

Predict  
Continuous  
Variable

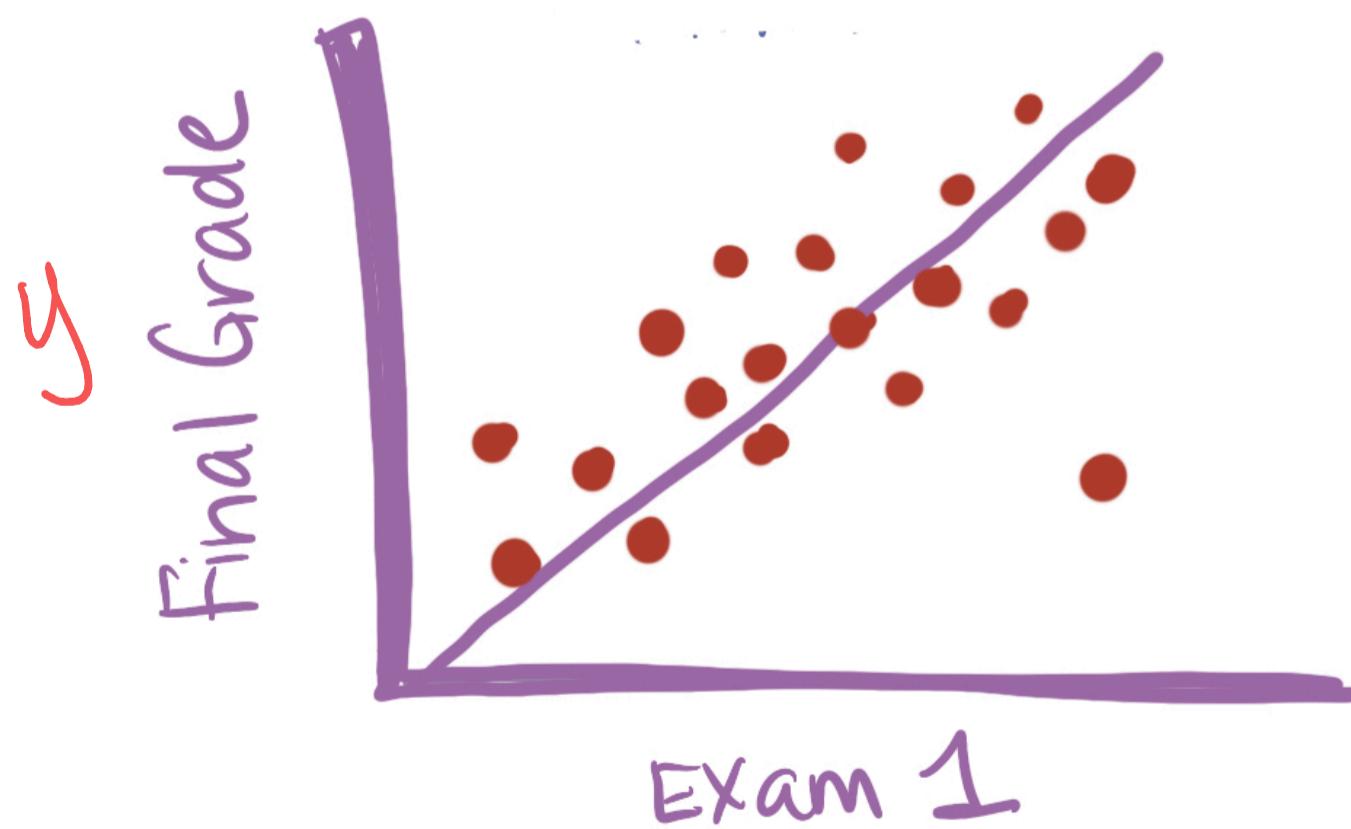
# REGRESSION

VS. CLASSIFICATION

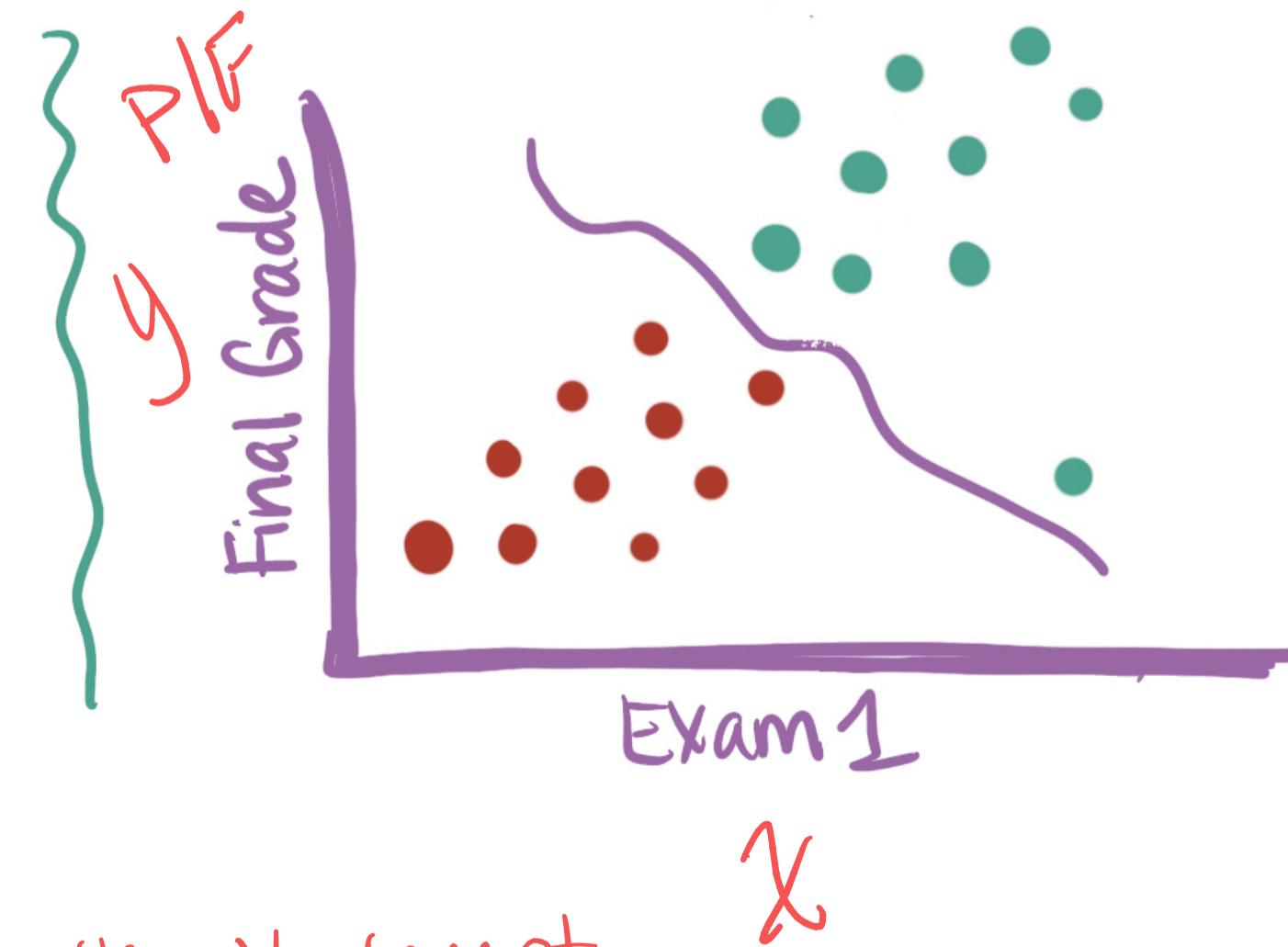
Predict  
Categorical  
Variable

How will each Student perform  
in the course?

