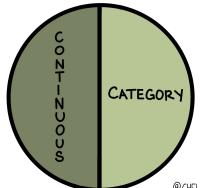


PREDICT



@CHELSEAPARLETT

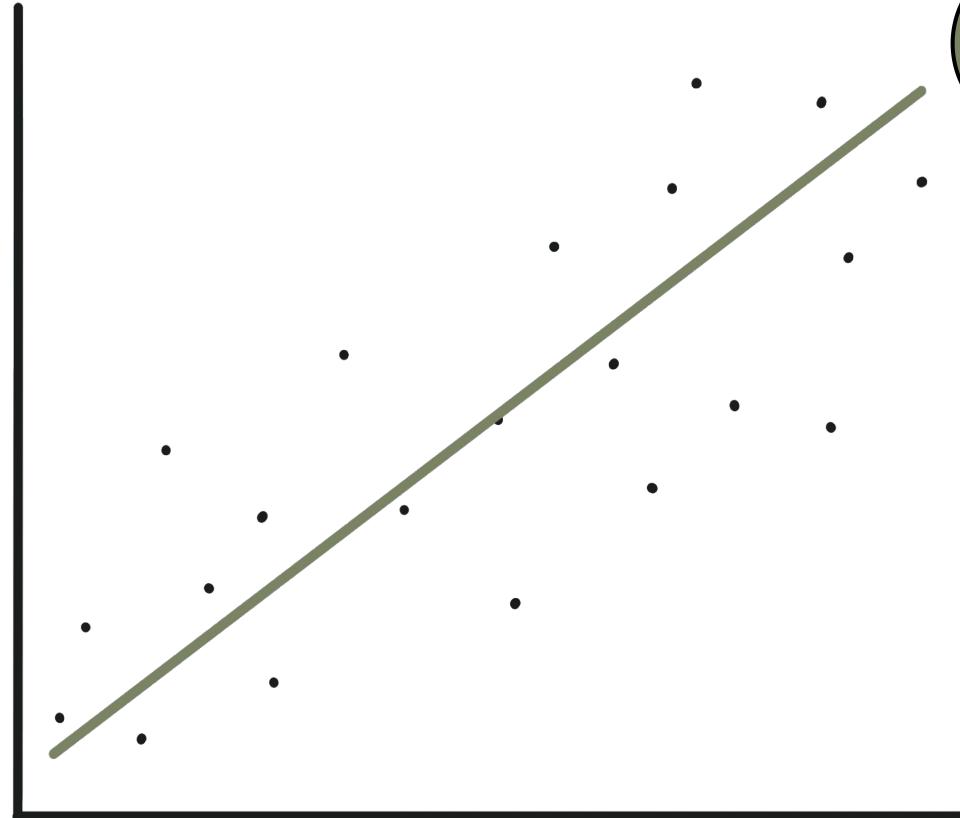
Linear Regression

Chelsea Parlett-Pelleriti

What

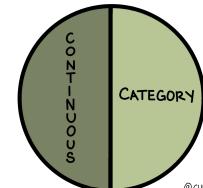
- Use **multiple variables** (can be continuous, categorical, or both) to predict a **continuous variable**.
- Use a line (or a plane) to describe the relationship between these variables.

