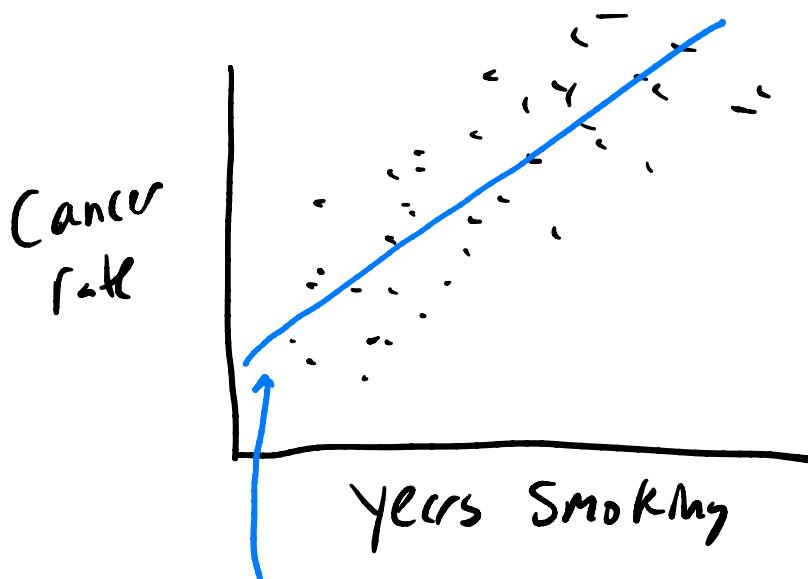


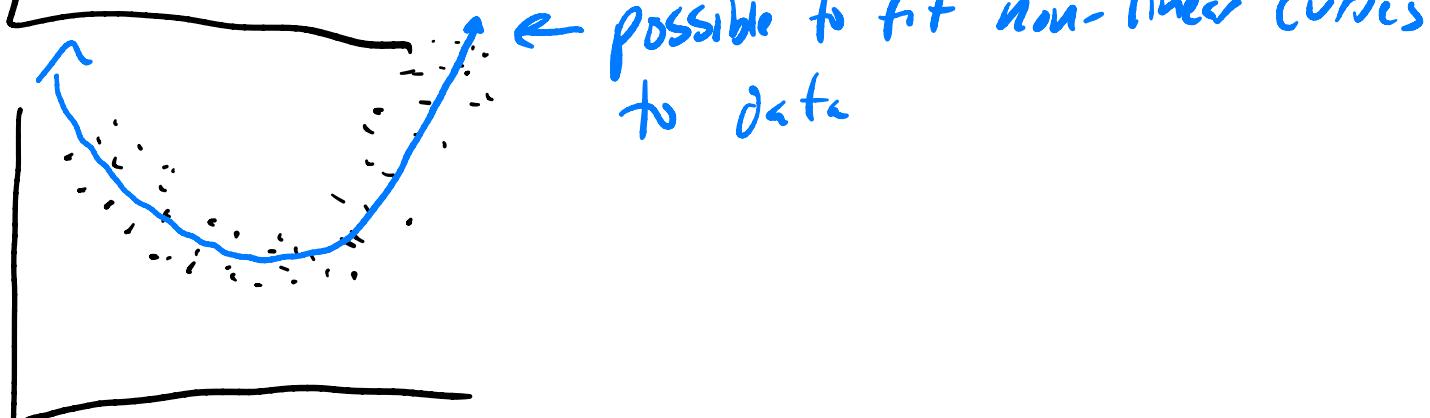
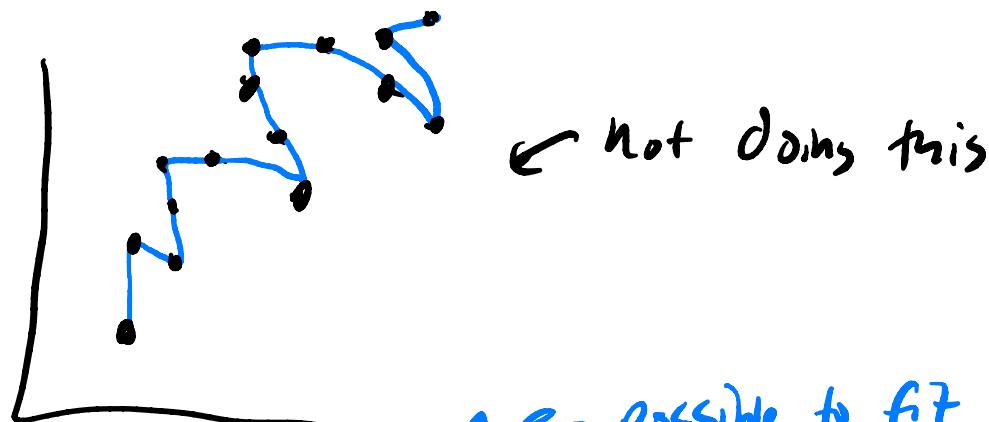
Lecture 11: Linear Regression



