  $\mathbf{w}$ 
  - So non-parabolic squared error function  $E$  w.r.t. parameters  $\mathbf{w}$
- Idea: force the model  $f$  be **linear** w.r.t. parameters
  - Local – around current guess for parameters
    - This means it must be possible to approximate a function with a line locally
  - Use truncated Taylor series
  - Solve linear system as with linear least squares
  - Update current guess and repeat
- Note:  $\mathbf{w}$  may have a very different dimensionality from  $\mathbf{x}$

# Local Linearization (1-dim cases)

residual:  $f(x_i; w) - y_i$



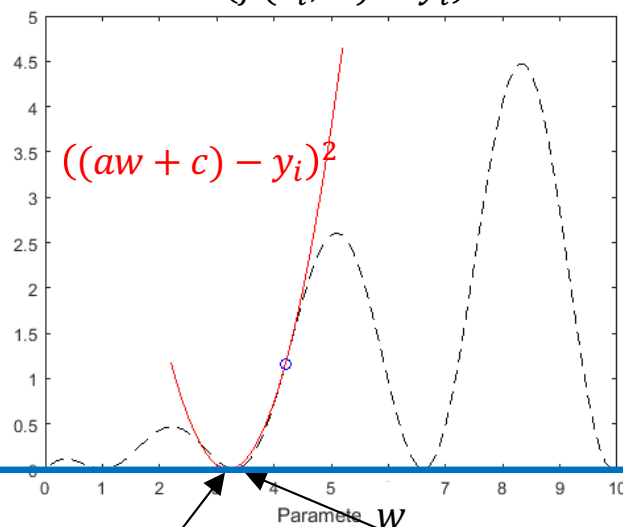
Initial guess:  $w' = 4.2$

Linearize around guess

Note:  $a, c$  will be known;  $w$  is to be solved

Want to find where the **red line**  
intersect the **blue line**!

$(f(x_i; w) - y_i)^2$



New guess:  
 $w' = 3.2031$

True minima:  
 $w^* = 3.3069$

Very close even after just one iteration!

You know how to  
find  $w'$  to minimize  
 $((aw + c) - y_i)^2$

# Local Linearization (a 1-dim case for $w$ and $x$ )

- Recall the model (per data point  $i$ )

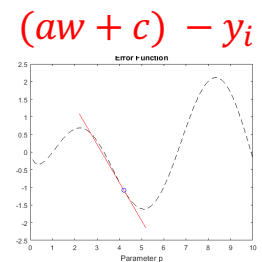
$$f(x_i; w)$$

← Read:  $f$  of  $x_i$  given parameters  $w$

- Taylor series expansion (assume  $w$  is one-dimensional)

$$f(x_i; \mathbf{w}) = \underbrace{f(x_i; \mathbf{w}')}_{\text{Constant term}} + \underbrace{\frac{df(x_i; \mathbf{w}')}{d\mathbf{w}} (\mathbf{w} - \mathbf{w}')}_{\text{Linear (first-order) term}} + \text{higher-order terms (H.O.T.)}$$

Variable to solve (points to  $\mathbf{w}$ )  
 Current guess (points to  $\mathbf{w}'$ )  
 Variation "around"  $\mathbf{w}'$  (points to  $\mathbf{w} - \mathbf{w}'$ )



- Linearize: drop higher-than-linear terms

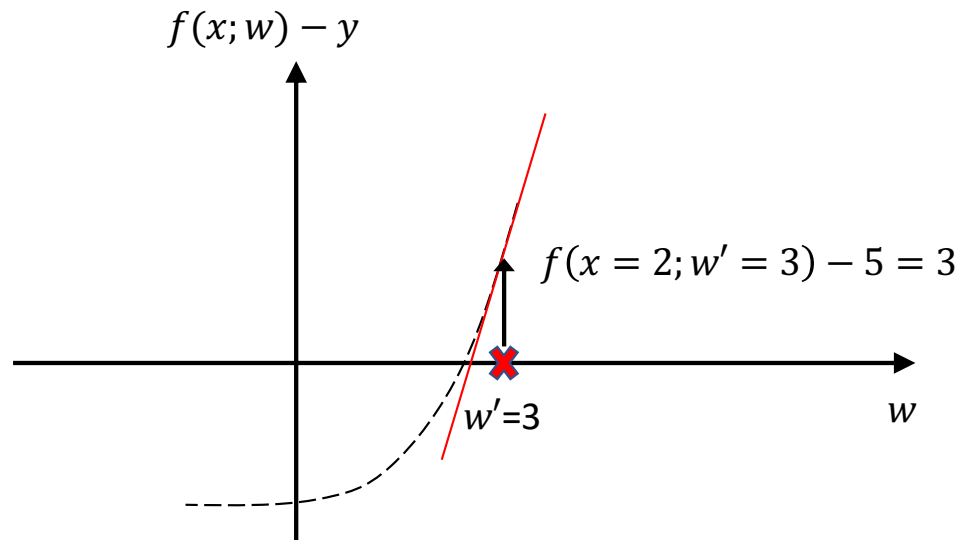
$$f(x_i; \mathbf{w}) - y_i \approx f(x_i; \mathbf{w}') + \frac{df(x_i; \mathbf{w}')}{d\mathbf{w}} (\mathbf{w} - \mathbf{w}') - y_i = J(\mathbf{w}') \Delta \mathbf{w} - \Delta y_i(\mathbf{w}')$$