numt



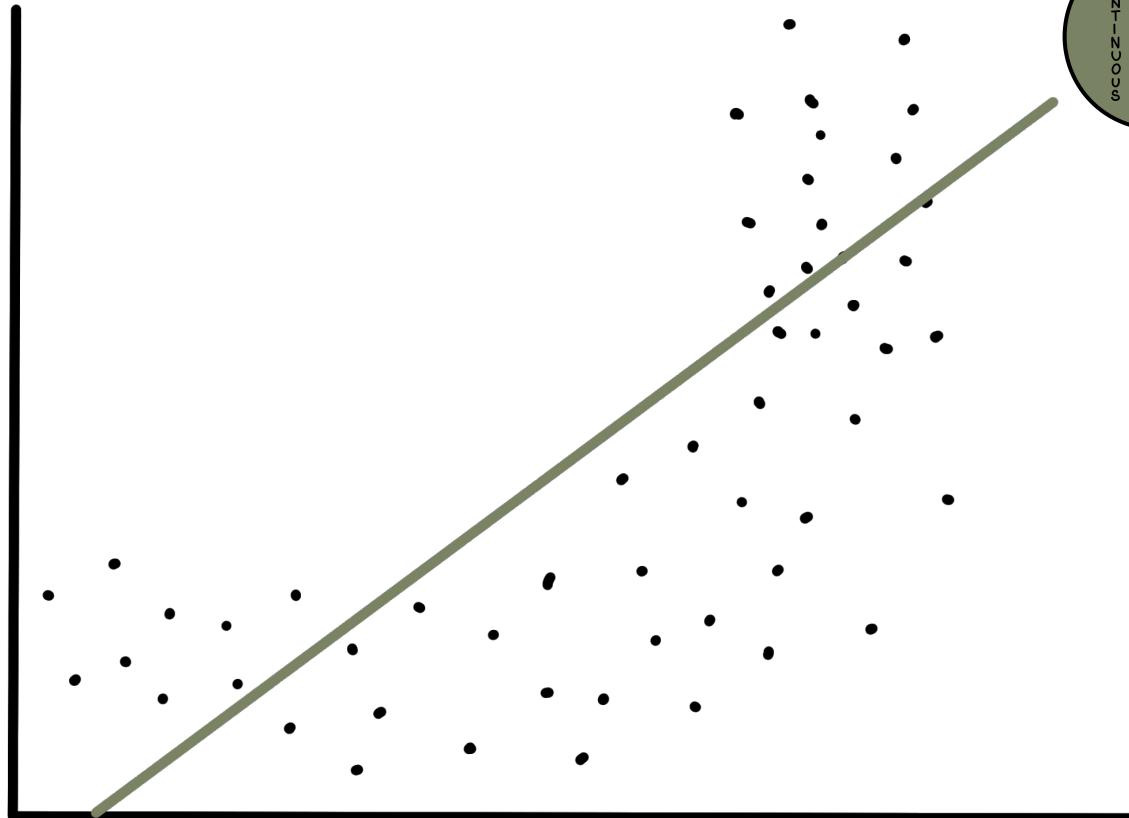
income

PREDICT

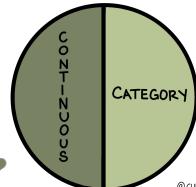


Assumptions

- The relationship between your variables is linear



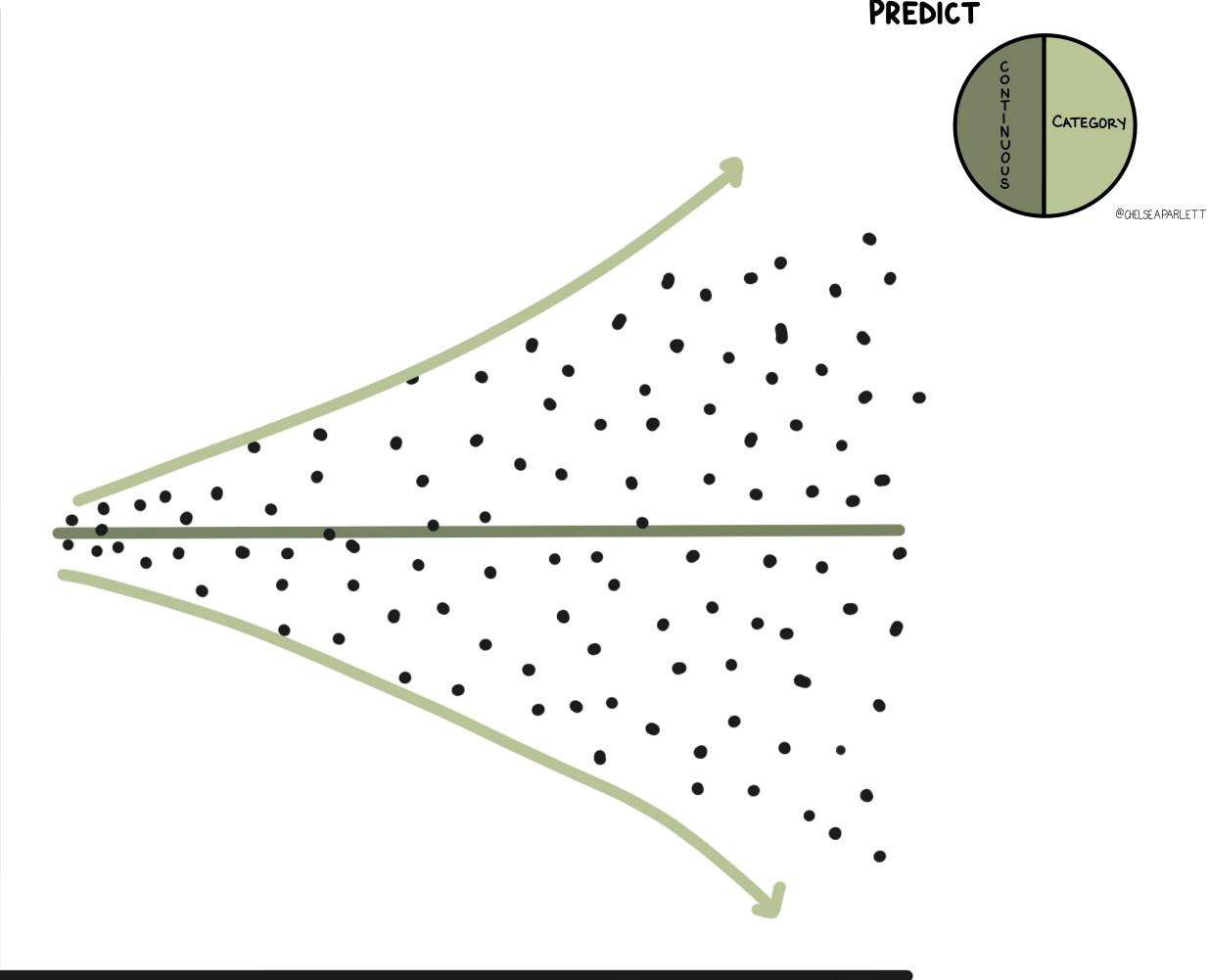
PREDICT

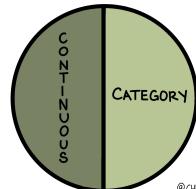


@CHELSEAPARLETT

Assumptions

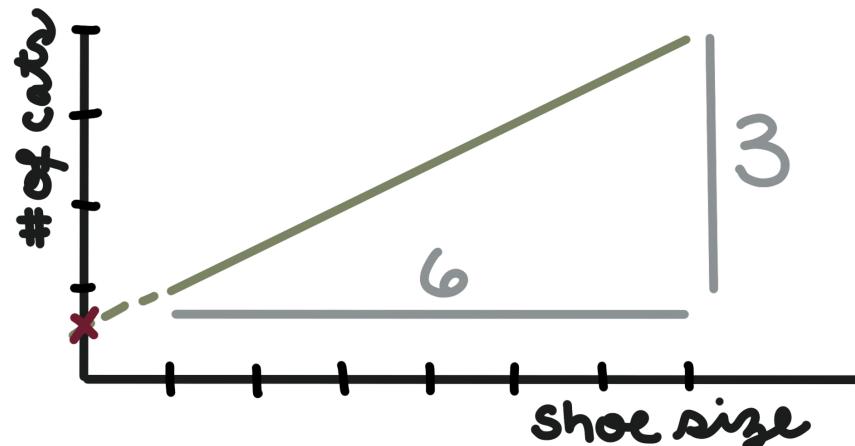
- Homoscedasticity
 - Is the mode worse in some areas than others?
- Normality of Errors



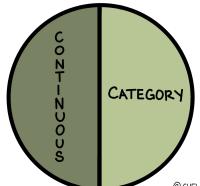


How

- $Y = mx + b$
- $Y = mx + nz + b$
- Slope tells you how variables change together
- Intercept tells you what would happen if all your predictors were 0.