linear regression :
quantify the strength
of association between
any 2 vars.

regression Curve: Trying to get as close as possible to all the data samples.

tolerates / good with noise, measurement error



• ← genetic protection from cancer

Fitting a Straight line to 2 Vars

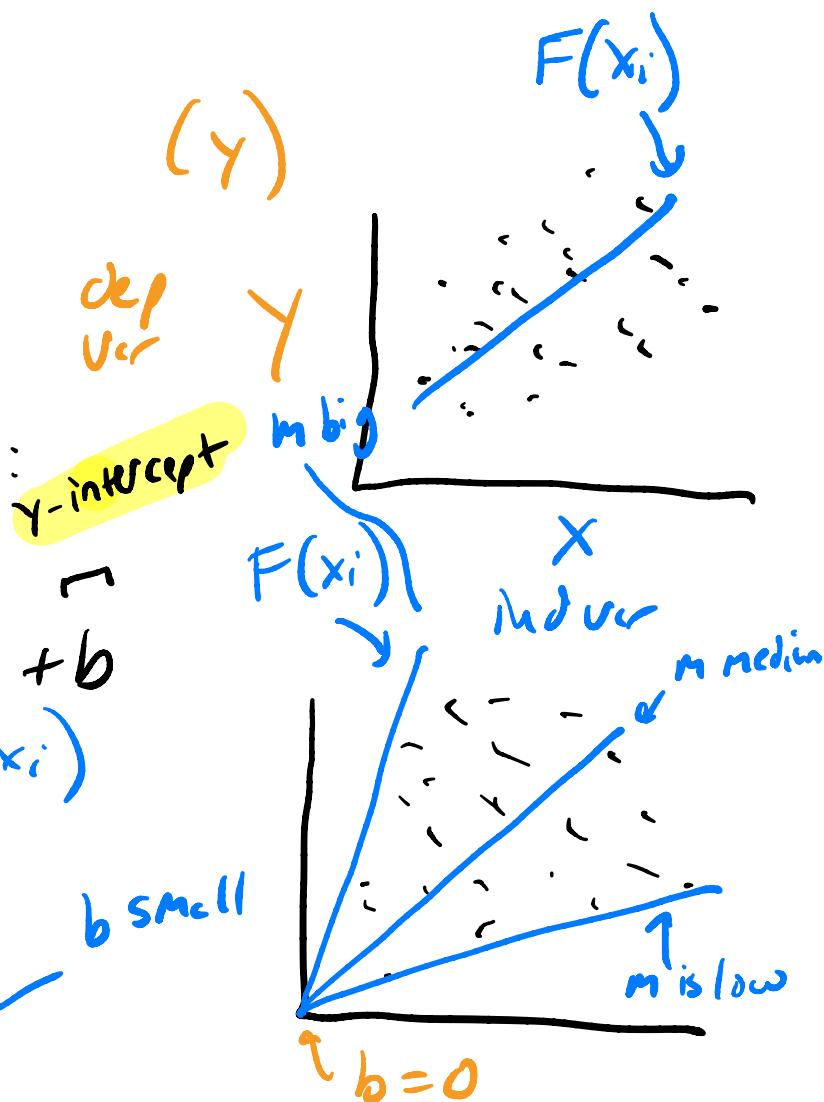
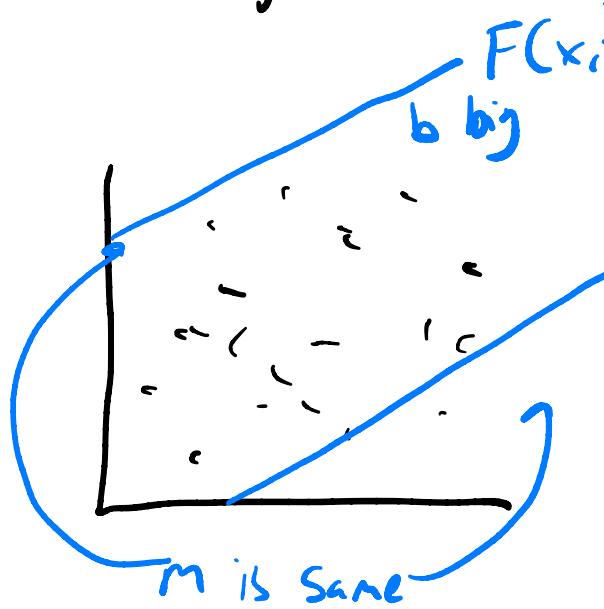
Linear regression w/ 2 Vars: Simple linear regression

