

## Project 3: Linear Regression

Goal

Beyond just "looking" at data in a scatter plot, we went to quantify relationships between 2+ variables and the strength of association.

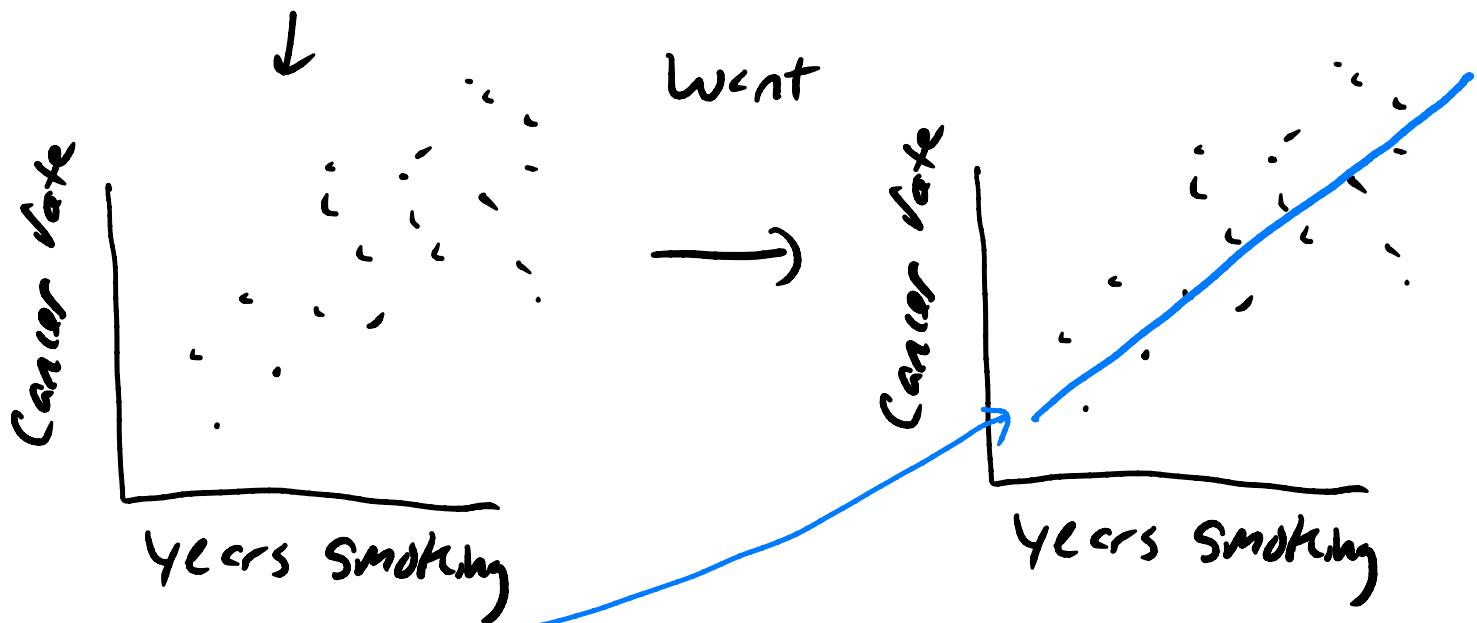
is Cancer associated with

- smoking?
- drinking?
- old age?
- ⋮

- weak?
- Strong?
- Unrelated?

Raw Data

=



regression: curve fit to data that tolerates error in measurement and noise

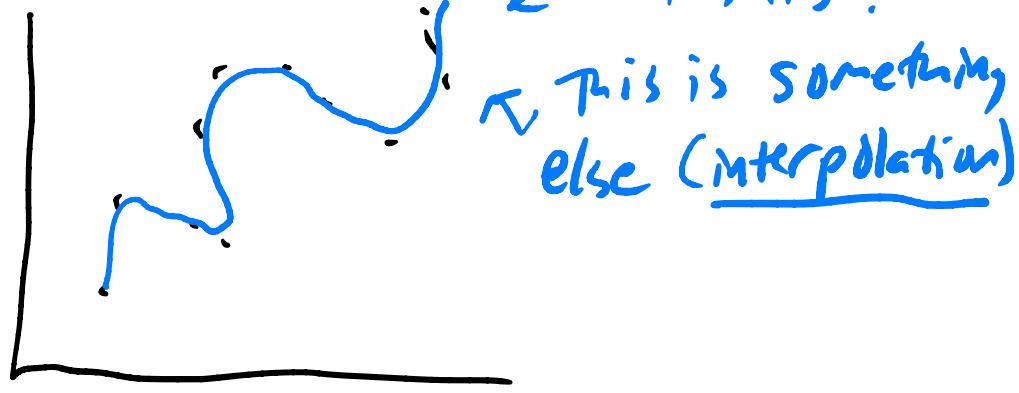
e.g. if Cancer data comes from Self-report Survey, people