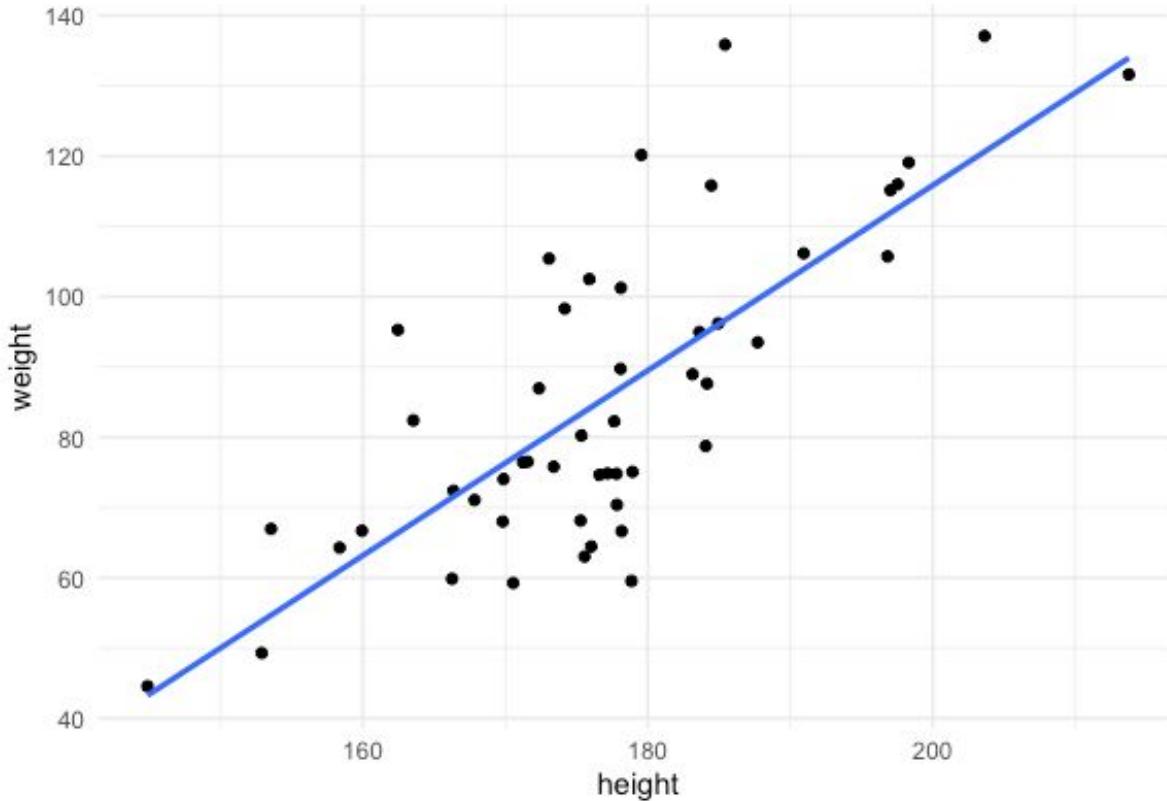


PREDICT



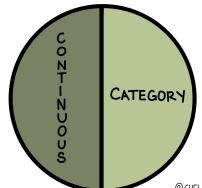
Simple example

Predict weight by height



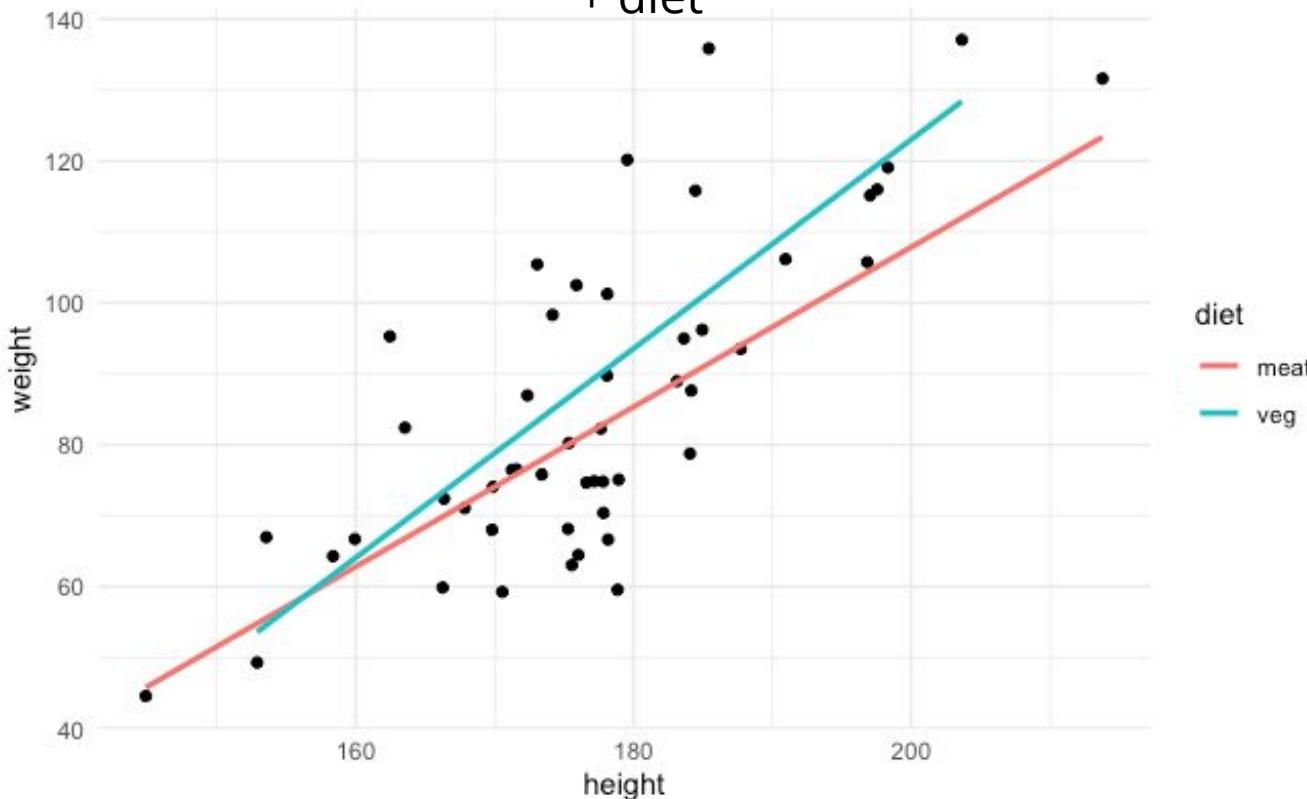
coef
Intercept -82.2887
height 0.9786

PREDICT



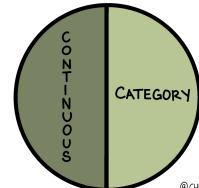
Simple example

Predict weight by height
+ diet



coef	
Intercept	-72.0358
diet[T.veg]	-7.6222
height	0.9420

PREDICT



Simple example

Predict weight by height
+ diet + age

coef	
Intercept	-57.4078
diet[T.veg]	-8.2640
height	0.8948