Regression: Predict  
Final Grade



Classification: Predict  
Pass or Fail

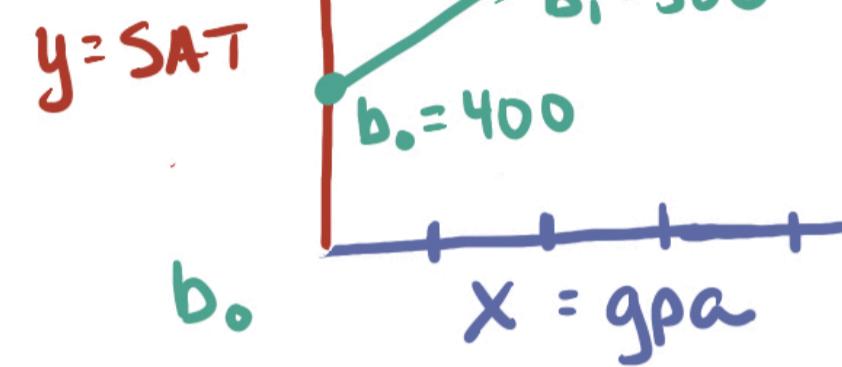


ONLY diff is the  $y_i$  target.  
Input features can be same as classification

## REGRESSION types

### SIMPLE

$$y = b_0 + bx$$

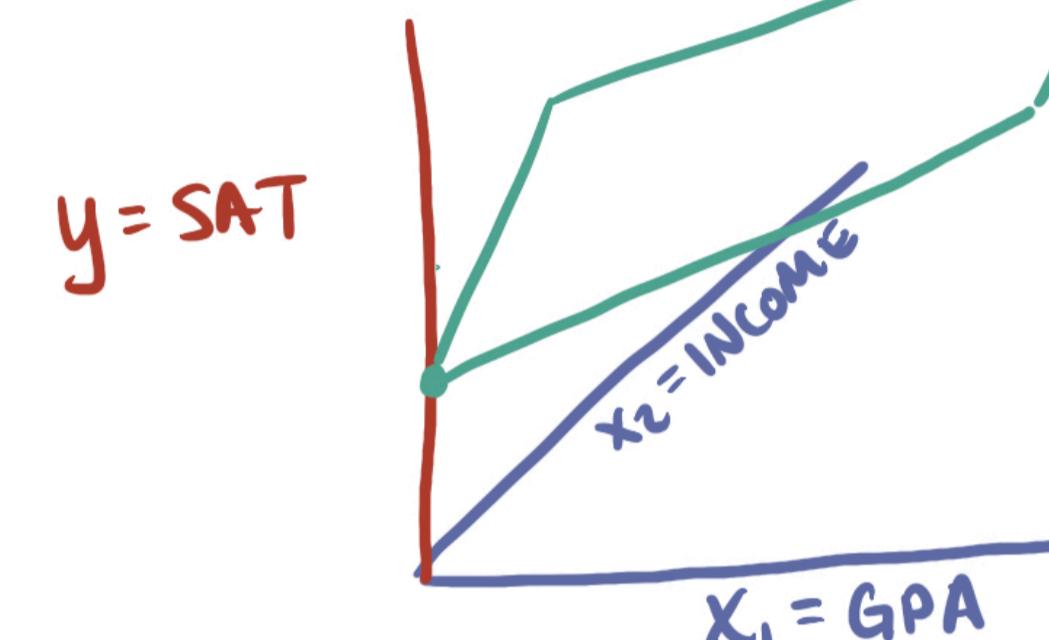


$$\text{SAT} = 400 + 500 \cdot \text{gpa}$$

UNIVARIATE

### MULTIPLE

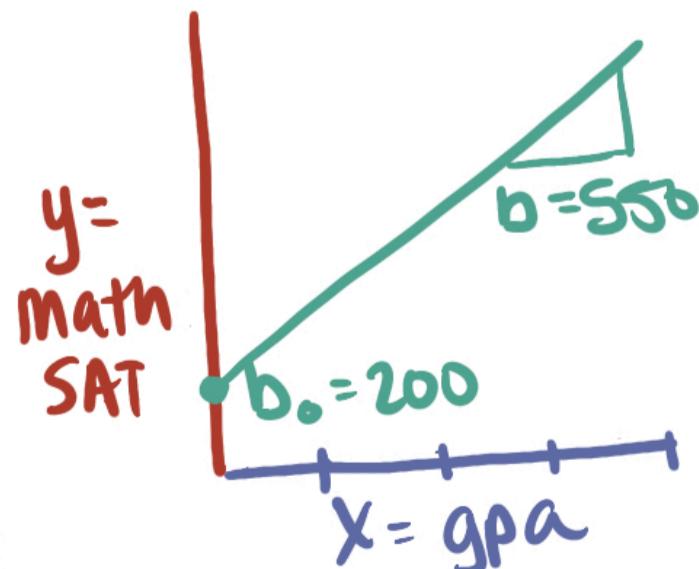
$$y = b_0 + b_1x_1 + \dots + b_nx_n$$



$$\text{SAT} = 400 + 500 \cdot \text{GPA} + .005 \cdot \text{INCOME}$$

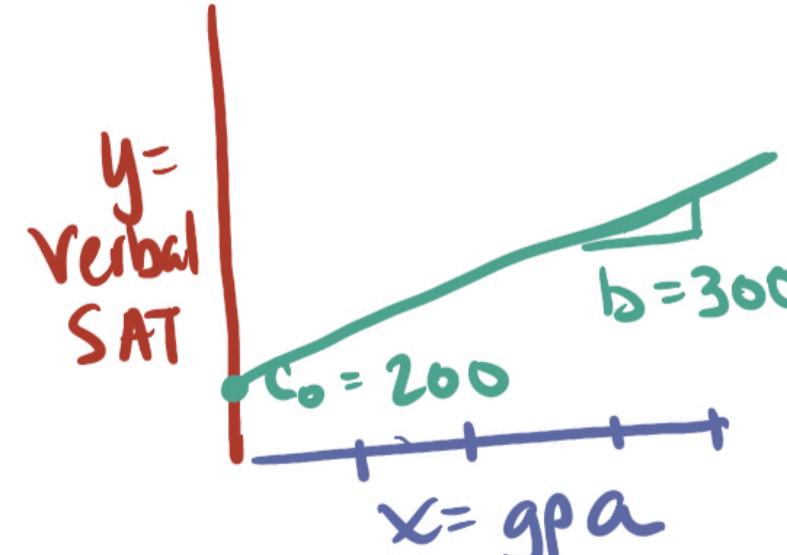
$$y_1 = b_0 + bx$$

$$y_2 = c_0 + cx$$



$$\text{Math SAT} =$$

$$200 + 550 \cdot \text{gpa}$$



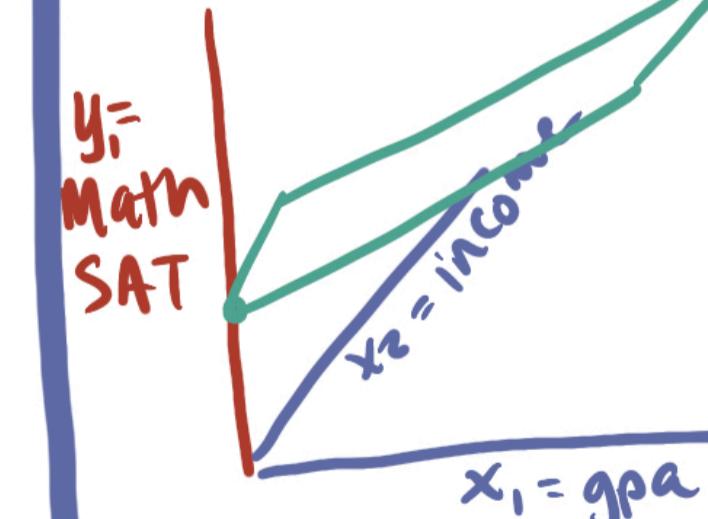