- make mistakes in filling out form
- misremember
- guess (b/c they forgot)

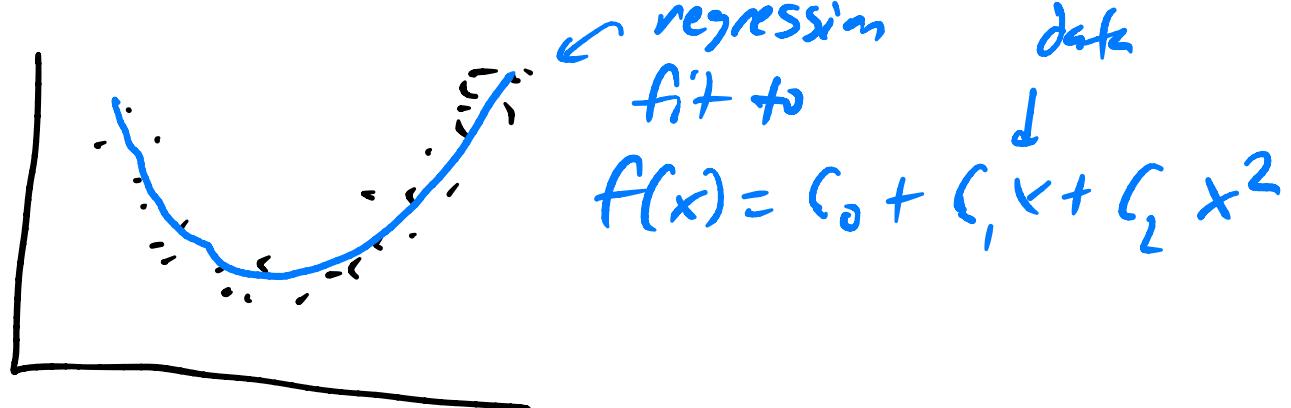
...

\* Try to get as close as possible to all samples with regression Curve.

- Does not need to pass thru all data Samples



- Does not need to be a line



## Fitting a Straight line to data

Let's start simple to gain intuition:

- Regression between 2 variables:

- 1 independent variable ( $x$ )

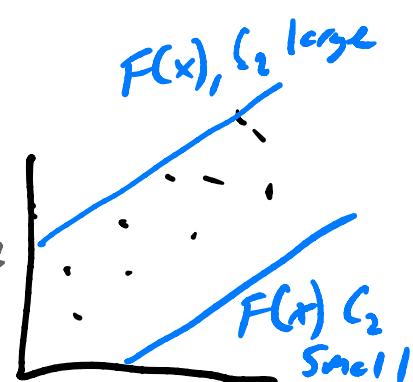
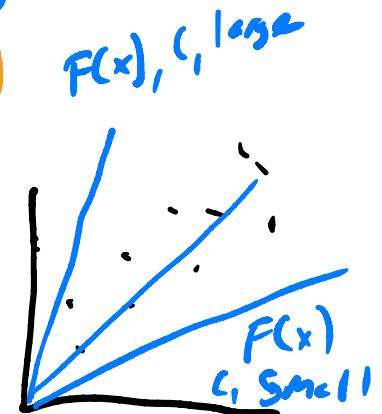
- 1 dependent variable ( $y$ )  $f(x), c_1 \text{ large}$

- Linear Fit function:

$$y_i = F(x_i) = c_0 + c_1 x_i$$

↑  
Unknown  
intercept  
of regression  
line.

