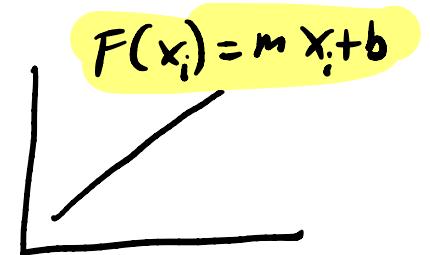


$$m = \frac{(\sum x_i)(\sum y_i) - N \sum x_i y_i}{(\sum x_i)^2 - N \sum x_i^2}$$



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

Now that we know m & b that minimize the least squares errors — let's write the linear regression equation in matrix form:

Scalar Equation

$$F(x_i) = b + m x_i = \underbrace{(1, x_i)}_{\text{dot product between data & coefficients}} \cdot \underbrace{(b, m)}_{\text{dot product}}$$

data matrix A :

$$\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \quad \vec{x}$$

independent Variable

Slope + intercept parameters:

$$\begin{bmatrix} b \\ m \end{bmatrix} \quad \vec{c}$$

dependent var vector

$$\vec{y}$$

For intercept
→ like homogeneous coord!!

(2, 1)

(N, 1)

↓

↓

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

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

To set us up for generality later, let's rename

intercept $b \rightarrow c_0$

slope $m \rightarrow c_1$

Our \vec{c} becomes : $\vec{c} = \begin{bmatrix} c_0 \\ c_1 \end{bmatrix}$

and the linear regression model equation becomes:

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

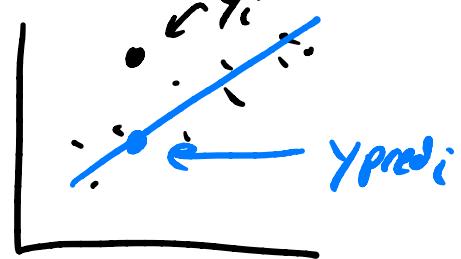
Regression Work flow

- 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}_{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.

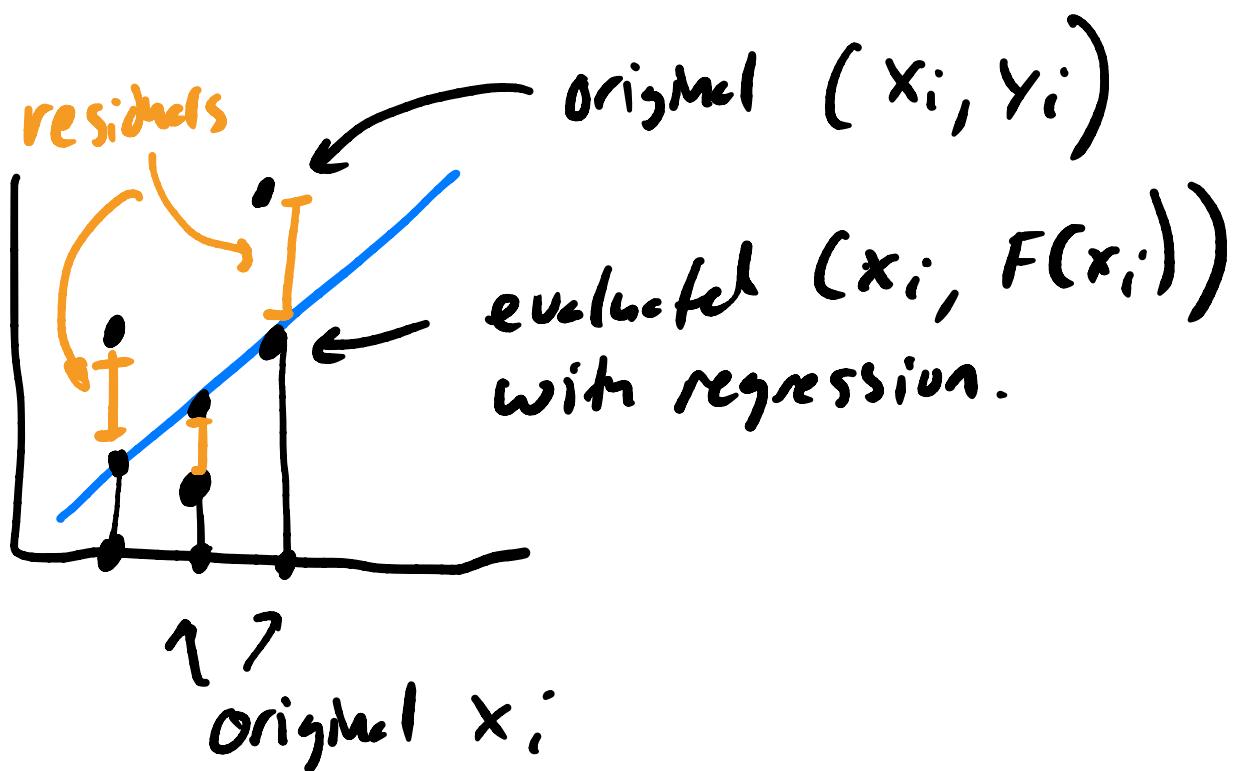
(m, b)