$$\text{Verbal SAT} =$$

$$200 + 300 \cdot \text{gpa}$$

MULTIVARIATE

$$y_1 = b_0 + b_1 + \dots + b_n$$

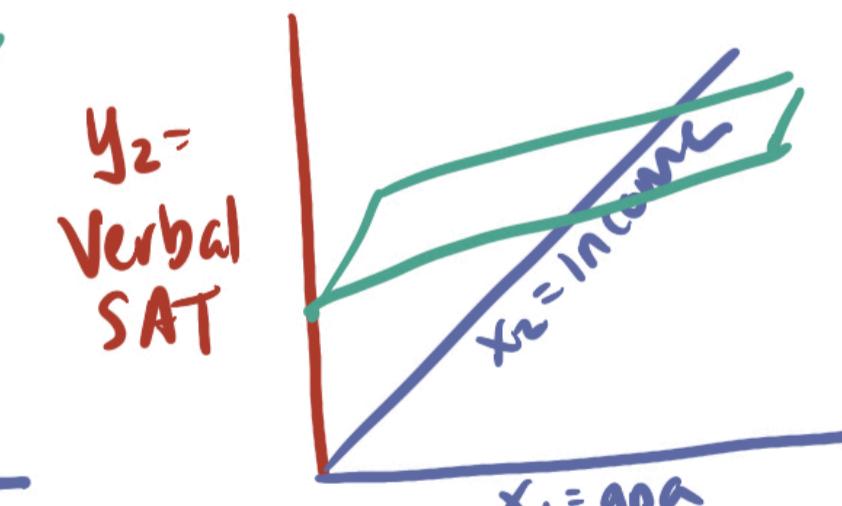
$$y_2 = c_0 + c_1 + \dots + c_n$$



$$\text{Math SAT} =$$

$$200 + 550 \cdot \text{gpa}$$

$$+ .003 \cdot \text{INCOME}$$



$$\text{Verbal SAT} =$$

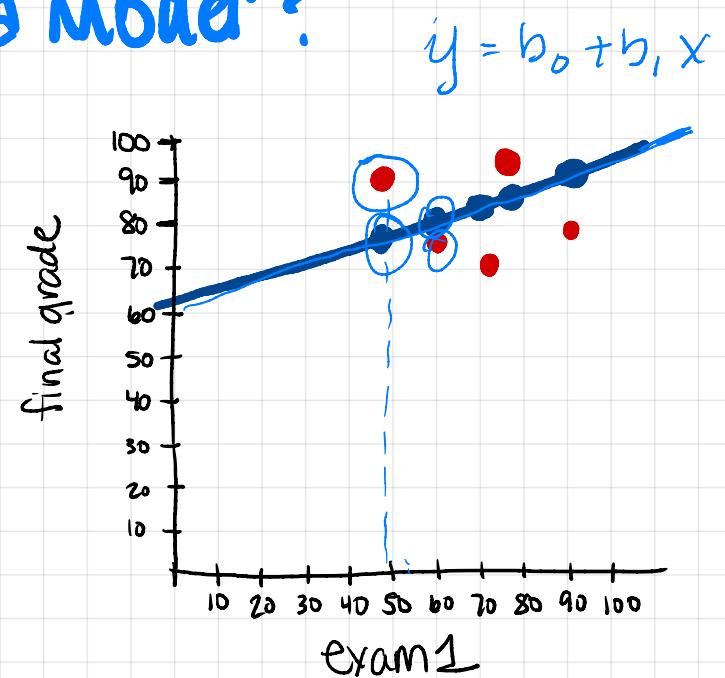
$$200 + 300 \cdot \text{gpa}$$

$$+ .002 \cdot \text{INCOME}$$

# evaluating regression Models

What does it mean to "evaluate a model"?

StudentID	Exam1 x	Final grade y	Predicted final grade $\hat{y}$
1	48	90	77
2	60	76	80
3	71	70	84
4	76	94	87
5	90	80	92



Evaluating  $\Rightarrow$

How well did my model predict the students' final grades?