↑  
Unknown Slope  
of regression line

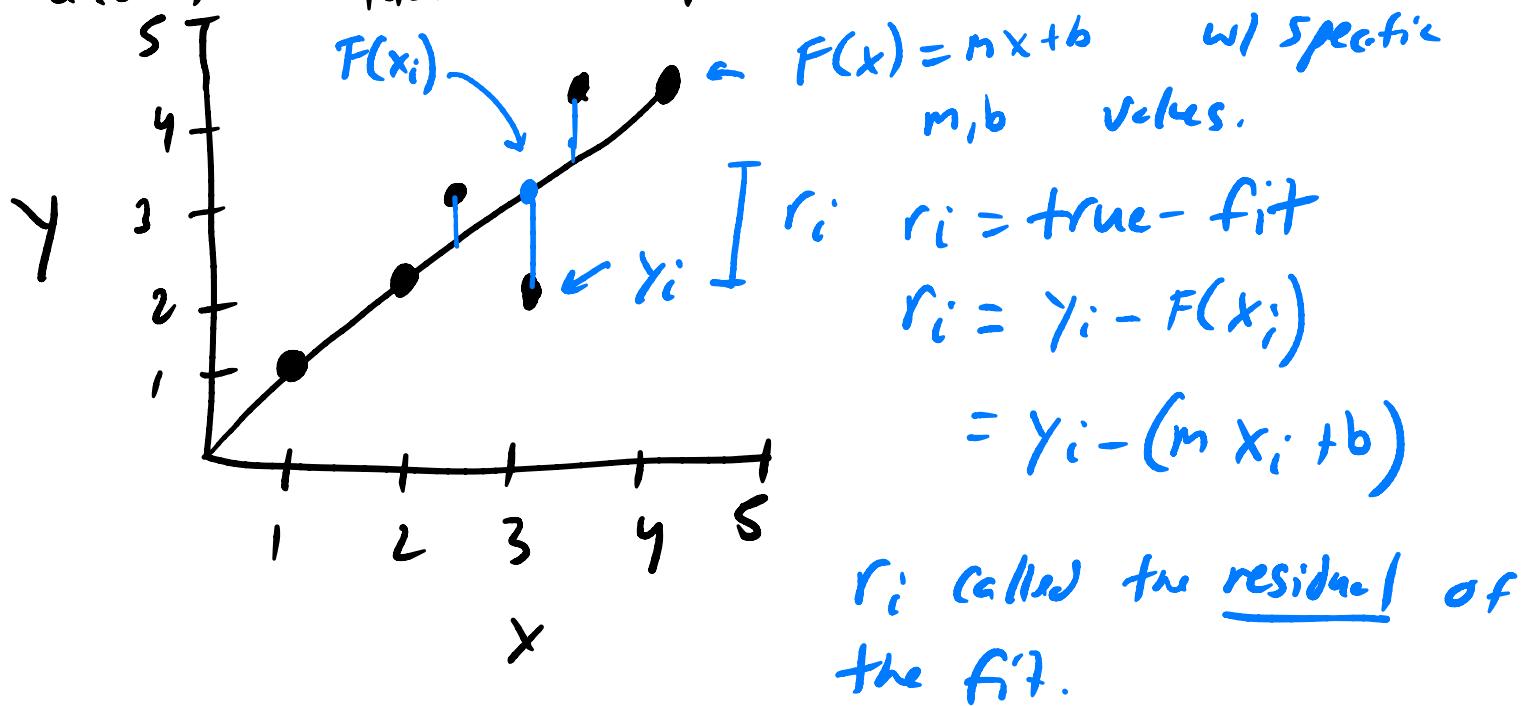


Let's rename  $c_1 = m$ ,  $c_0 = b$

$$\Rightarrow y_i = F(x_i) = m x_i + b$$

**Goal** Find best  $m$  &  $b$  values to fit data.

. Error in regression fit to data  $\rightarrow$  vertical distance between line fit at each  $x_i$ :  $F(x_i)$  and the actual data  $y_i$ :



problem: Can't just sum  $r_i$  to get total error in fit — why??

Some  $r_i \geq 0$ , some  $r_i \leq 0$

$\Rightarrow$  error would cancel!



Sum = 0, but there is error!

possible workarounds:

Absolute distance

$$\sum_i |r_i|$$

Squared distance

$$\sum_i r_i^2$$

- Squared dist usually selected b/c Math easier to compute best m, b.
- With Squared dist, linear regression goal is to find coefficients m, b that minimize the average error over all data pts
  - ↳ Called mean sum of squares error (MSSE)

$$MSSE = \rho = \frac{1}{N} \sum_{i=1}^N r_i^2 = \frac{1}{N} \sum_{i=1}^N (y_i - F(x_i))^2$$

- Minimize MSSE equation with respect to m, b — take partial derivatives set to 0, solve!
- Approach = least squares

$$m = \frac{(\sum x_i)(\sum y_i) - N \sum x_i y_i}{(\sum x_i)^2 - N \sum x_i^2}$$

$$F(x_i) = m x_i + b$$

$$b = \frac{(\sum x_i)(\sum x_i y_i) - (\sum x_i^2)(\sum y_i)}{(\sum x_i)^2 - N \sum x_i^2}$$

$$\Rightarrow \begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & N \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix} = \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

$x_{\text{independent}}$   $\rightarrow A$        $\vec{c}_{\text{unknown coefficients}}$        $\vec{y}_{\text{dependent Variable}}$

$$\Rightarrow A \vec{c} = \vec{y}$$

## Regression Work flow

• ok so now what? — we have  $m, b$ .