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

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

- can use to assess fit error

(e.g. MSSE)

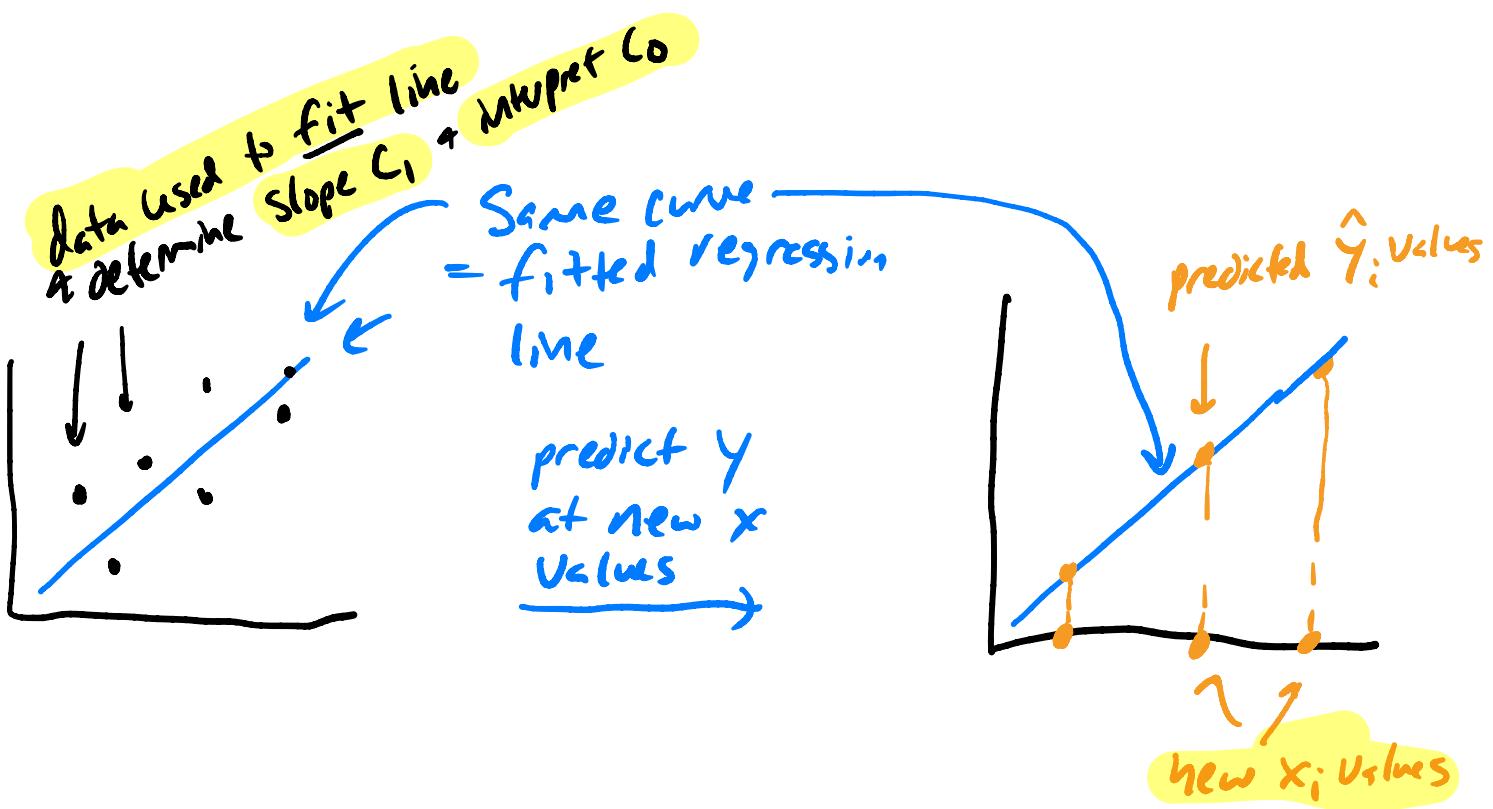


2) Plug in \vec{c} and use x_i values
— A_{new} — then those used to solve for \vec{c} .

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

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

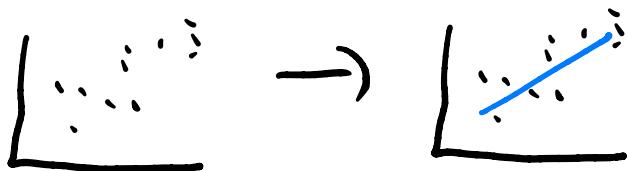
* Tests regression generalization