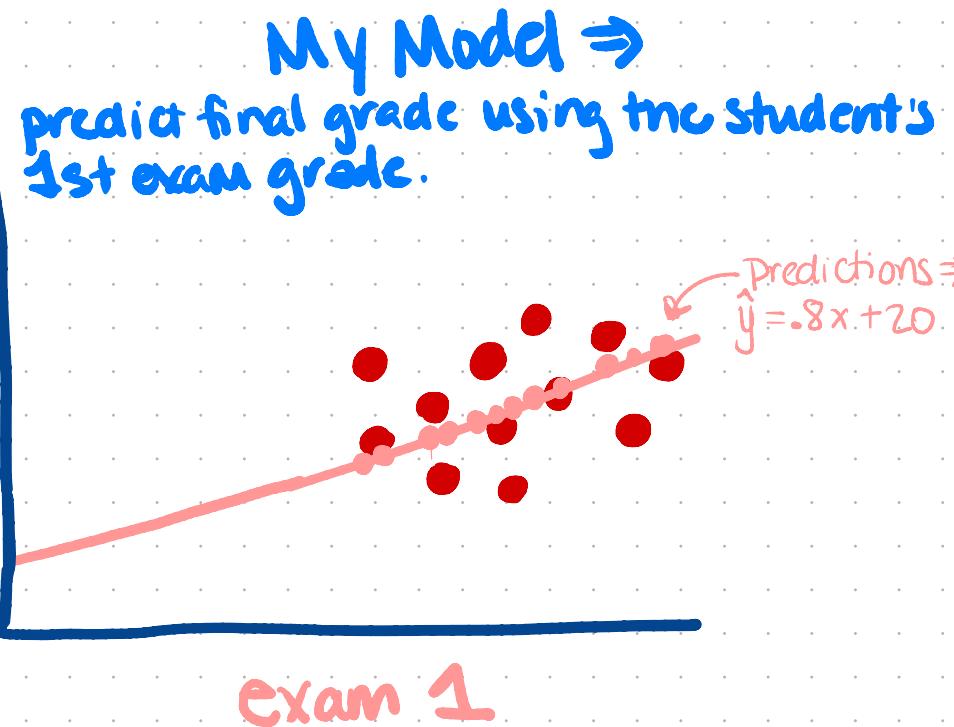
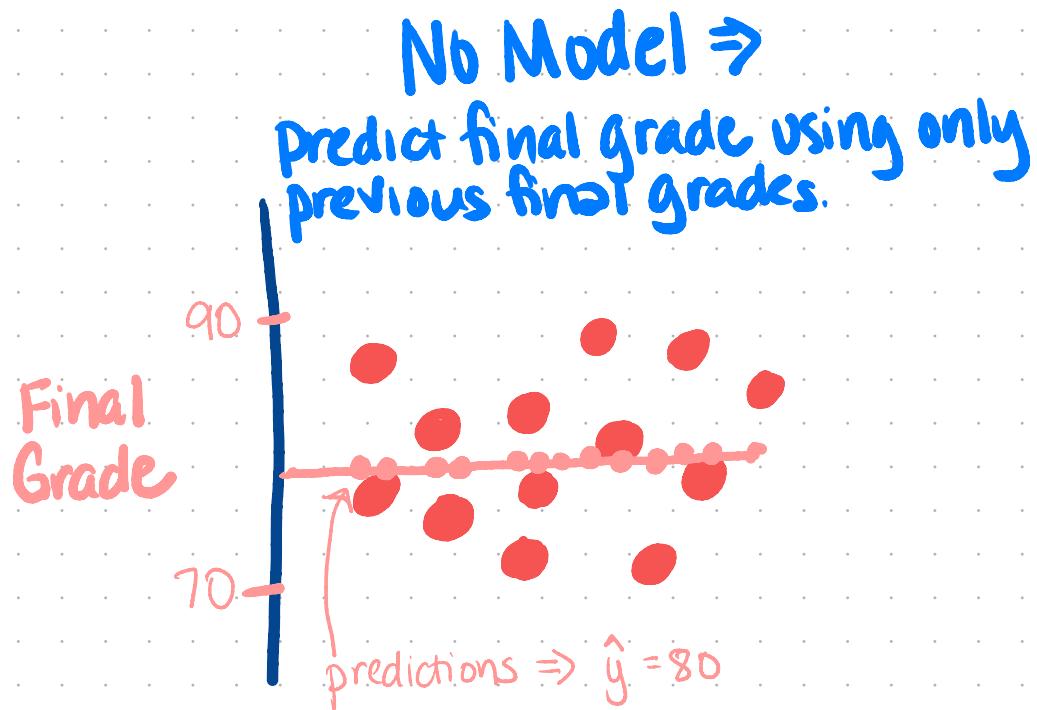
# Why evaluate models?

Q1 Does this model add ANY value?

Q2 Which of these 2 models is better?

Q3 How much confidence should I have in this model?

**Q1** Does my model add ANY value?  
Is it "Good Enough"?



Which is better?

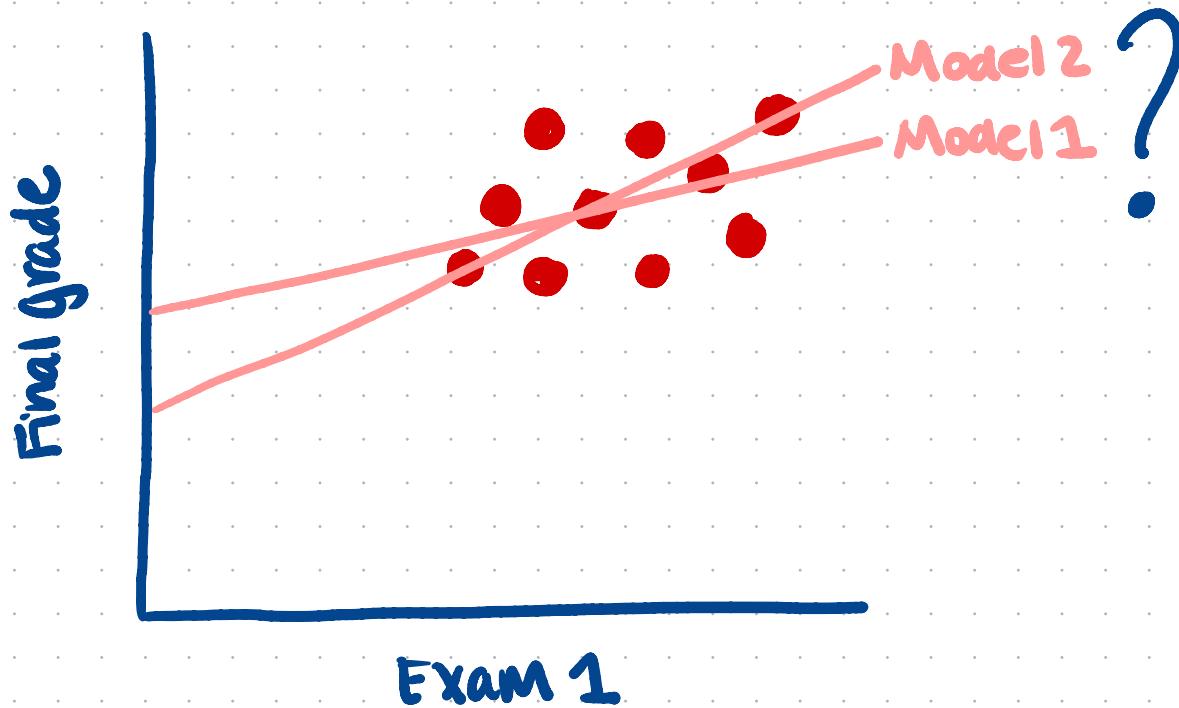
Predict final grade  
for each student to be 80! ??

Predict the final grade for each  
student to be  
80% of their 1st exam score  
+ 20 points.