\* This gives us a **predictive model** of the data:

use  $A$  & fitted  $\vec{c}$  to compute  $\vec{y}_{\text{pred}}$ ,  
the predicted  $y$  values:  
 $A \vec{c} = \vec{y}_{\text{pred}}$



## Options

1) plug in Same independent Variable  $x_i$ : values used to fit regression Coefficients.

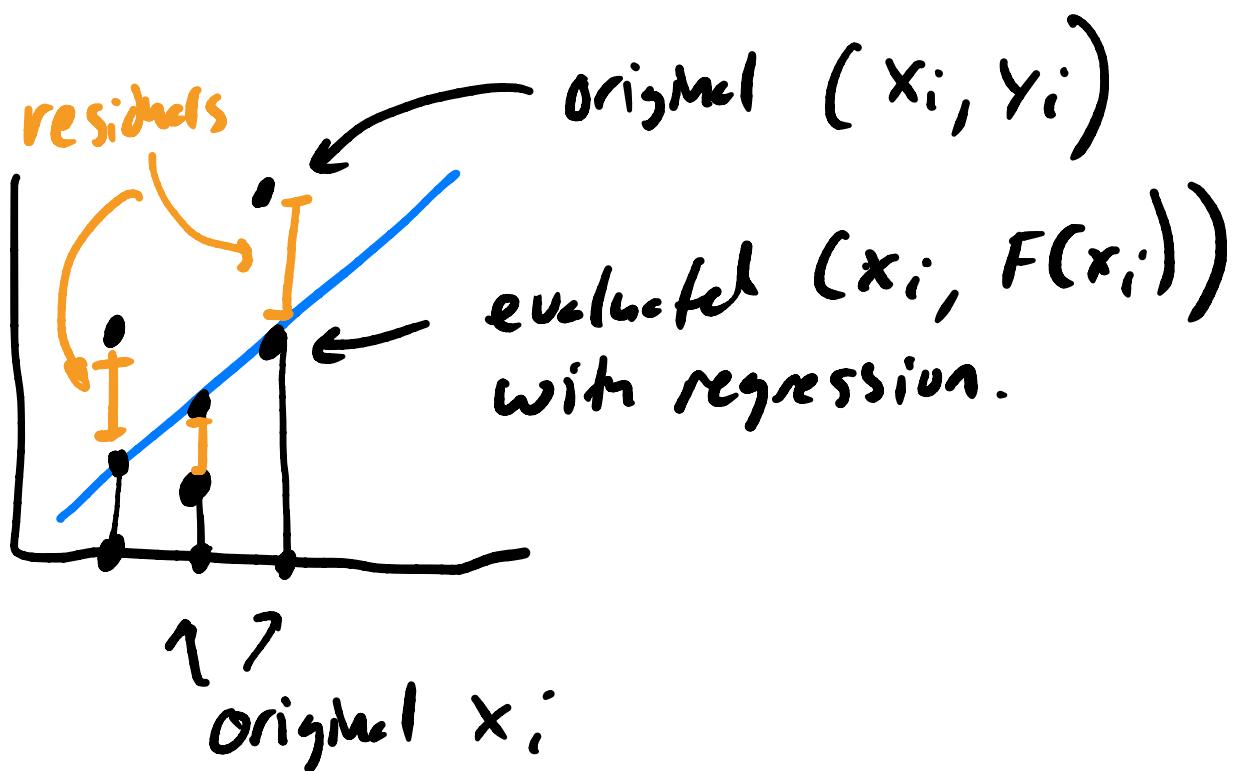
( $m, b$ )

$$A \vec{c} = \vec{y}$$

$A$  is same as used to solve for  $\vec{c}$ .

- can use to assess fit error

(e.g. MSSE)



2) Plug in  $\vec{c}$  and use  $x_i$  values  
—  $A_{\text{new}}$  — then those used to solve for  $\vec{c}$ .

⇒ use  $\vec{c}$  to get predictions —  $\hat{y}_i$

$$A_{\text{new}} \vec{c} = \hat{y}_{\text{pred}}$$

\* Tests regression generalization  
to new data.