1) Independent Variable (x)

2) Dependent Var

Linear regression equation:

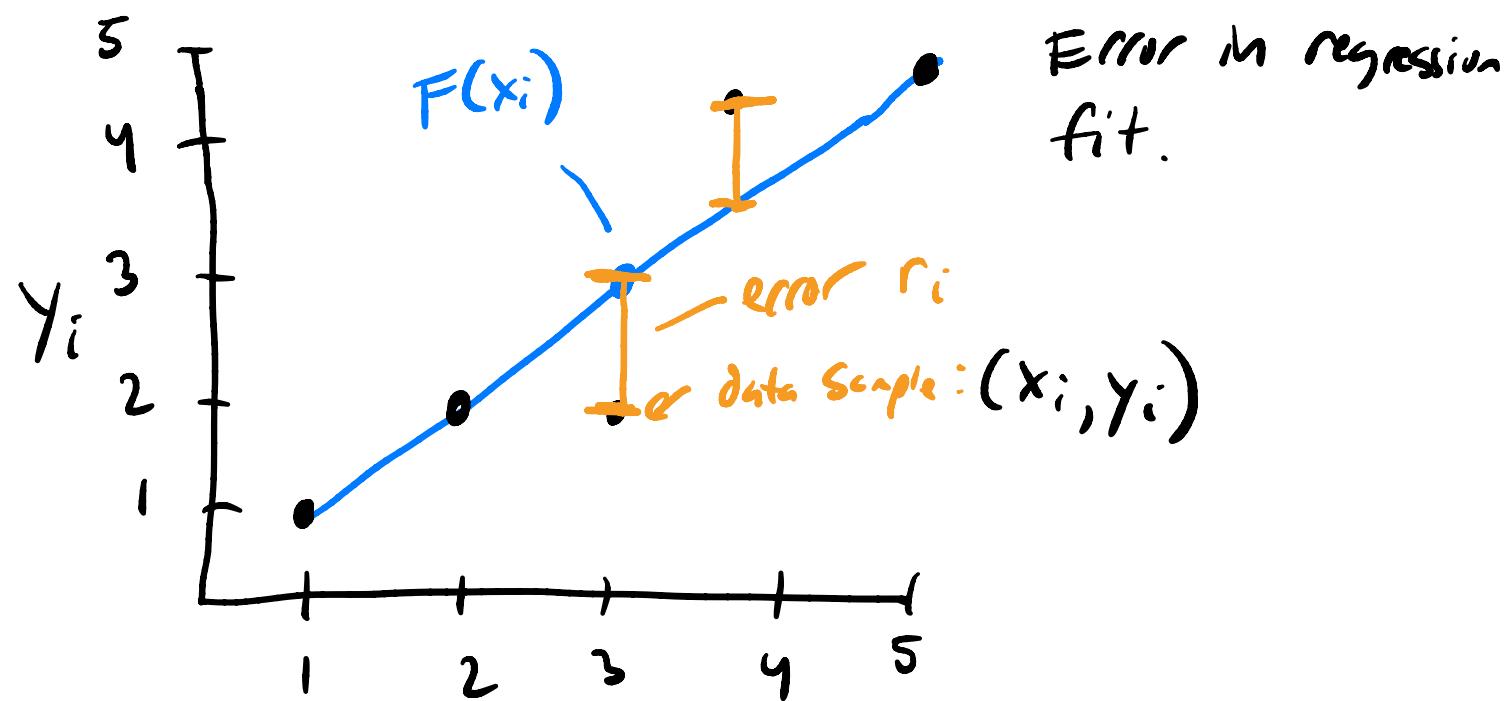
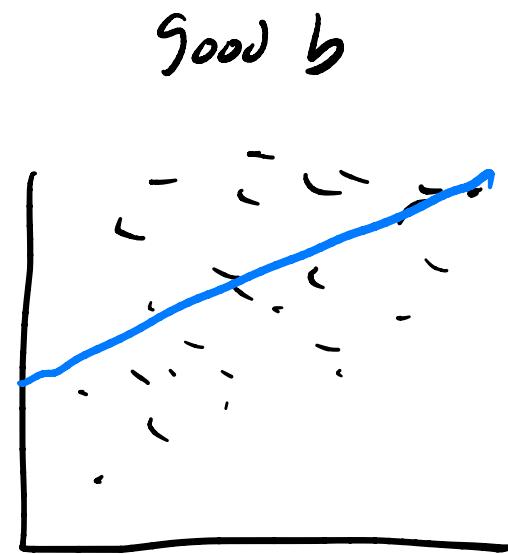
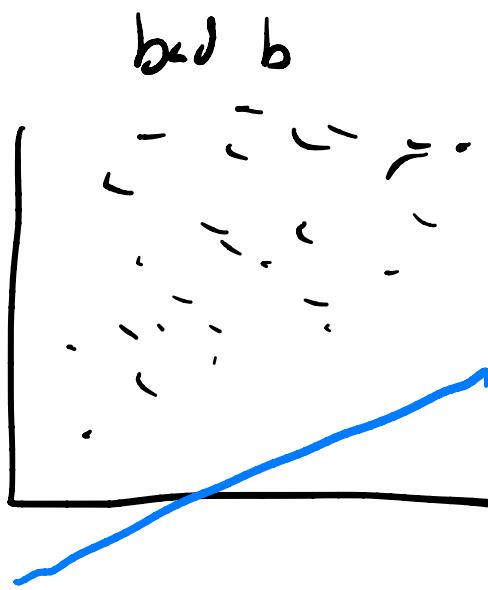
$$Y_i = F(x_i) = m x_i + b$$