age	-0.1298

Who is the GOAT?



378 three-pointers

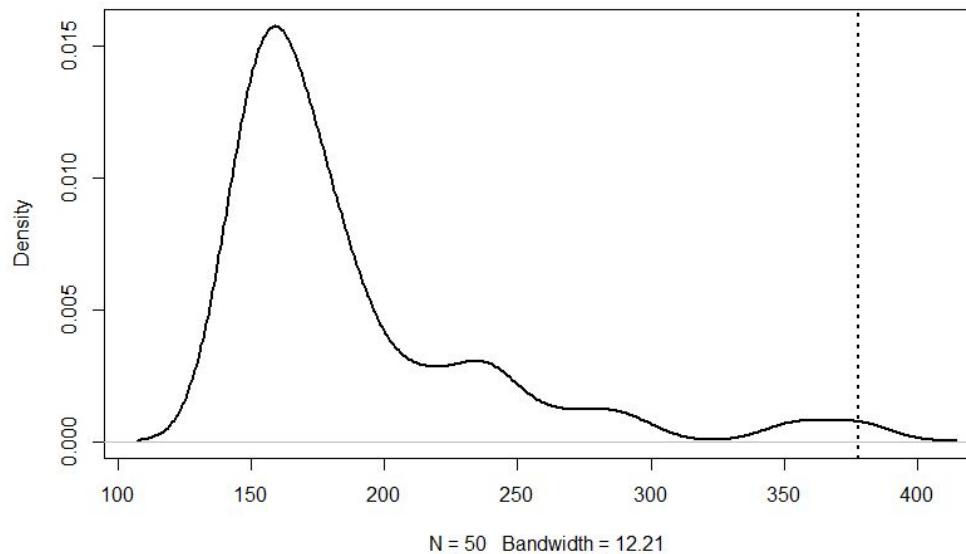


53 home-runs

Who is the GOAT?



Basketball



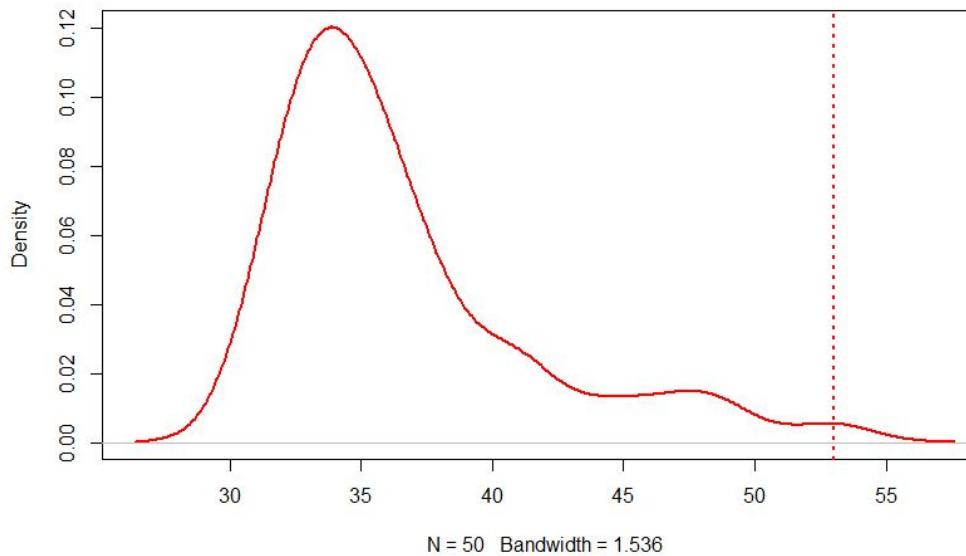
[2018-19 NBA Regular Season: Total 3-Pointers Made Leaders](#)

[2019 MLB Player Batting Stats | Home Runs](#)

Who is the GOAT?



Baseball



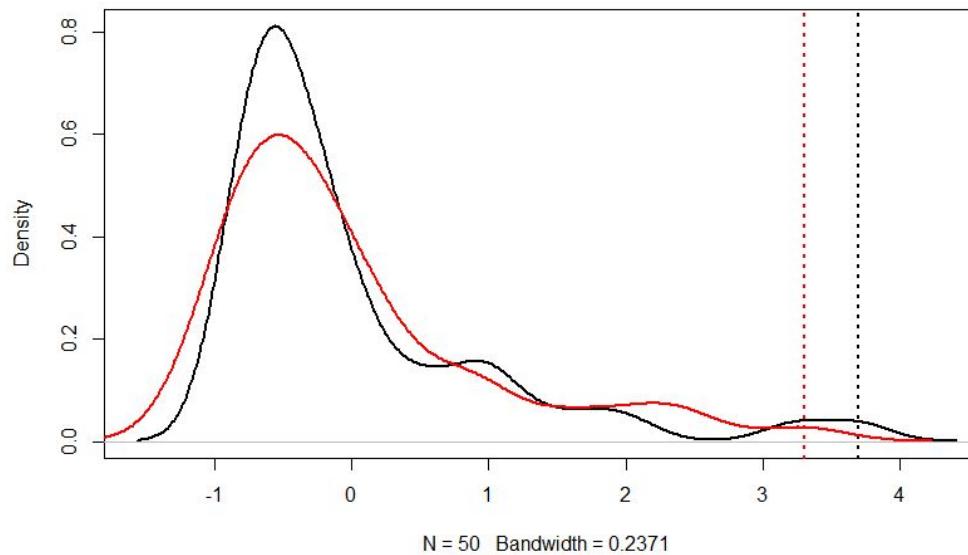
[2018-19 NBA Regular Season: Total 3-Pointers Made Leaders](#)

[2019 MLB Player Batting Stats | Home Runs](#)

Who is the GOAT?