$$\Delta y_i(\mathbf{w}') = y_i - f(x_i; \mathbf{w}') \quad \Delta \mathbf{w} = \mathbf{w} - \mathbf{w}' \quad \text{Jacobian: } J_i(\mathbf{w}') = \frac{df(x_i; \mathbf{w}')}{d\mathbf{w}}$$

# Local Linearization (a 1-dim case for $w$ and $x$ )

- Example:

- $f(x; w) = x^w$
- $(x, y) = (2, 5)$



- Initial guess

- $w' = 3$
- $\frac{df(x=2; w'=3)}{dw} = \log(2) \times 2^3 = 5.545$
- $J(w')\Delta w - \Delta y_i(w') = 5.545(w - 3) - (5 - 8) = 5.545(w - 3) + 3$

# Minimizing the “sum of square error”

---

- Then, sum of squared error:

$$E'(\Delta w) = \sum_i [J_i(w')\Delta w - \Delta y_i(w')]^2 = \sum_i [J_i\Delta w - \Delta y_i]^2$$

Note, to simplify, we are not stating the dependencies on the current guess  $w'$ ,  
but it is always there!

- “Least squares” approximation (take derivative and set it to 0)

$$0 = \frac{d}{d\Delta w} E'(\Delta w) = \sum_i \frac{d}{d\Delta w} [J_i\Delta w - \Delta y_i]^2$$

$$= 2 \sum_i [J_i\Delta w - \Delta y_i] \frac{d}{d\Delta w} [J_i\Delta w - \Delta y_i]$$

$$= 2 \sum_i [J_i\Delta w - \Delta y_i] J_i$$

# Minimizing the sum of square error (cont.)

---

- “Least squares” approximation (cont.)

$$\sum_i J_i \Delta y_i = \sum_i J_i J_i \Delta w$$

- Rewriting in matrix form gives:

$$J^T \Delta \mathbf{y} = (J^T J) \Delta \mathbf{w} \quad \Delta \mathbf{w} = (J^T J)^{-1} J^T \Delta \mathbf{y}$$

$$J = \begin{bmatrix} J_1 \\ \vdots \\ J_N \end{bmatrix}, \quad \Delta \mathbf{y} = \begin{bmatrix} \Delta y_1 \\ \vdots \\ \Delta y_N \end{bmatrix}$$

Reminder: These should be  $J_i(w')$  and  $\Delta y_i(w')$  as they change with the current guess  $w'$ !

- Note: Compare with linear least squares
  - LLSQ reminder:  $\tilde{\mathbf{X}} \mathbf{y} = (\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T) \tilde{\mathbf{w}}$

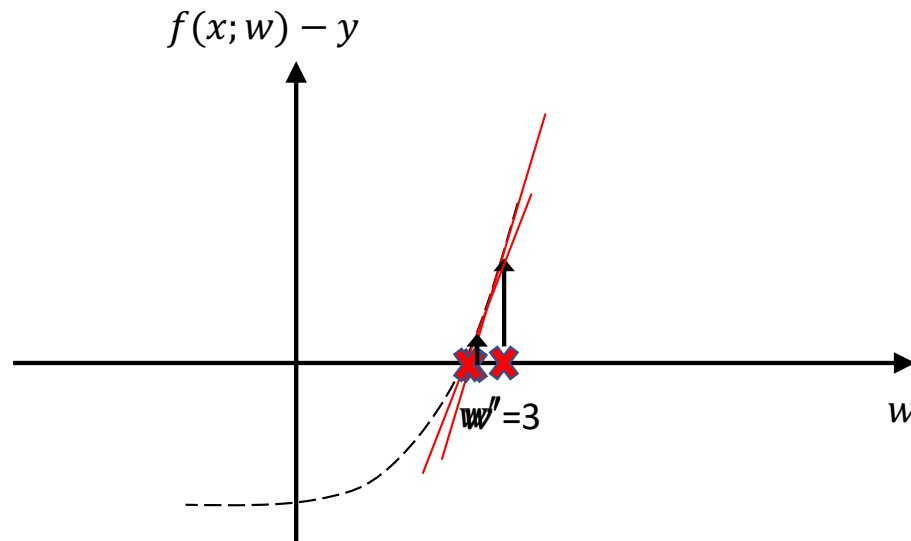
# Local Linearization (a 1-dim case for $w$ and $x$ )

- Example:

- $f(x; w) = x^w$
- $(x, y) = (2, 5)$

- Algorithm

- Initial:  $w' = 3$
- If not converge:
  - $w' := w' + \Delta w$





# Gauss-Newton Algorithm

---

1. Get initial guess for parameters  $\hat{w}^{(1)}$
2. For  $t=1...\infty$  (or until we decide to stop)

1. Find error for current guess ( $\hat{w}^{(t)}$ )

$$\Delta y_i = y_i - f(\mathbf{x}_i; \hat{w}^{(t)})$$

2. Find Jacobian around current guess

$$J_i = \frac{df(\mathbf{x}_i; \hat{w}^{(t)})}{dw}$$

3. Calculate change in guess

$$\Delta \mathbf{w} = (J^T J)^{-1} J^T \Delta \mathbf{y}$$

4. Make new guess

$$\hat{w}^{(t+1)} = \hat{w}^{(t)} + \Delta \hat{w}$$

# Example

---

- Log-linear model

$$f(x; \mathbf{w}) = e^{wx}$$

- Data points

- $x_1 = 0.5, y_1 = 1$

- $x_2 = 1.5, y_2 = 2$

- $x_3 = 3, y_3 = 3$

## Example - Jacobian

---

- Equation is simple, so closed-form Jacobian is possible

$$\frac{df(x; \mathbf{w})}{dw} = x e^{wx}$$

$$J(\mathbf{w}) = \begin{bmatrix} \frac{df(x_1; \mathbf{w})}{dw} \\ \frac{df(x_2; \mathbf{w})}{dw} \\ \frac{df(x_3; \mathbf{w})}{dw} \end{bmatrix} = \begin{bmatrix} 0.5 e^{0.5w} \\ 1.5 e^{1.5w} \\ 3 e^{3w} \end{bmatrix}$$

# Example

---

- Initial  $\hat{w}^{(1)} = 1$

$$J(\hat{w}^{(1)}) = \begin{bmatrix} \frac{df(x_1; \hat{w}^{(1)})}{dw} \\ \frac{df(x_2; \hat{w}^{(1)})}{dw} \\ \frac{df(x_3; \hat{w}^{(1)})}{dw} \end{bmatrix} = \begin{bmatrix} 0.5 e^{0.5 \hat{w}^{(1)}} \\ 1.5 e^{1.5 \hat{w}^{(1)}} \\ 3 e^{3 \hat{w}^{(1)}} \end{bmatrix} = \begin{bmatrix} 0.5 e^{0.5} \\ 1.5 e^{1.5} \\ 3 e^{3} \end{bmatrix}$$

$$\begin{aligned} x_1 &= 0.5, y_1 = 1 \\ x_2 &= 1.5, y_2 = 2 \\ x_3 &= 3, y_3 = 3 \end{aligned}$$

$$\Delta \mathbf{y}(\hat{w}^{(1)}) = \begin{bmatrix} y_1 - f(\mathbf{x}_1; \hat{w}^{(1)}) \\ y_2 - f(\mathbf{x}_2; \hat{w}^{(1)}) \\ y_3 - f(\mathbf{x}_3; \hat{w}^{(1)}) \end{bmatrix} = \begin{bmatrix} 1 - e^{0.5 \hat{w}^{(1)}} \\ 2 - e^{1.5 \hat{w}^{(1)}} \\ 3 - 3e^{3 \hat{w}^{(1)}} \end{bmatrix} = \begin{bmatrix} 1 - e^{0.5} \\ 2 - e^{1.5} \\ 3 - 3e^3 \end{bmatrix}$$

# Local Linearization (let $\mathbf{w}$ be $K$ -dim)

---

- Recall the model (per data point  $i$ )

$$f(\mathbf{x}; \mathbf{w})$$

- Taylor series expansion

$$\begin{aligned} f(\mathbf{x}; \mathbf{w}) &= f(\mathbf{x}_i; \mathbf{w}') + \sum_{k=1}^K \frac{\partial f(\mathbf{x}_i; \mathbf{w}')}{\partial w[k]} (w[k] - w'[k]) + \text{higher-order terms (H.O.T.)} \\ &= f(\mathbf{x}_i; \mathbf{w}') + \left[ \frac{\partial f(\mathbf{x}_i; \mathbf{w}')}{\partial w[1]} \quad \dots \quad \frac{\partial f(\mathbf{x}_i; \mathbf{w}')}{\partial w[K]} \right] \begin{bmatrix} (w[1] - w'[1]) \\ \vdots \\ (w[K] - w'[K]) \end{bmatrix} + \text{H.O.T.} \\ &= f(\mathbf{x}_i; \mathbf{w}') + \nabla_{\mathbf{w}} f(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}') + \text{H.O.T.} \end{aligned}$$