to new data.



Think of regression as a 2-step workflow:

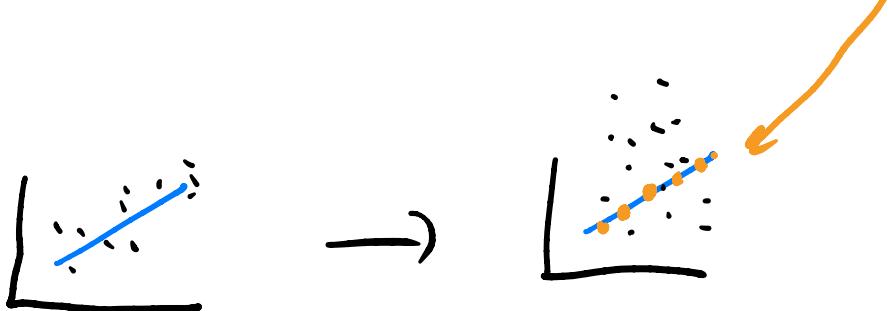
1) **Fit** regression model to data \Rightarrow get **model coefficients**

* solve $A\vec{c} = \vec{y}$ for \vec{c} . (e.g. m, b)
ie $[m, b]$



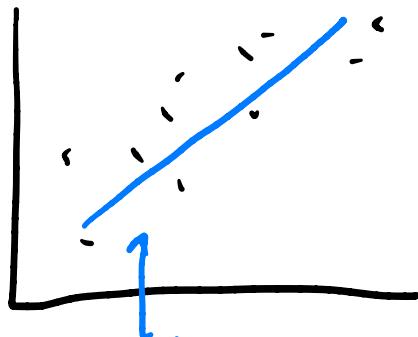
2) use fitted coefficients to **predict** Y value
of new data [not in original dataset used to
fit coefficients] — how well does regression
Generalize to new data?