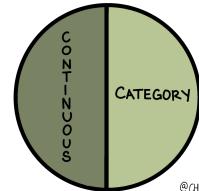
Both Std.



[2018-19 NBA Regular Season: Total 3-Pointers Made Leaders](#)

[2019 MLB Player Batting Stats | Home Runs](#)

PREDICT

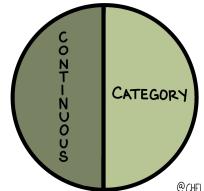


Simple example

Predict weight by height
+ diet + age

	coef
Intercept	93.6861
diet[T.veg]	-8.2640
height	13.4689
age	-2.5245

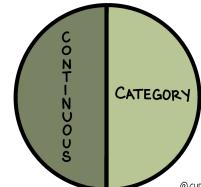
PREDICT



Standardizing variables

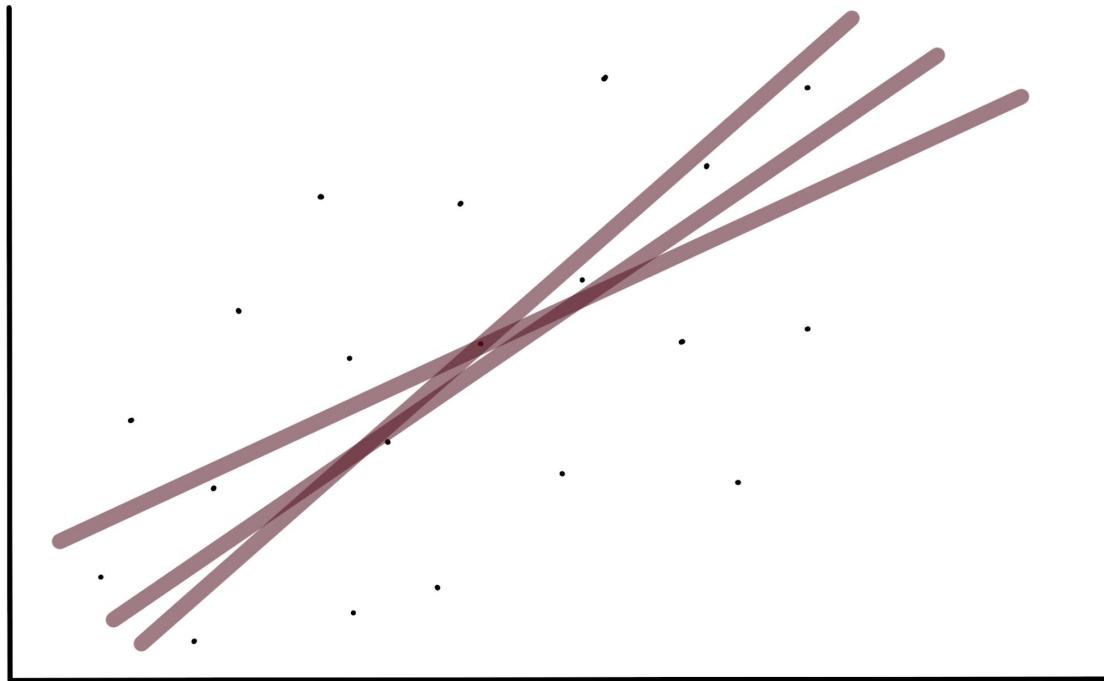
For understanding and for model convergence