# Local Linearization (let $\mathbf{w}$ be $K$ -dim)

---

- Linearize: drop higher-than-linear terms

$$\begin{aligned} f(\mathbf{x}; \mathbf{w}) - y_i &\approx f(\mathbf{x}; \mathbf{w}') + \left\{ \sum_{k=1}^K \frac{\partial f(\mathbf{x}_i; \mathbf{w}')}{\partial w[k]} (w[k] - w'[k]) \right\} - y_i \\ &= \left\{ \sum_{k=1}^K J_{ik}(\mathbf{w}') \Delta \mathbf{w}[k] \right\} - \Delta y_i(\mathbf{w}') \end{aligned}$$

$$\Delta y_i(\mathbf{w}') = y_i - f(\mathbf{x}; \mathbf{w}')$$

$$\Delta \mathbf{w} = \mathbf{w} - \mathbf{w}'$$

$$\text{Jacobian matrix: } \mathbf{J} = \begin{bmatrix} \frac{\partial f(\mathbf{x}_1; \mathbf{w}')}{\partial w[1]} & \dots & \frac{\partial f(\mathbf{x}_1; \mathbf{w}')}{\partial w[K]} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{x}_N; \mathbf{w}')}{\partial w[1]} & \dots & \frac{\partial f(\mathbf{x}_N; \mathbf{w}')}{\partial w[K]} \end{bmatrix} \quad J_{ik} = J[i, k] = \frac{\partial f(\mathbf{x}_i; \mathbf{w}')}{\partial w[k]}$$

## Minimizing the “sum of square error” (let $\mathbf{w}$ be $K$ -dim)

---

- Then, sum of squared error:

$$E'(\Delta \mathbf{w}) = \sum_i \left[ \left\{ \sum_{k=1}^K J_{ik}(\mathbf{w}') \Delta \mathbf{w}[k] \right\} - \Delta y_i(\mathbf{w}') \right]^2 = \sum_i \left[ \left\{ \sum_{k=1}^K J_{ik} \Delta \mathbf{w}[k] \right\} - \Delta y_i \right]^2$$

- “Least squares” approximation (take partial derivatives and set them to 0)

$$\begin{aligned} 0 &= \frac{\partial}{\partial \Delta \mathbf{w}[k]} E'(\Delta \mathbf{w}) = \sum_i \frac{\partial}{\partial \Delta \mathbf{w}[k]} \left[ \left\{ \sum_{k=1}^K J_{ik} \Delta \mathbf{w}[k] \right\} - \Delta y_i \right]^2 \\ &= 2 \sum_i \left[ \left\{ \sum_{j=1}^K J_{ij} \Delta \mathbf{w}[k] \right\} - \Delta y_i \right] \frac{\partial}{\partial \Delta \mathbf{w}[k]} \left[ \left\{ \sum_{k=1}^K J_{ik} \Delta \mathbf{w}[k] \right\} - \Delta y_i \right] \\ &= 2 \sum_i \left[ \left\{ \sum_{j=1}^K J_{ij} \Delta \mathbf{w}[k] \right\} - \Delta y_i \right] J_{ik} \end{aligned}$$

# Minimizing the “sum of square error” (let $\mathbf{w}$ be $K$ -dim)

---

- “Least squares” approximation (cont.)

$$\sum_i J_{ik} \Delta y_i = \sum_i \sum_{j=1}^K J_{ik} J_{ij} \Delta \mathbf{w}[k]$$

- Rewriting in matrix form gives:

$$\mathbf{J}^T \Delta \mathbf{y} = (\mathbf{J}^T \mathbf{J}) \Delta \mathbf{w}$$

Reminder: These should be  $J_i(\mathbf{w}')$  and  $\Delta y_i(\mathbf{w}')$  as they change with the current guess  $\mathbf{w}'$ !



# Gauss-Newton Algorithm (let $\mathbf{w}$ be $K$ -dim)

---

1. Get initial guess for parameters  $\hat{\mathbf{w}}^{(1)}$
2. For  $t=1...\infty$  (or until we decide to stop)

1. Find error for current guess (  $\hat{\mathbf{w}}^{(t)}$  )

$$\Delta y_i = y_i - f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)})$$

2. Find Jacobian around current guess

$$J_{ik} = \frac{\partial f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)})}{\partial w[k]}$$

3. Calculate change in guess

$$\Delta \mathbf{w} = (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \Delta \mathbf{y}$$