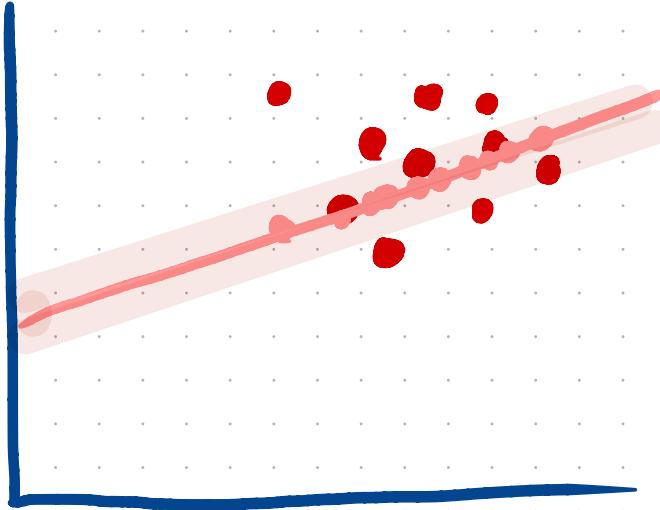
Q2:

Which Model is better?

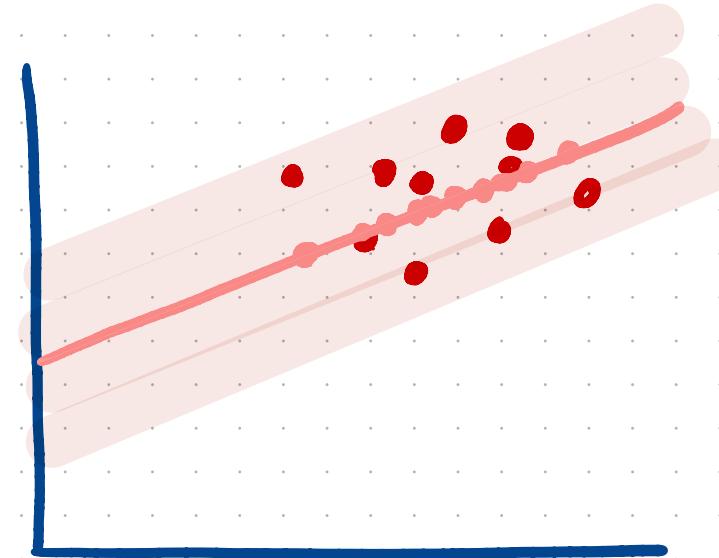
How can I determine which predicts better?



### Q3: How much Confidence Should I have in this Model?



It's gonna be close!



Yeah, we really have no idea ...

[Does my model add value?] [Which Model is better?]

To Answer, we need a **single metric** that measures how well a model predicts.

Common Metrics : [Root Mean Squared Error]

Mean Absolute Error

Median Absolute Error

[Mean Squared Error]

RMSE (via MSE)

Manual calculation

RESIDUALS

Error of each prediction  $\Rightarrow$  square each error  $\Rightarrow$  Avg of the squared errors  $\Rightarrow$  Sqrt of the avg.

$$\hat{y}_1 - y_1 = r_1$$

$$\hat{y}_2 - y_2 = r_2$$

:

$$\hat{y}_n - y_n = r_n$$

residuals

$n = \# \text{ of observations / rows}$

$r_i = \text{residual/error for the } i\text{th row/observation}$

$\hat{y}_i = \text{predicted value for the } i\text{th row/observation}$

$y_i = \text{actual value for the } i\text{th row/observation}$

SSE

then sum

MSE

Avg of the squared errors

RMSE

Sqrt of the avg.

$$r_1^2 + r_2^2 + \dots + r_n^2 = SSE \rightarrow \frac{SSE}{n} = MSE \rightarrow \sqrt{MSE} = RMSE$$

RMSE in plain English:

"Our predictions are off by an average of 5 pts on the final grades."

RMSE : `sklearn.metrics.mean_squared_error  
math.sqrt`

RMSE =  
`sqr((mean_squared_error(df.y, df.yhat)))`

# Does my Model work better than having no Model? [guessing? pred. mean? pred mode?]

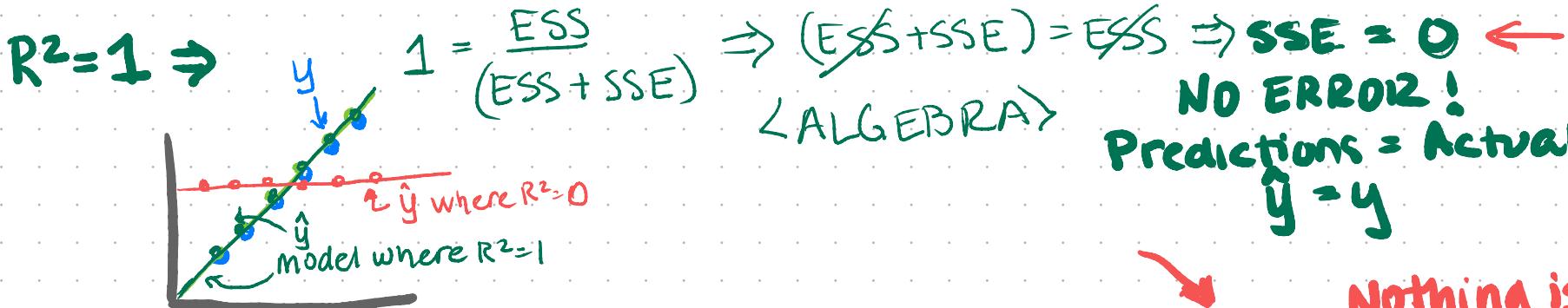
\* Applies to ordinary least squares algorithms such as  
statsmodels.api.OLS & sklearn.linear\_model.LinearRegression

**Answer:** run an f-regression test on  $R^2$   
if the p-value of the f-statistic  
is  $< \alpha (n, 05)$ , then our model is  
better than having no model.

$R^2$  = the explained variance score  
= the coefficient of determination  
=(Pearson's R)<sup>2</sup> (when simple  
univariate regression)

**$R^2$**

$$R^2 = \frac{ESS}{TSS} = \frac{\text{Explained Sum of Squares}}{\text{Total Sum of Squares}} = \frac{ESS}{(ESS + SSE)} = \frac{\text{Explained SS}}{(\text{Explained SS} + \text{Error})}$$



$R^2 = 0 \Rightarrow 0 = \frac{ESS}{(ESS + SSE)} \Rightarrow (ESS + SSE) \cdot 0 = ESS \Rightarrow 0 = ESS$

Nothing is explained!  
All error!