PREDICT

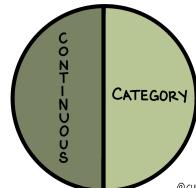


@CHELSEAPARLETT

Choosing the line of best fit

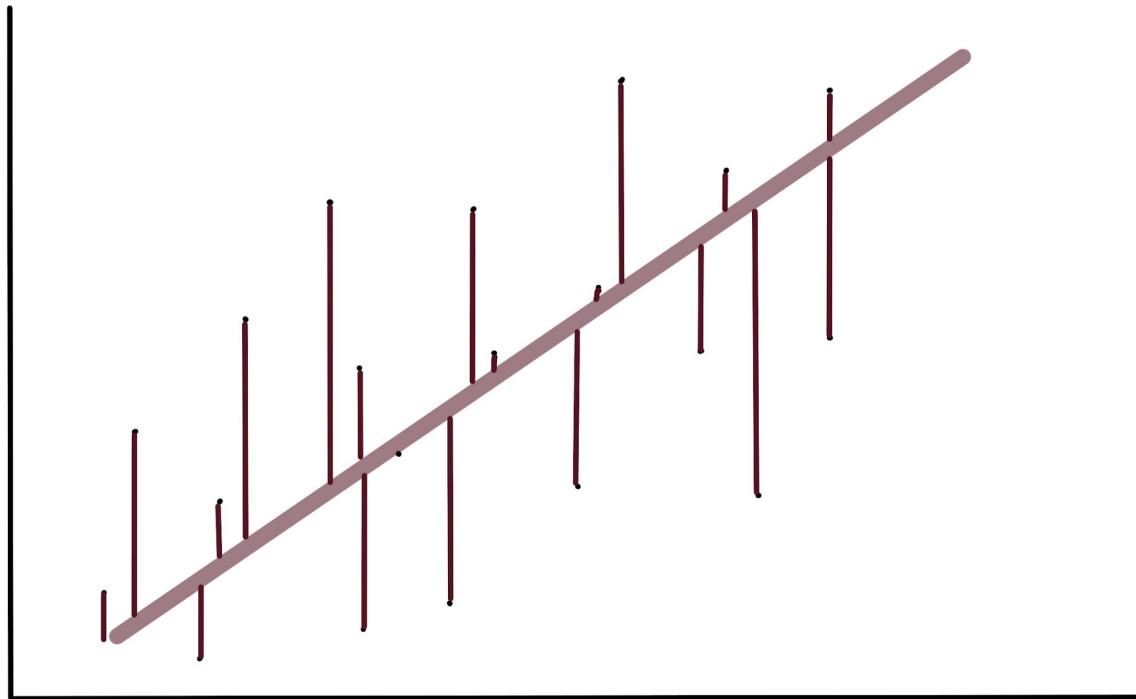


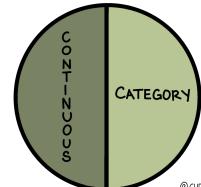
PREDICT



@CHELSEAPARLETT

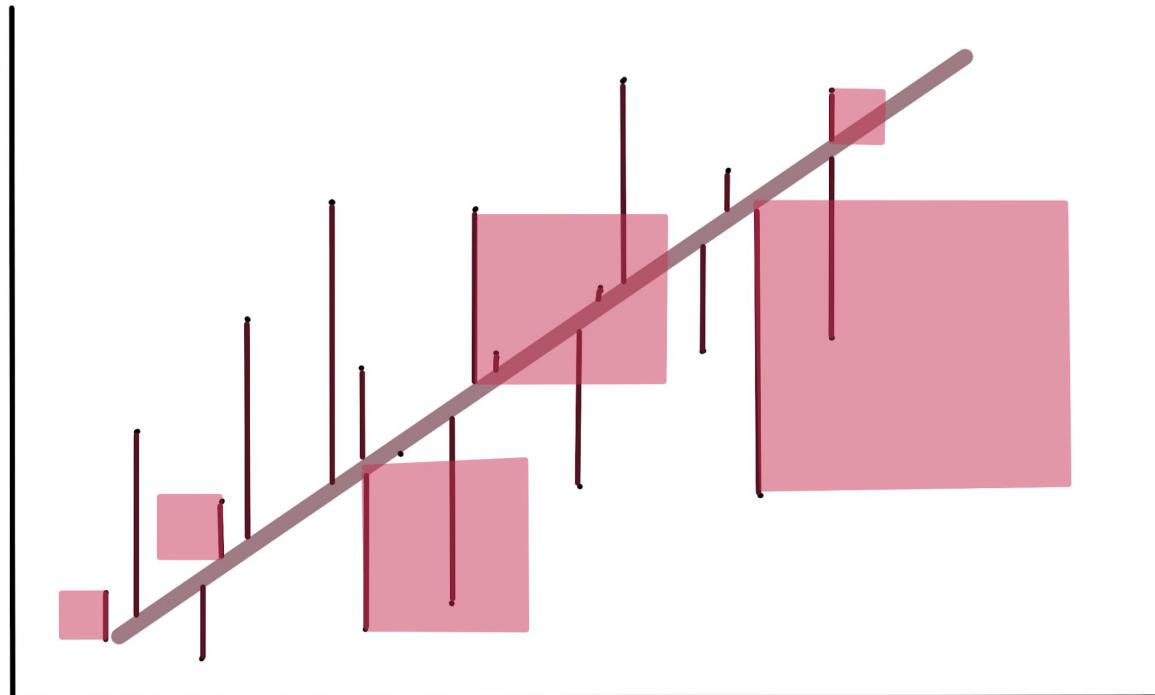
Choosing the line of best fit



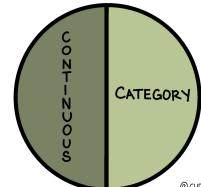


Choosing the line of best fit

- Sum of Squared Errors