Goal of regression

Find "best" values for m and b to fit data

Example: Fix m . & Find best b



$$r_i = \text{true value} - \text{fit value} = Y_i - F(x_i)$$

$$= Y_i - [mx_i + b] = \text{error}$$

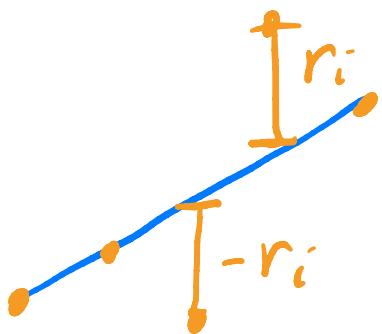
r_i : residual error for one sample i

W_{ent}: Overall error across all data samples.

Total error \rightarrow Sum up all r_i values.

Average them — error on average
across all samples.

problem:



$$\text{total error} = r_i - r_i = 0$$

Fixes:

1) Take absolute values

$$\text{total error} = \sum_{i=1}^N |r_i|$$

More common

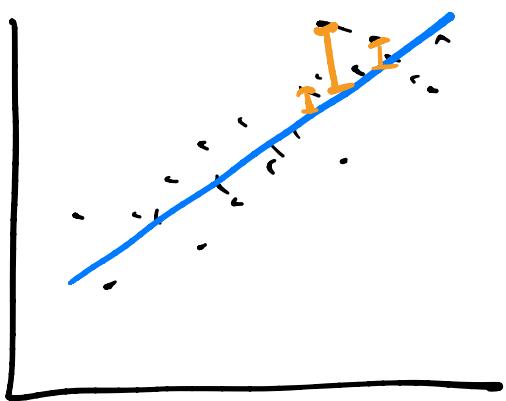
2) Squared error:

$$\text{error} = \sum_{i=1}^N r_i^2$$

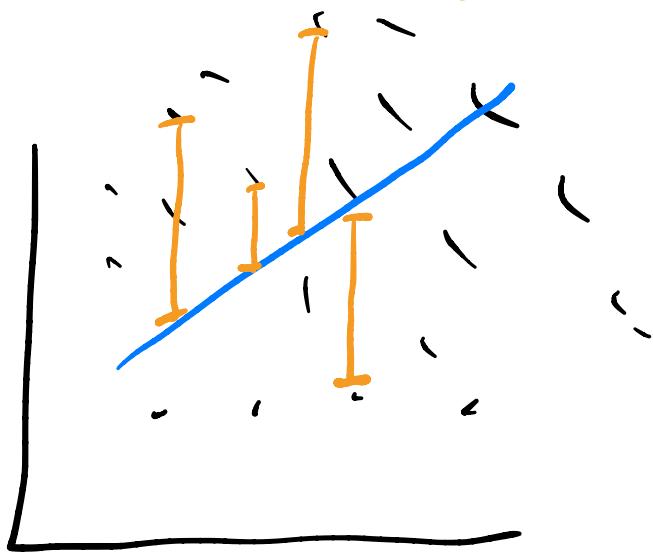
total average error = Mean Sum of Squared error
MSSE.

$$MSSE = \frac{1}{N} \sum_{i=1}^N r_i^2 = \frac{1}{N} \sum_{i=1}^N [y_i - (mx_i + b)]^2$$

↑
"Least Squares"
of error r_i MSSE is Small



MSSE is big

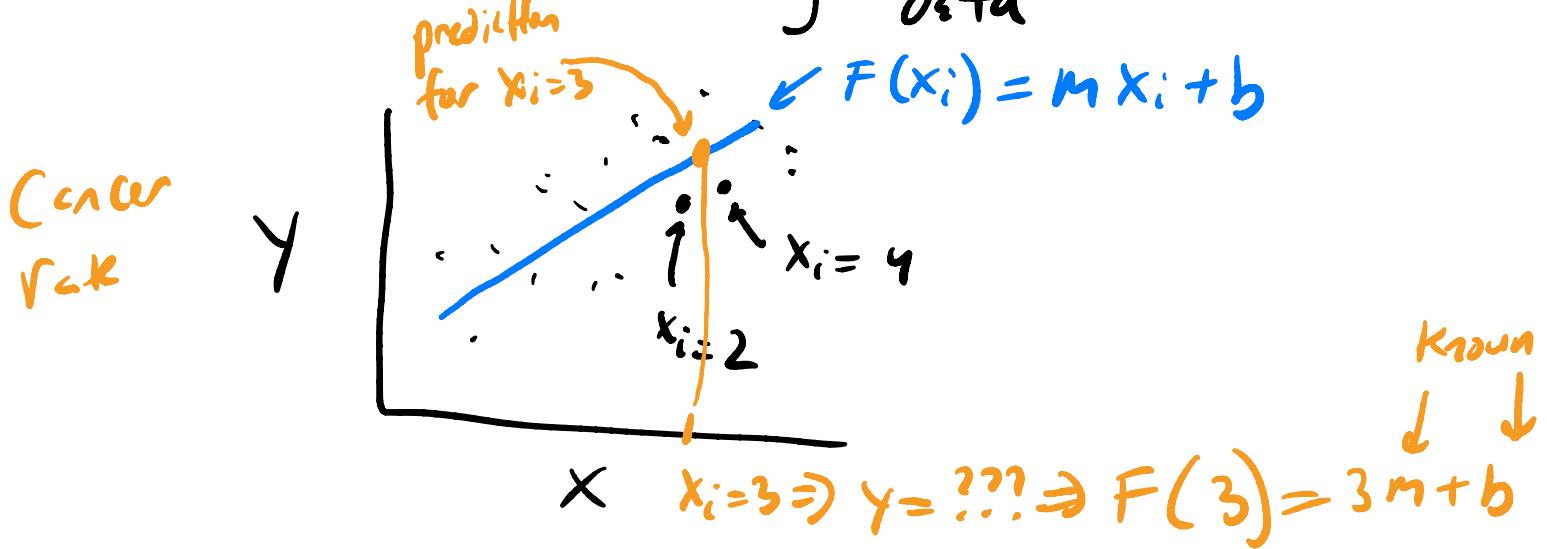


- 1) Set up data appropriately
- 2) Scipy function to give us $m+b$.
"Least Squares" function
- 3) Given m, b Values, we can calculate MSSE,
plot regression line

Regression Walk flow

Have m, b — now what?