4. Make new guess

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

# Example

---

- Log-linear model

$$f(x; \mathbf{w}) = e^{w[1]x + w[2]}$$

- Data points

- $x_1 = 0.5, y_1 = 1.6487$

- $x_2 = 1.5, y_2 = 0.6065$

- $x_3 = 3, y_3 = 0.1353$

## Example - Jacobian

---

- Equation is simple, so closed-form Jacobian is possible

$$\frac{\partial f(x; \mathbf{w})}{\partial w[1]} = x e^{w[1]x+w[2]}$$

$$\frac{\partial f(x; \mathbf{w})}{\partial w[2]} = e^{w[1]x+w[2]}$$

$$J(\mathbf{w}) = \begin{bmatrix} \frac{\partial f(x_1; \mathbf{w})}{\partial w[1]} & \frac{\partial f(x_1; \mathbf{w})}{\partial w[2]} \\ \frac{\partial f(x_2; \mathbf{w})}{\partial w[1]} & \frac{\partial f(x_2; \mathbf{w})}{\partial w[2]} \\ \frac{\partial f(x_3; \mathbf{w})}{\partial w[1]} & \frac{\partial f(x_3; \mathbf{w})}{\partial w[2]} \end{bmatrix} = \begin{bmatrix} 0.5 e^{0.5w[1]+w[2]} & e^{0.5w[1]+w[2]} \\ 1.5 e^{1.5w[1]+w[2]} & e^{1.5w[1]+w[2]} \\ 3 e^{3w[1]+w[2]} & e^{3w[1]+w[2]} \end{bmatrix}$$

# Problems

---

- Similar problems to linear least squares

- Matrix inverse is expensive
- Outliers!

$$\Delta \mathbf{w} = (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \Delta \mathbf{y}$$

- Initial guess

- Very sensitive to choice of values
- Cannot guarantee global minimum
- Where do we get it from?
  - Expert knowledge
  - OR, try random values (multiple times!)

# Topics

---

- Brief overview of machine learning: Part 3
  - Training vs. testing
  - Parameter estimation
- Non-linear least squares
  - Estimate parameters for non-linear functions
- Gradient descent
  - Alternate solution methods
  - Not limited to SSE error functions

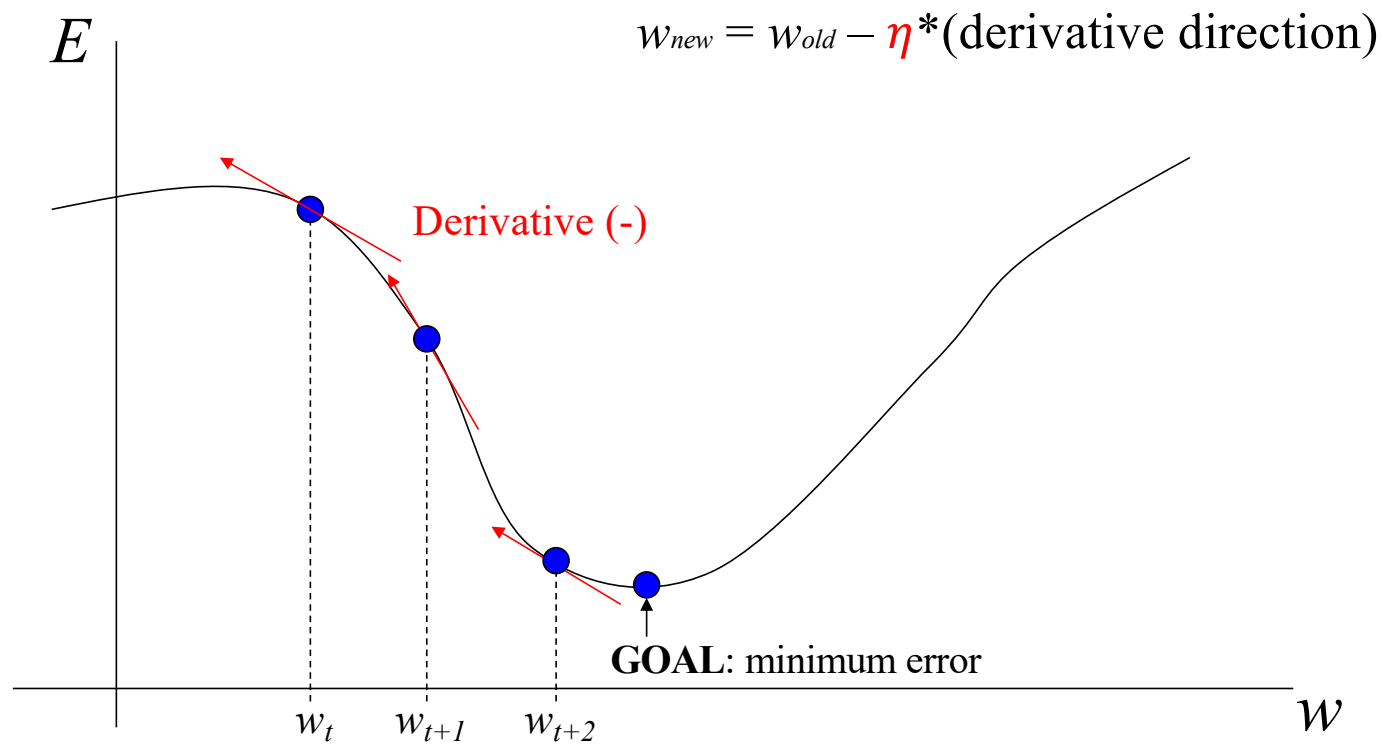
# Gradient Descent (GD)