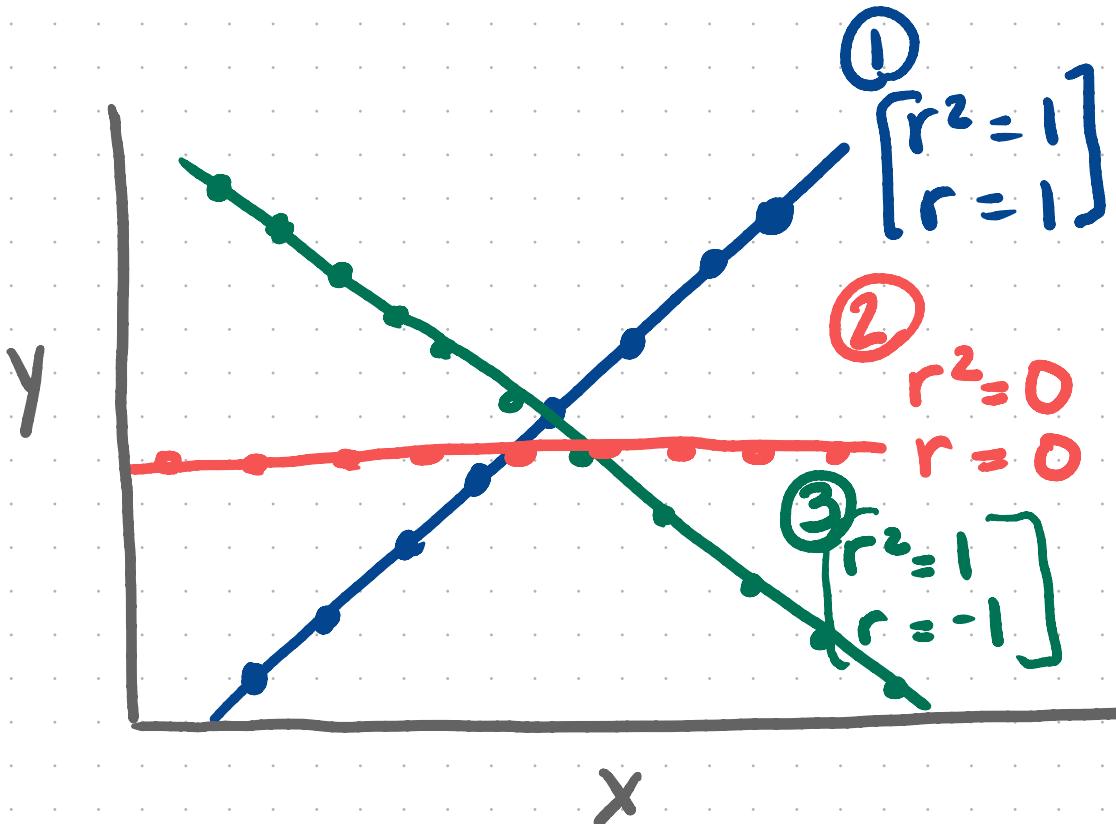
sklearn.metrics.explained\_variance\_score(df.y, df.yhat)

statsmodels.api.OLS => MyOLSModel.rsquared

P-Value

myOLSModel.f\_pvalue | sklearn.feature\_selection.  
f\_regression

$$r \neq r^2$$



1. As  $x \uparrow, y \uparrow$  by a constant value
2. As  $x \uparrow, y \downarrow$  by a constant value
3. As  $x \uparrow, y$  doesn't change  $\Rightarrow$   
 $x$  does not affect  $y$  in ANY way

- $r \rightarrow$  • Perfect positive correlation
- $r^2 \rightarrow$  • The variance in  $x$  explains ALL of the variance in  $y$
- $r \rightarrow$  • Perfect negative correlation
- $r^2 \rightarrow$  • The variance in  $x$  explains ALL of the variance in  $y$
- $r \rightarrow$  No correlation
- $r^2 \rightarrow$  The variance in  $x$  explains None of the variance in  $y$ .

## Evaluate Models

- Always use
  - 1. A metric related to the errors/residuals (RMSE, MSE, MAE)
  - 2. R<sup>2</sup> and f-regression test
    - USE ONLY WITH OLS, linear regression models.
- Sometimes use
  - USED MORE IN PRACTICE
  - more options for other types of regression

Feature Selection  $\Rightarrow$  Evaluate Features

Algorithmic way to evaluate the impact features have on a target

### 1. Recursive Feature Elimination (RFE)



[ML Model = Linear-Regression  
# of features = 2]

### 2. Select K Best

- run a statistical test, like f-regression, on each [feature + target] combination
- Select the K features with the best score ( $\rightarrow$  lowest p-value, e.g.)

[Statistic Score, K = # of features]

# Select K Best

## Select the best 'k' features

**K?** # of desired features

**Best?** best is defined by a statistical test -  
the test selected depends on the type of problem.

### Test / score functions

#### REGRESSION:

- f\_regression : test correlation (LINEAR\*)
- mutual\_info\_regression : test dependence
  - 0  $\Rightarrow$  completely independent
  - 1  $\Rightarrow$  completely dependent

#### CLASSIFICATION:

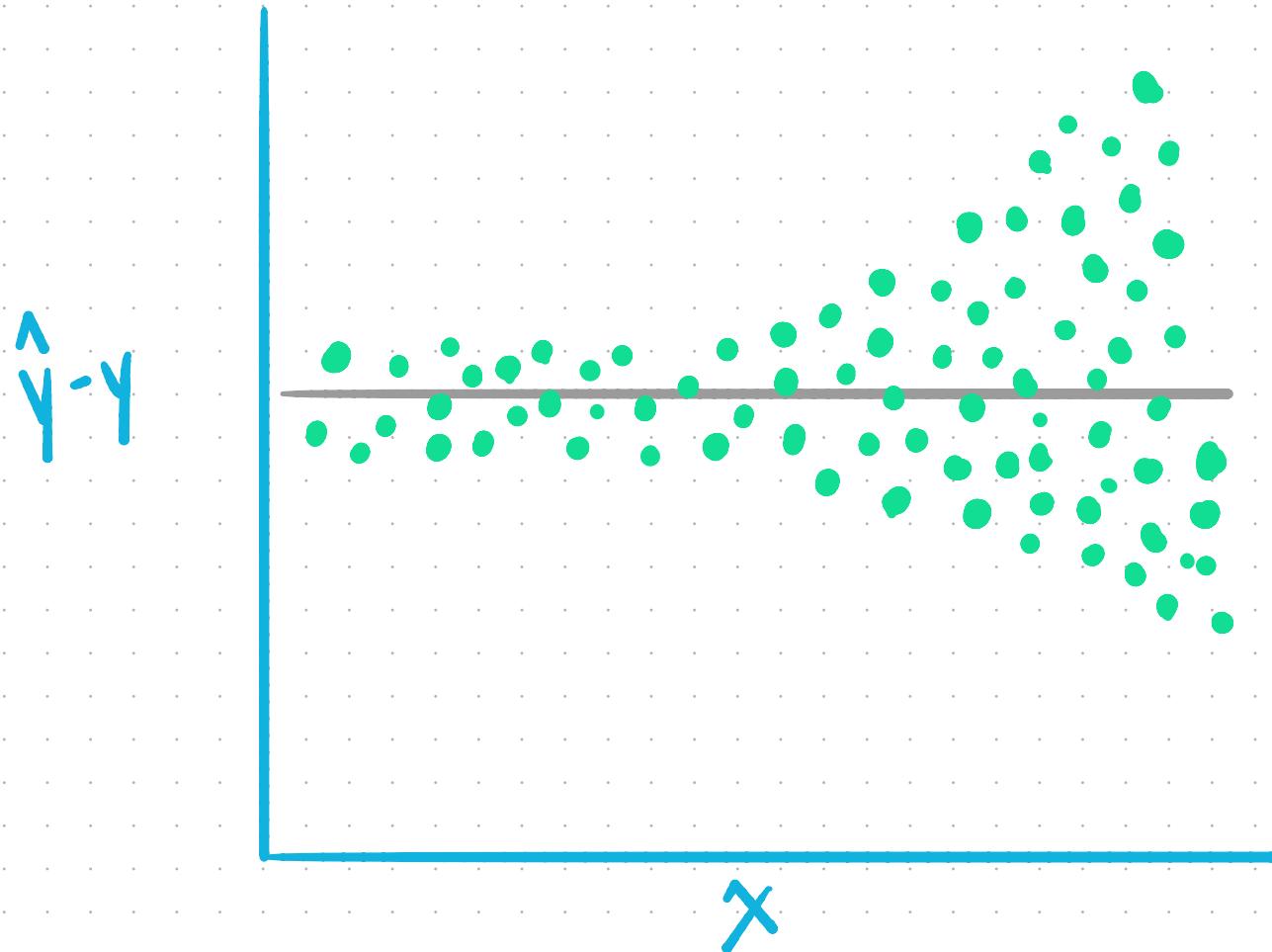
- f\_classif
- mutual\_info\_classif
- chi2
- SelectFPR