- Mean Squared Error

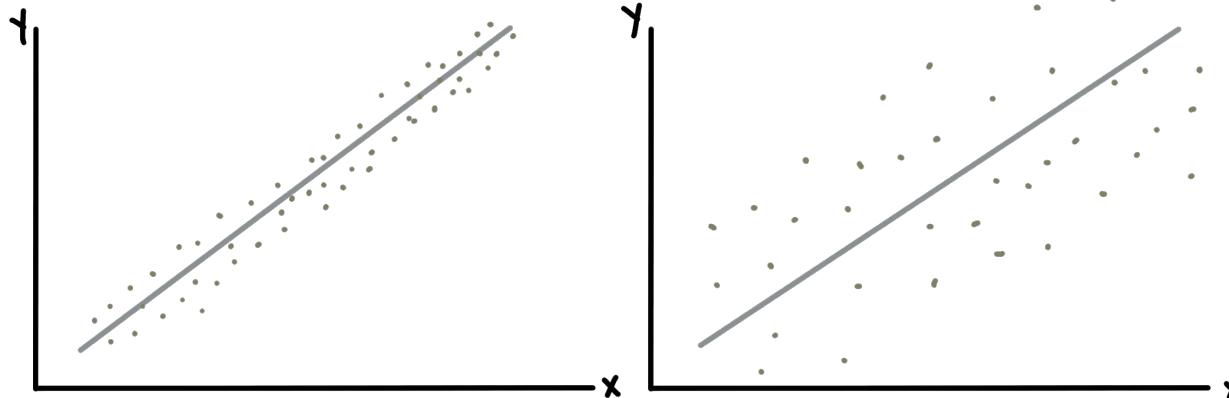


PREDICT



Choosing the line of best fit

- Sum of Squared Errors
- **Mean Squared Error**

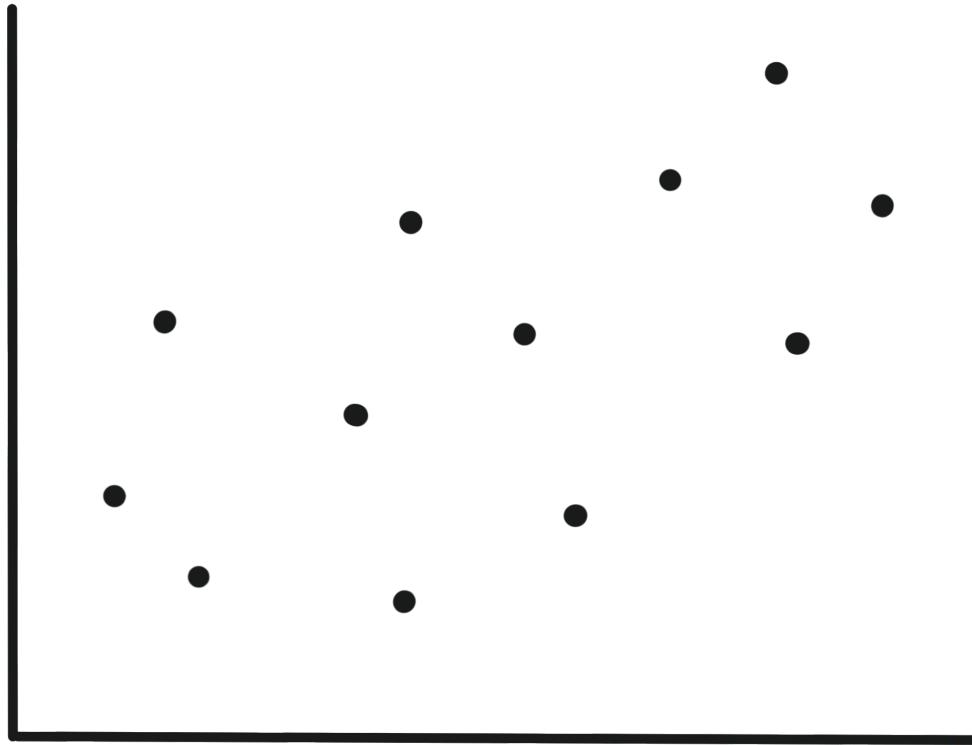


How to Measure Model Success

MSE:

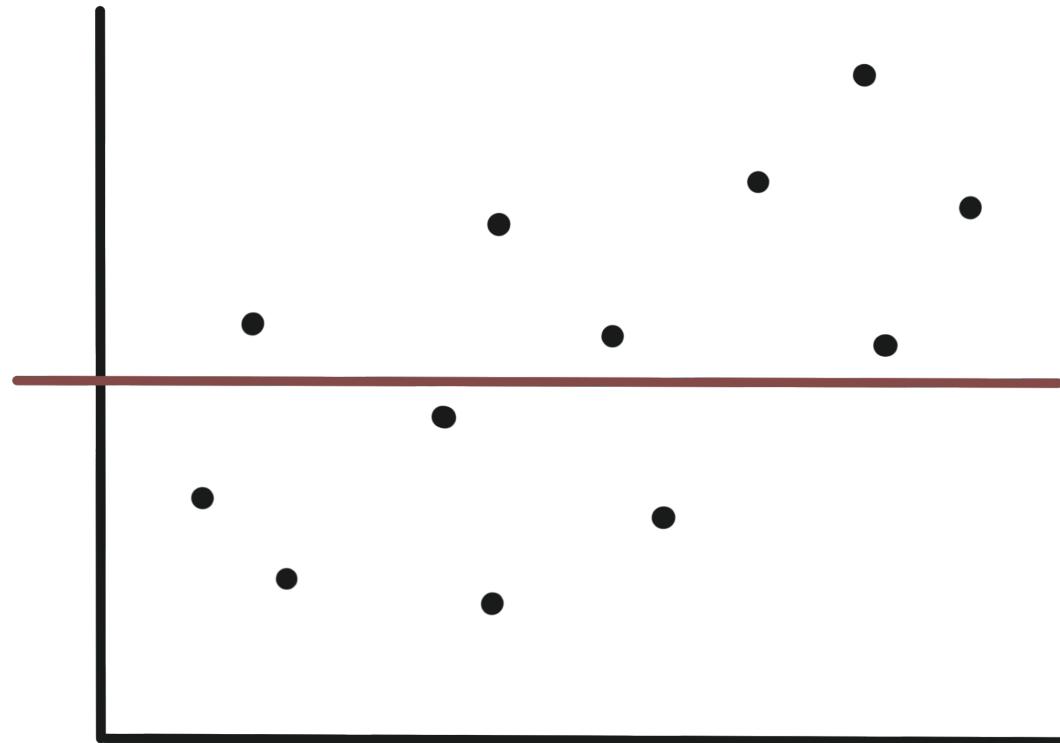
How to Measure Model Success

R^2 :



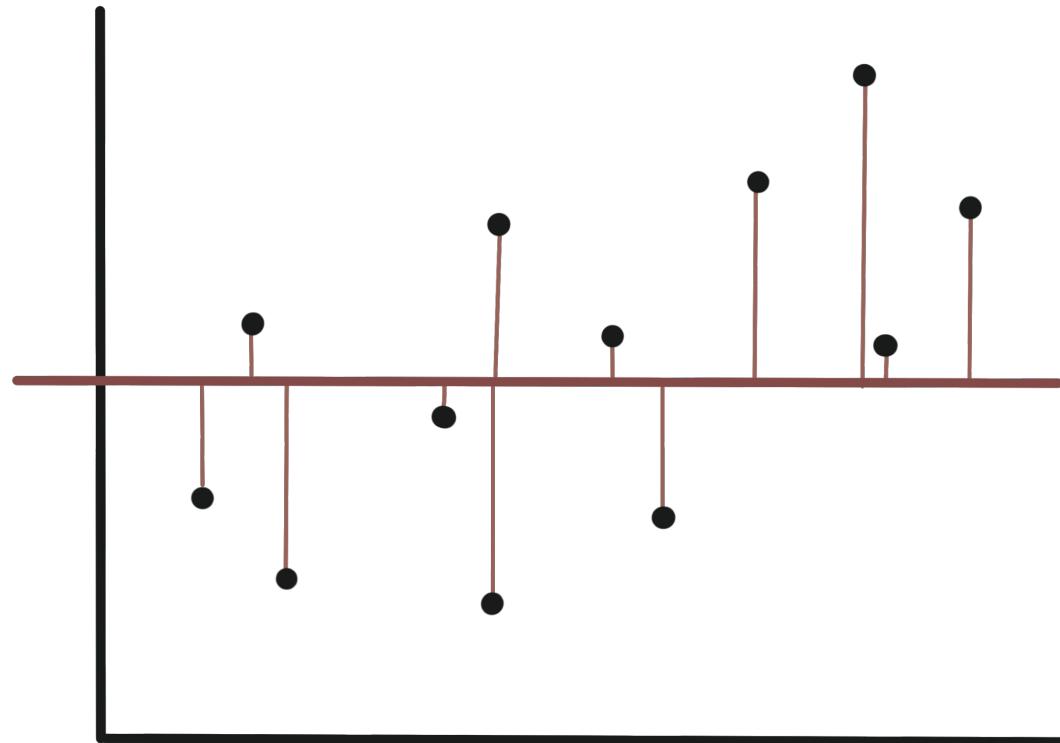
How to Measure Model Success

R^2 :