---

- Simple (and faster?) method for optimization
  - $f(\mathbf{x}; \mathbf{w})$  linear or nonlinear w.r.t. parameters  $\mathbf{w}$
  - Sum or square losses or other losses
- A general method to find local minimum for a multi-parameter error function
  - Search through “parameter space” to find minimum error
- Main idea
  - Evaluate error function for current guess of parameters
  - Determine change in parameters that decrease error
    - From gradient of error function
  - Update parameters in this new change
  - Repeat process until converge (or hit max iterations)

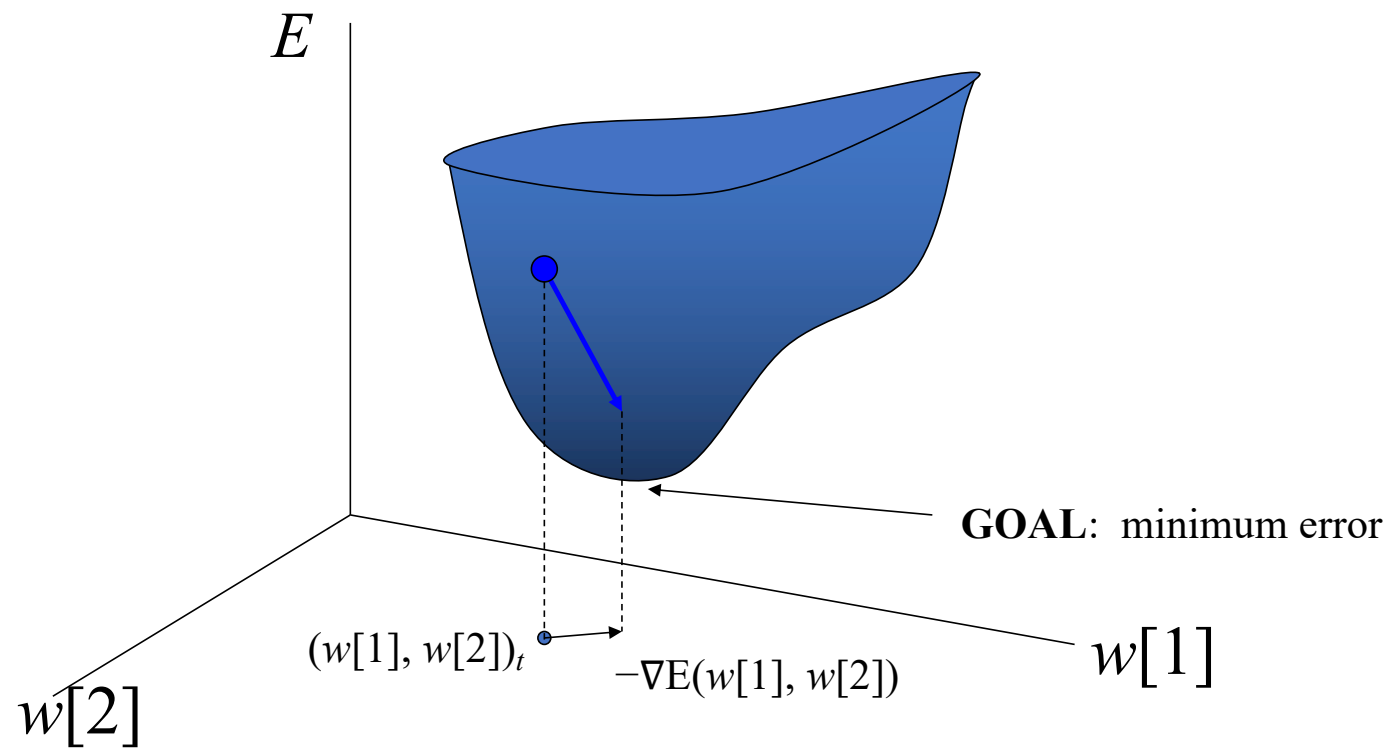
# Problem Visualization

---



# Problem Visualization

---





# Why does gradient descent work?

---

- Recall the Taylor expansion

- But now apply it to the loss function  $E(\mathbf{w})$  directly, not to  $f(\mathbf{x}_i; \mathbf{w})$

$$\begin{aligned} E(\mathbf{w}) &= E(\mathbf{w}') + \sum_{k=1}^K \frac{\partial E(\mathbf{w}')}{\partial w[k]} (w[k] - w'[k]) + \text{higher-order terms (H.O.T.)} \\ &= E(\mathbf{w}') + \nabla_{\mathbf{w}} E(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}') + \text{H.O.T.} = E(\mathbf{w}') + \nabla_{\mathbf{w}} E(\mathbf{w}')^T \Delta \mathbf{w} + \text{H.O.T.} \end{aligned}$$