...

$X_{new} = \text{Sklearn.feature-Selection.SelectKBest(f\_regression, k=3)}$

# Heteroskedasticity

Non-constant variance  
in ERRORS



Heteroskedasticity will bias standard errors

# MEAN ABSOLUTE ERROR

	$y$	$\hat{y}$	Error
1	$y_1$	$\hat{y}_1$	$y_1 - \hat{y}_1 = r_1$
2	$y_2$	$\hat{y}_2$	$y_2 - \hat{y}_2 = r_2$
:			:
$n$	$y_n$	$\hat{y}_n$	$y_n - \hat{y}_n = r_n$

$\hookrightarrow n$  observations

$$MAE = \frac{|r_1| + |r_2| + \dots + |r_n|}{n}$$

$$= \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$y_i$  = actual value for observation / row i

$\hat{y}_i$  = predicted target value for observation / row i

$r_i$  = residual, error, difference between actual & predicted value.

# MEAN SQUARED ERROR

	$y$	$\hat{y}$	Error
1	$y_1$	$\hat{y}_1$	$y_1 - \hat{y}_1 = r_1$
2	$y_2$	$\hat{y}_2$	$y_2 - \hat{y}_2 = r_2$
:			:
$n$	$y_n$	$\hat{y}_n$	$y_n - \hat{y}_n = r_n$

→  $n$  observations

$$\text{MSE} = \frac{r_1^2 + r_2^2 + \dots + r_n^2}{n}$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$y_i$  = actual value for observation/row  $i$

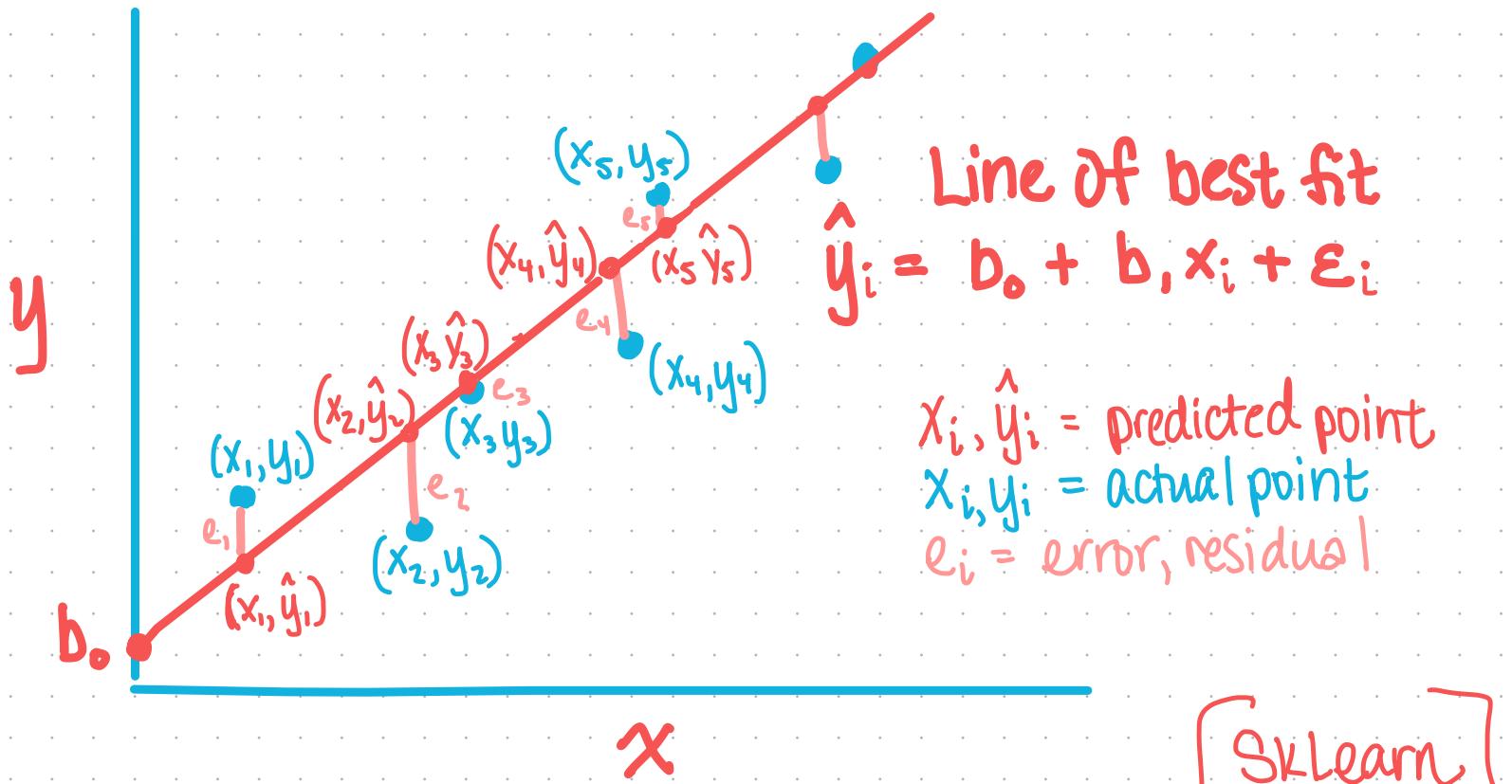
$\hat{y}_i$  = predicted target value for observation/row  $i$

$r_i$  = residual, error, difference between actual & predicted value.