$$y = F(x_i) = mx_i + b \quad] \quad \text{predictive model of data}$$



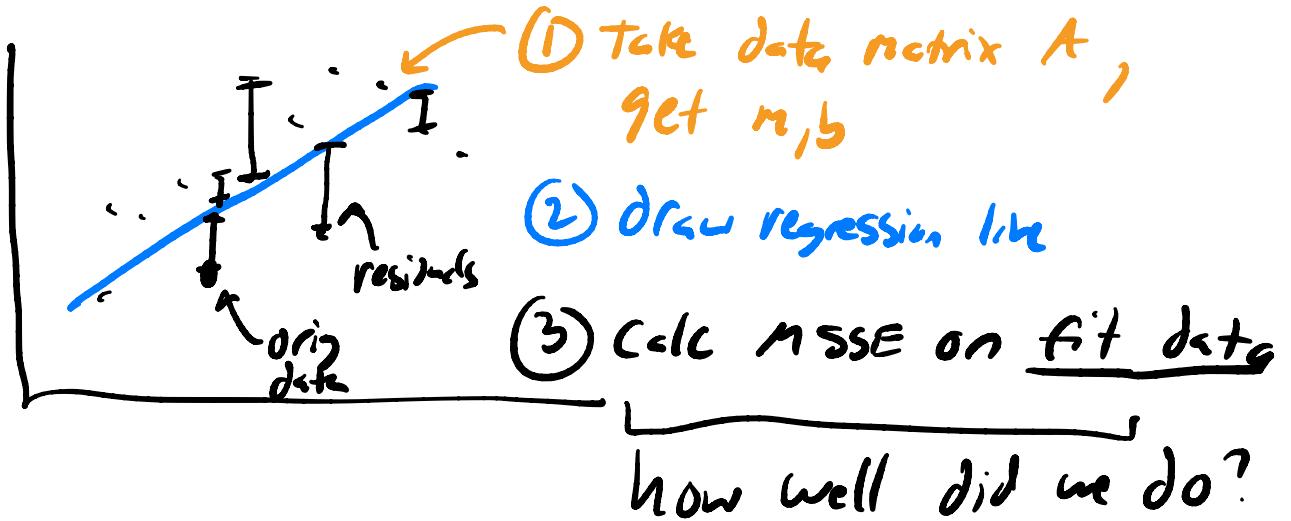
known ↓

predicted

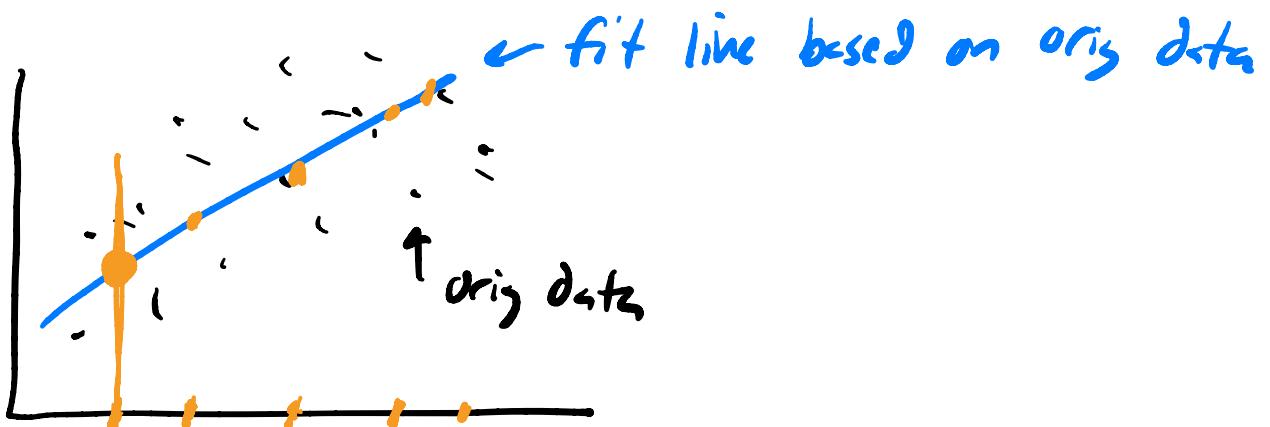
Cancer rate
for $x_i = 3$

Power: predict new values.

- 1) Plug in some independent var values used to fit regression.
- ⇒ Calculate how well regression curve fits original



- 2) plus m new data into regression equation
 → get novel predictions



$x_{\text{New}_i} \rightarrow$ predicted y values