* Solve $A_{\text{new}} \vec{c} = \vec{y}_{\text{pred}}$ for \vec{y}_{pred}

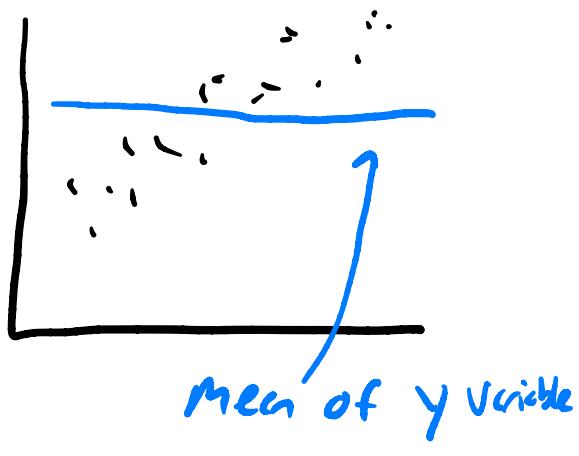


② Quality of fit

- Is all our work worth it? Is the regression a better explanation of the data than the simple mean?



vs.



Easy to Compute!

- Let \hat{y}_i be the predicted y-value by the regression model:

$$\hat{y}_i = c_0 + c x_i$$

- Let \bar{y} be the \bar{y} data mean (of the dependent variable)

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

regression
 model
 predictions
 ↓
 mean
 ↓

How much of an improvement are \hat{y}_i vs. \bar{y} ??

We define :

$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y}_.)^2}$$

$\hat{y} = y$ value on
 fitted line

] error of **model**
 our mean

] error of **data**
 our mean

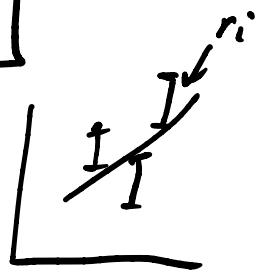
i.e how much is our linear model buying us?

Can be written in terms of **Sum-of-Squares error (SSE)**

$$R^2 = 1 - \frac{\|\vec{r}\|_2^2}{\sum (y_i - \bar{y})^2}$$

SSE = Sum of Squared error

where $\vec{r} = (\underbrace{r_1, r_2, r_3, \dots, r_N}_{\text{residual errors}})$



$\|\vec{r}\|_2$ = Euclidean distance

$$\text{Squared} \quad = \sqrt{r_1^2 + r_2^2 + r_3^2 + \dots + r_N^2}$$

$$\|\vec{r}\|_2^2 = r_1^2 + r_2^2 + r_3^2 + \dots + r_N^2$$

\uparrow Sum of Squared errors (SSE)

- like MSSE - but total

(not mean) error

- If $\|r\|_2^2 \approx 0$, $R^2 = 1$ [line passes thru all pts]

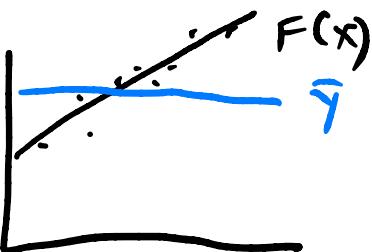


 $\underbrace{\|r\|_2^2}_{\text{error}}$

- If $\|r\|_2^2 \ll \sum (y_i - \bar{y})^2$

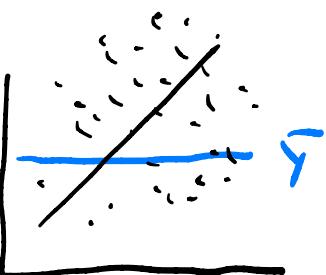
$$R^2 \approx 1.$$

Then big improvement over mean



- If $\|r\|_2^2 \approx \sum (y_i - \bar{y})^2$

$$R^2 \approx 0.$$

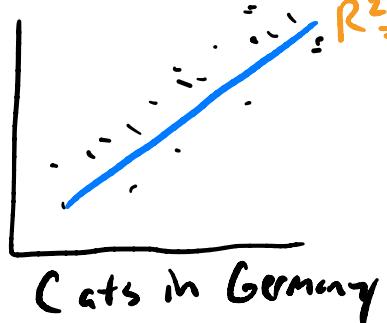


Then fit no better than mean

Important

Just b/c $R^2 \approx 1$, does not mean independent Variable X caused dependent Variable Y

Cancer rate in USA



$$R^2 = 0.95$$

Just b/c the Cat population in Germany is correlated with the Cancer rate in USA does not mean the Cats caused Cancer far away!
 \Rightarrow Can be a coincidence.