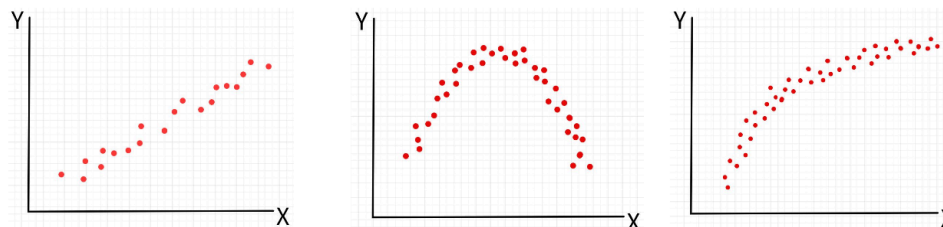


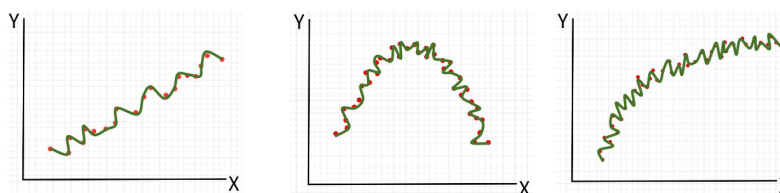
## Linear Regression

## 1. Motivation

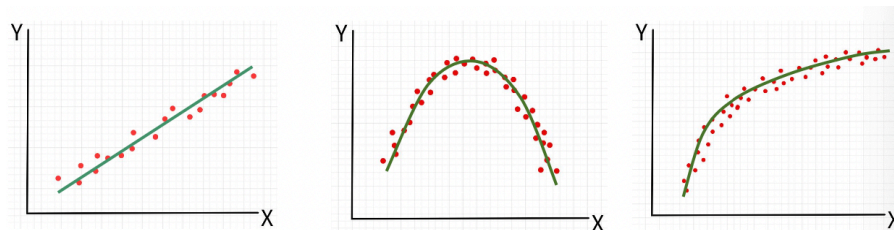
- a. Given  $n$  samples / data points  $(y_i, x_i)$ .  $Y$  is a continuous variable (as opposed to classification)



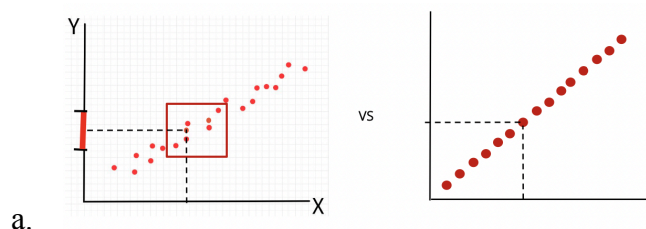
- b. Understand/explain how  $y$  varies as a function of  $x$  (i.e. find a function  $y = h(x)$  that best fits our data)
- d. Should  $h$  be the curve that goes through the most samples? I.e. do we want  $h(x_i) = y_i$  for the maximum number of  $i$ ?

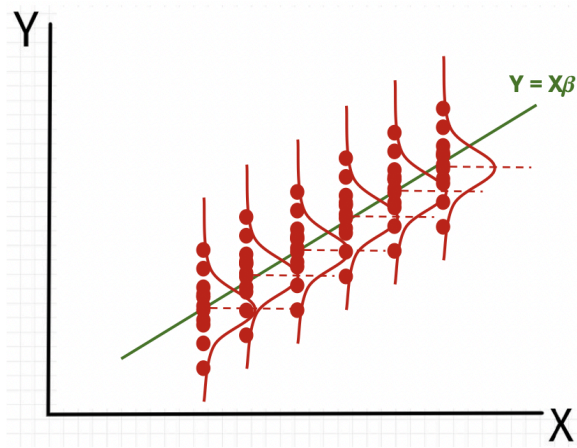


- i.  $h$  may be too complex
- iii. overfitting - may not perform well on unseen data
- e. The following curves seem the most intuitive “best fit” to our samples. How can we define this best fit mathematically? Is it just about finding the right distance function?