- Approximate  $E(\mathbf{w})$  by the linear terms:  $E'(\Delta \mathbf{w}) = E(\mathbf{w}') + \nabla_{\mathbf{w}} E(\mathbf{w}')^T \Delta \mathbf{w}$
- Which direction, i.e.,  $\Delta \mathbf{w}$  with  $\|\Delta \mathbf{w}\|_2^2 = 1$ , can minimize  $E'(\Delta \mathbf{w})$ ?
  - Answer: **direction of**  $\Delta \mathbf{w} = -\frac{\nabla_{\mathbf{w}} E(\mathbf{w}')}{\|\nabla_{\mathbf{w}} E(\mathbf{w}')\|_2}$
  - Interpretation: move in the direction of  $-\nabla_{\mathbf{w}} E(\mathbf{w}')$  is the most efficient for minimization
  - Only work for **small  $\Delta \mathbf{w}$  (so a scale  $\eta$ )** due to the nature of linear Taylor expansion

# Algorithm

---

Initialize parameters:  $\hat{\mathbf{w}}^{(1)} = [w[1], w[2], \dots, w[K]]^T$

Error function:  $E(\mathbf{w})$

---

Compute gradients:  $\nabla E(\hat{\mathbf{w}}^{(t)}) = \begin{bmatrix} \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[1]} \\ \vdots \\ \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[K]} \end{bmatrix}$

Update parameters:  $\hat{\mathbf{w}}^{(t+1)} = \hat{\mathbf{w}}^{(t)} - \eta \cdot \nabla E(\hat{\mathbf{w}}^{(t)})$

[Repeat until converges]

Learning rate



# Compute Gradients: Least Squares ( $w$ is 1-dim)

---

Error function:  $E(\mathbf{w}) = \sum_i [f(\mathbf{x}_i; \mathbf{w}) - y_i]^2$

Error gradient:  $\nabla E(\mathbf{w})$

$$\begin{aligned} \frac{dE(\hat{\mathbf{w}}^{(t)})}{dw} &= \frac{\partial}{\partial w} \sum_i [f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i]^2 \\ &= 2 \sum_i (f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i) \frac{df(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)})}{dw} \end{aligned}$$

Recall: Gauss-Newton

$$2 \sum_i (f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i) \frac{df(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)})}{\partial w} = 2 \sum_i -\Delta y_i J_i(\hat{\mathbf{w}}^{(t)}) \quad \longrightarrow \quad \nabla E(\mathbf{w}) = -2 J^T(\mathbf{w}) \Delta \mathbf{y}(\mathbf{w})$$

# Compute Gradients: Least Squares ( $\mathbf{w}$ is $K$ -dim)

---

Error function:  $E(\mathbf{w}) = \sum_i [f(\mathbf{x}_i; \mathbf{w}) - y_i]^2$

Error gradient:  $\nabla E(\mathbf{w})$

$$\begin{aligned} \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[k]} &= \frac{\partial}{\partial w[k]} \sum_i [f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i]^2 \\ &= 2 \sum_i (f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i) \frac{\partial f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)})}{\partial w[k]} \end{aligned}$$

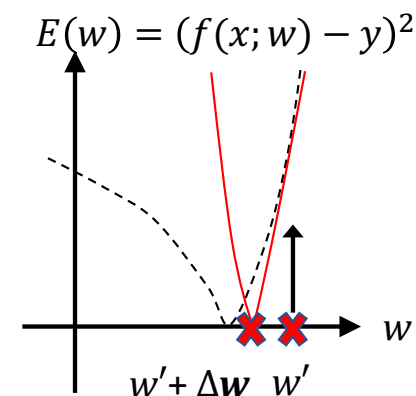
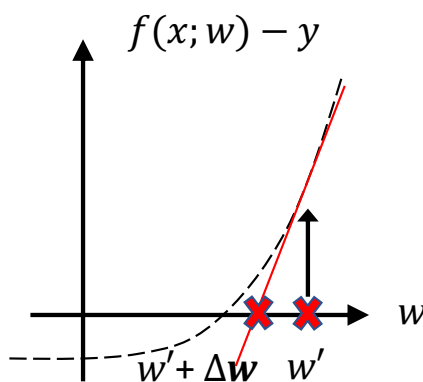
Recall: Gauss-Newton

$$2 \sum_i (f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i) \frac{\partial f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)})}{\partial w[k]} = 2 \sum_i -\Delta y_i J_{ik}(\hat{\mathbf{w}}^{(t)}) \quad \longrightarrow \quad \nabla E(\mathbf{w}) = -2 J^T(\mathbf{w}) \Delta \mathbf{y}(\mathbf{w})$$

# Gradient descent vs Gauss-Newton work

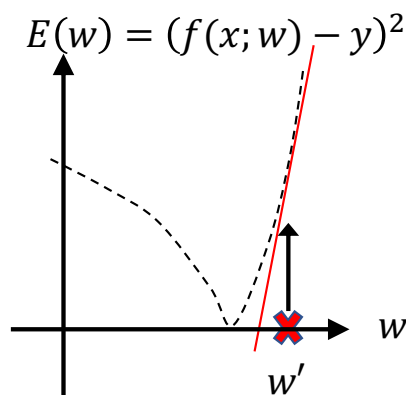
- Gauss-Newton

- For square loss  $\sum_i [f(\mathbf{x}_i; \mathbf{w}) - y_i]^2$
- Linearization on  $f(\mathbf{x}; \mathbf{w})$  at  $\mathbf{w}'$
- Solve the approximated least square with the closed-form solution to get  $\Delta \mathbf{w}$
- Update  $\mathbf{w}'$  by  $\mathbf{w}' + \Delta \mathbf{w}$