# Ordinary Least Squares

## OLS

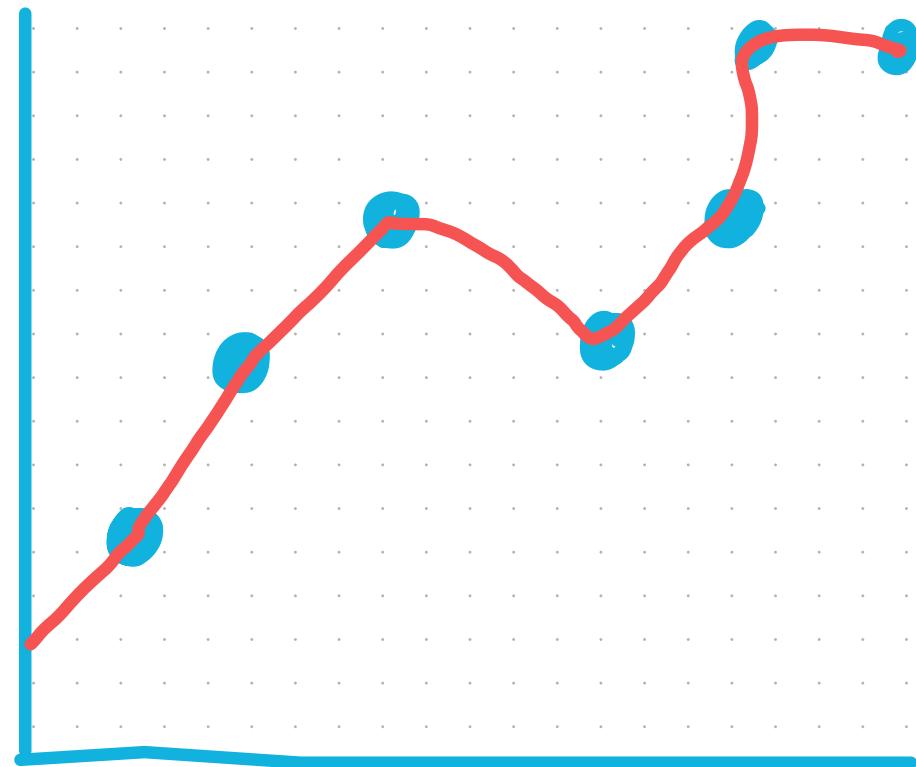
- Linear regression algorithm that finds the line of best fit.  $e_1^2 + \dots + e_n^2$  that returns the smallest SSE =  $(e_1^2 + e_2^2 + \dots + e_n^2)$
- Often used as a baseline.



`y-pred = linear_model.LinearRegression().fit(x,y).predict(x)`

[SKLearn]

# OVERFITTING

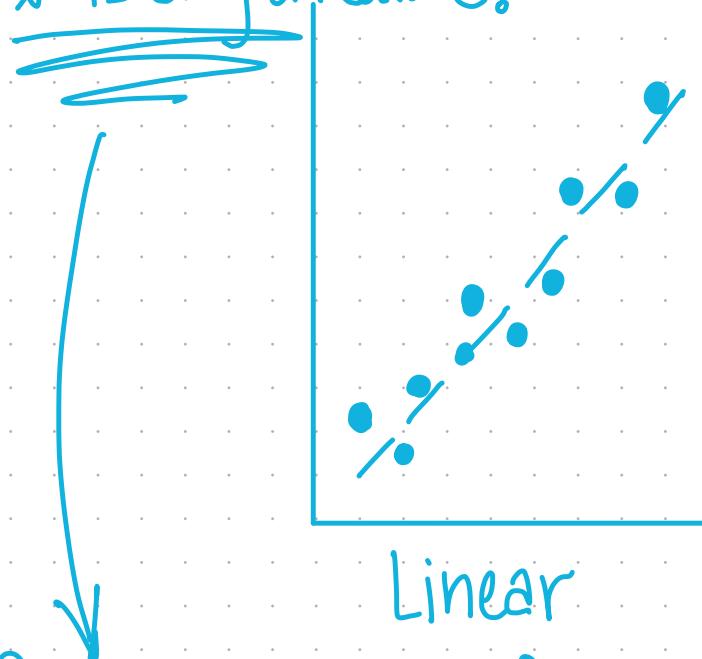


When a model "memorizes" specifics of training data  
And then loses the ability to generalize.

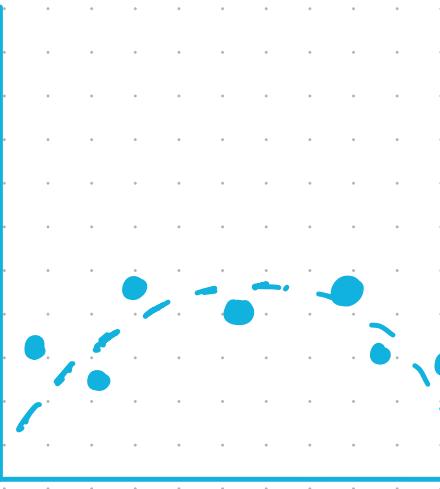
# POYNOMIAL REGRESSION

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_d x_i^d + e_i$$

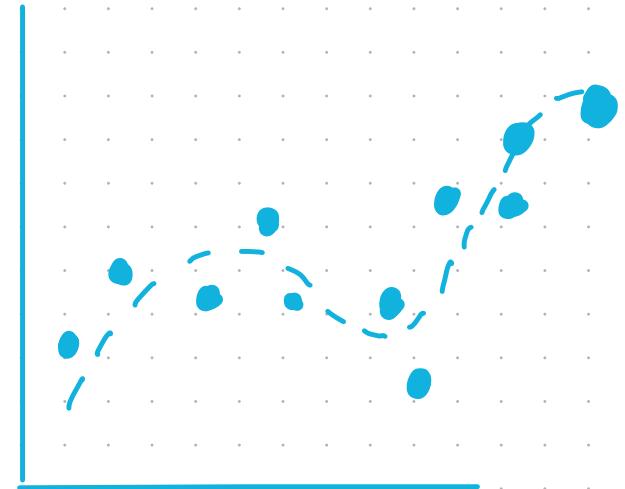
Still Considered to be a linear model → the coeff are still linear!  
 $x^2$  is only a feature.



Linear



Polynomial



Polynomial

So... the original features must be converted into polynomial features.

[  
x\_poly = preprocessing.PolynomialFeatures(degree=2).fit\_transform(x)  
y\_poly\_pred = linear\_model.LinearRegression().fit\_predict(x\_poly, y)  
]

Sklearn