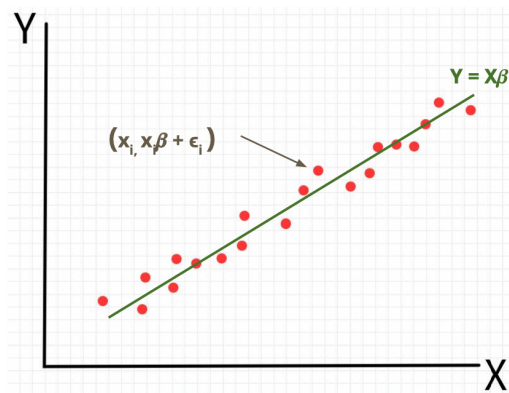


## 2. Assumptions

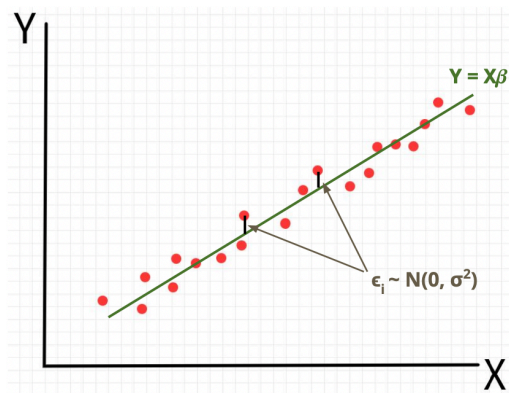




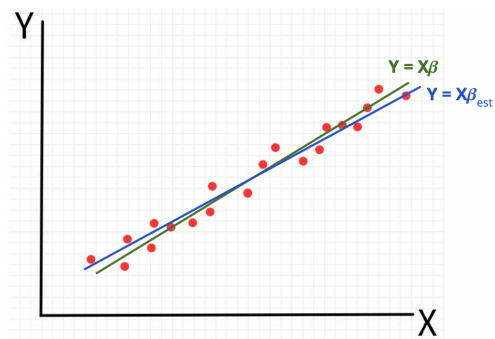
b.



c.



d.



e.

f. Our data was generated by a linear function plus some noise:

i.  $y = h_x(B) + e$

where  $h$  is linear in a parameter  $B$

where  $e$  are independent  $N(0, \sigma^2)$  distribution

### 3. Cost Function

a. Given our data:  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

b. Suppose we are given a curve  $y = h(x)$ , how can we evaluate whether it is a good fit to our data?

c. Compare  $h(x_i)$  to  $y_i$  for all  $i$

d. Goal: For a given distance function  $d$ , find  $h$  where  $L$  is smallest

$$L(h) = \sum_i d(h(x_i), y_i)$$

e.

### 4. Assumptions

a. The relation between  $x$  (independent variable) and  $y$  (dependent variable) is linear in a parameter  $B$

b.  $e$  are independent, identically distributed random variables following a  $N(0, \sigma^2)$  distribution

### 5. Goal

a. Given these assumptions, let's try to minimize the cost function defined earlier

b. What parameter(s) are we trying to learn / estimate?

i.  $A: B$

### 6. Least Squares

$$\beta_{LS} = \arg \min_{\beta} \sum_i d(h_{\beta}(x_i), y_i)$$

$(x_i, y_i)$  are from our dataset

a.

$$\beta_{LS} = \arg \min_{\beta} \sum_i d(h_{\beta}(x_i), y_i)$$

$$= \arg \min_{\beta} \|\vec{y} - h_{\beta}(X)\|_2^2$$

$$= \arg \min_{\beta} \|y - X\beta\|_2^2$$

b.

$$\begin{aligned}
\frac{\partial}{\partial \beta} (y - X\beta)^T (y - X\beta) &= 0 \\
\frac{\partial}{\partial \beta} (y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta) &= 0 \\
\frac{\partial}{\partial \beta} (y^T y - 2\beta^T X^T y + \beta^T X^T X\beta) &= 0 \\
-2X^T y + X^T X\beta &= 0 \\
X^T X\beta &= X^T y \\
\boxed{\beta_{LS} = (X^T X)^{-1} X^T y}
\end{aligned}$$

c.

## 7. Assumptions

- a. Our data was generated by a linear function plus some noise:

$$\vec{y} = h_X(\beta) + \vec{\epsilon}$$

- b. Where h is linear in a parameter B.  
c. Which functions below are linear in B?
- $h(B) = B_1 * x$
  - $h(B) = B_0 + B_1 x$
  - $h(B) = B_0 + B_1 x + B_2 * x^2$
  - $h(B) = B_1 \log(x) + B_2 * x^2$
  - $h(B) = B_0 + B_1 x + B_1^2 * x$

## 8. Maximum Likelihood

- a. Another way to define this problem is in terms of probability  
b. Define  $P(Y|h)$  as the probability of observing Y given that it was sampled from h  
c. Goal: Find h that maximizes the probability of having observed our data  
d. Maximize  $L(h) = P(Y|h)$   
e. Since  $e \sim N(0, \sigma^2)$  and  $Y = XB + e$  then  $Y \sim N(XB, \sigma^2)$

$$\begin{aligned}
\beta_{MLE} &= \arg \max_{\beta} \frac{1}{\sqrt{(2\pi)^n \sigma^n}} \exp\left(-\frac{\|y - X\beta\|_2^2}{2\sigma^2}\right) \\
&= \arg \max_{\beta} \exp\left(-\frac{\|y - X\beta\|_2^2}{2\sigma^2}\right) \\
&= \arg \max_{\beta} -\|y - X\beta\|_2^2 \\
&= \arg \min_{\beta} \|y - X\beta\|_2^2 \\
&= \beta_{LS} = (X^T X)^{-1} X^T y
\end{aligned}$$

f.

9. An Unbiased Estimator

- a. BLS is an unbiased estimator of the true  $\beta$ . That is  $E[\beta_{LS}] = \beta$

$$\begin{aligned} E[\beta_{LS}] &= E[(X^T X)^{-1} X^T y] \\ &= (X^T X)^{-1} X^T E[y] \\ &= (X^T X)^{-1} X^T E[X\beta + \epsilon] \\ &= (X^T X)^{-1} X^T X\beta + E[\epsilon] \\ &= \beta \end{aligned}$$

- b.