- Gradient descent

- For any kind of loss  $E(\mathbf{w})$
- Linearization in  $E(\mathbf{w})$  at  $\mathbf{w}'$
- Set  $\Delta \mathbf{w}$  to the direction of  $-\nabla_{\mathbf{w}} E(\mathbf{w}')$
- Update  $\mathbf{w}'$  by  $\mathbf{w}' + \eta \Delta \mathbf{w}$



# Gradient descent vs Gauss-Newton work

---

- Gauss-Newton can be seen as “quadratic approximation” of  $E(\mathbf{w})$

$$\begin{aligned} E'(\Delta \mathbf{w}) &= \sum_i [\nabla_{\mathbf{w}} E(\mathbf{w}')^T \Delta \mathbf{w} - \Delta y_i(\mathbf{w}')]^2 \\ &= \Delta \mathbf{w}^T \left( \sum_i \nabla_{\mathbf{w}} E(\mathbf{w}') \nabla_{\mathbf{w}} E(\mathbf{w}')^T \right) \Delta \mathbf{w} - 2 \sum_i \Delta y_i(\mathbf{w}')^T \nabla_{\mathbf{w}} E(\mathbf{w}')^T \Delta \mathbf{w} + \sum_i \Delta y_i(\mathbf{w}')^T \Delta y_i(\mathbf{w}') \end{aligned}$$

- As we use the quadratic term, the approximation is good for a larger  $\Delta \mathbf{w}$ 
  - Thus we can directly update  $\mathbf{w}'$  by  $\mathbf{w}' + \Delta \mathbf{w}$  (without a learning rate  $\eta$ )
  - In practice, you can also do  $\mathbf{w}'$  by  $\mathbf{w}' + \eta \Delta \mathbf{w}$

# Direct 2-order Taylor expansion approximation

---

- Taylor series expansion ( $\Delta \mathbf{w} = \mathbf{w} - \mathbf{w}'$ )

$$E(\Delta \mathbf{w}) = E(\mathbf{w}') + \nabla_{\mathbf{w}} E(\mathbf{w}')^T \Delta \mathbf{w} + \frac{1}{2} \Delta \mathbf{w}^T H(\mathbf{w}') \Delta \mathbf{w} + \text{H.O.T.}$$

- Hessian matrix  $H(\mathbf{w}')$

$$H(\mathbf{w}') = \begin{bmatrix} \frac{\partial^2 E(\hat{\mathbf{w}}_t)}{\partial w[1] \partial w[1]} & \cdots & \frac{\partial^2 E(\hat{\mathbf{w}}_t)}{\partial w[1] \partial w[K]} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 E(\hat{\mathbf{w}}_t)}{\partial w[K] \partial w[1]} & \cdots & \frac{\partial^2 E(\hat{\mathbf{w}}_t)}{\partial w[K] \partial w[K]} \end{bmatrix}$$

- Newton's method

- For any kind of loss  $E(\mathbf{w})$
- Quadratic approximation on  $E(\mathbf{w})$  at  $\mathbf{w}'$
- Solve the quadratic approximation to get  $\Delta \mathbf{w} = -H^{-1}(\mathbf{w}') \nabla_{\mathbf{w}} E(\mathbf{w}')$
- Update  $\mathbf{w}'$  by  $\mathbf{w}' + \Delta \mathbf{w}$

- Pros: faster convergence; Cons:  $H(\mathbf{w}')$  and  $H^{-1}(\mathbf{w}')$  are time consuming

# Benefits, Issues, Helpers

---

- Benefits of gradient descent
  - No expensive matrix operations (i.e., inversion)
  - Non-SSE error functions (as long as gradient is well formed)
  - Easy to parallelize (supercomputers!)
- Issues
  - Can be misled by local minima
    - Get stuck at smaller valleys
  - Can get stuck on a flat plain
    - Flat region in state-space function
  - Can overshoot back and forth
    - Oscillate from side to side
    - Learning rate too large
- Helpers
  - Use random-restart to handle some problems
    - Do process multiple times with random initial parameters
  - Methods exist for selecting appropriate (best) learning rate at each iteration



# Summary

---

- Gauss-Newton non-linear least squares
  - Iterative approximation
  - Linearize model around current parameter guess
  - Apply Linear Least Squares methods to linearized model
- Gradient descent
  - “Valley descending” (opposed to “hill climbing”)
  - Make the best update in parameters
    - In direction of negative gradient of error function