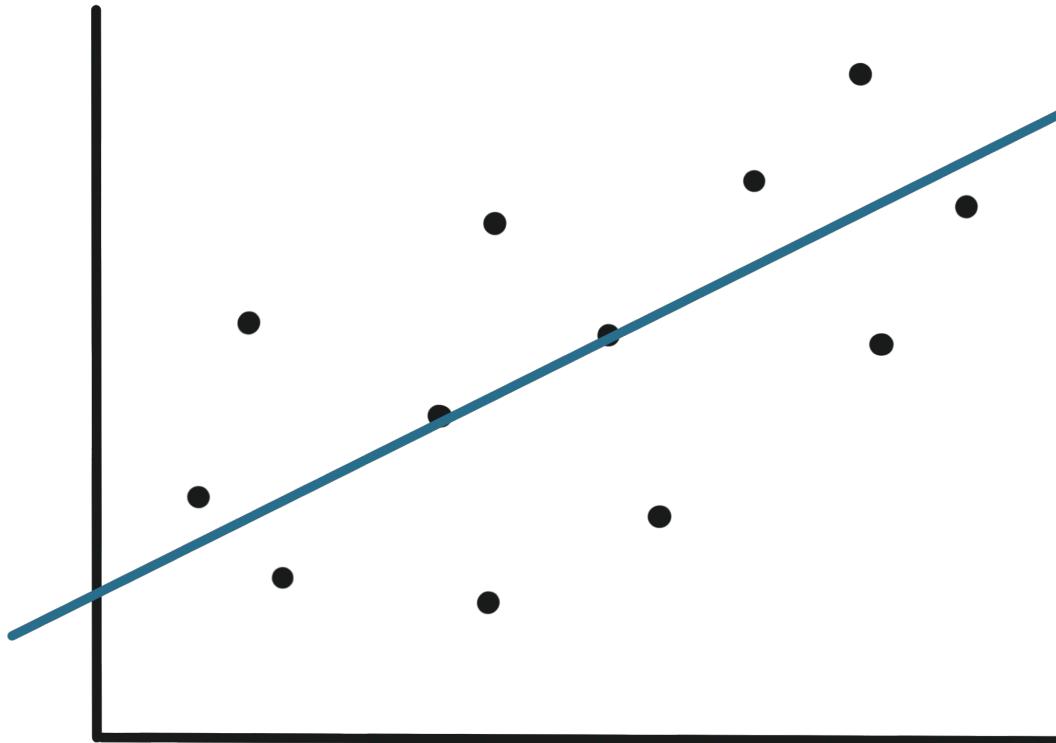
How to Measure Model Success

R^2 :



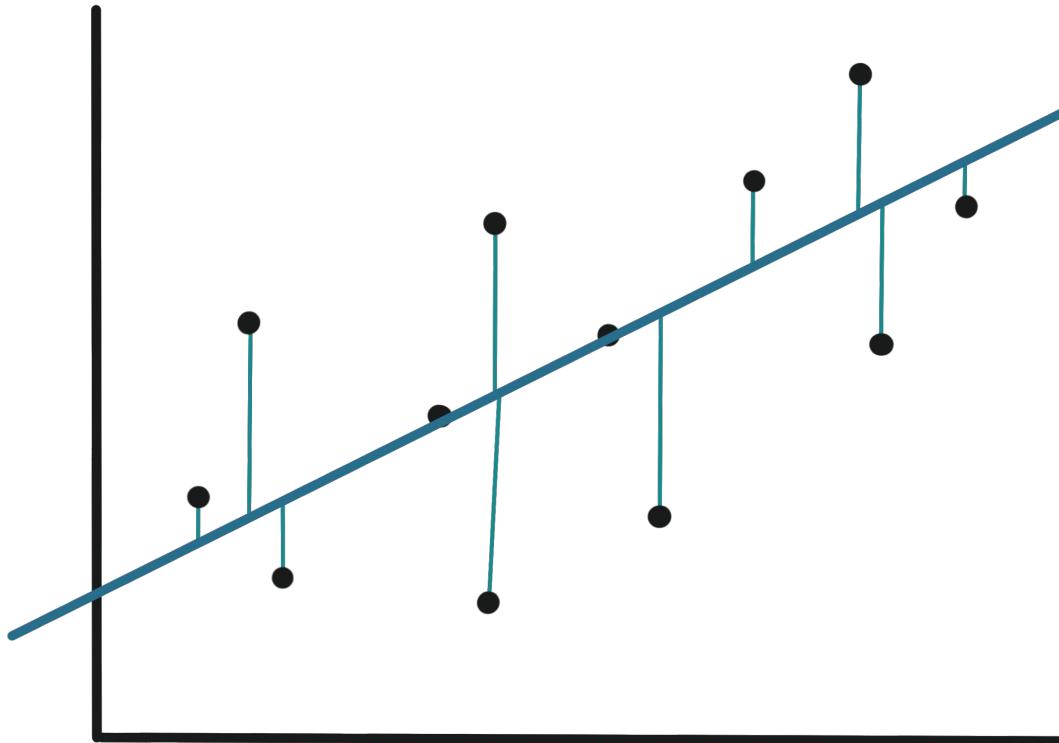
How to Measure Model Success

R^2 :



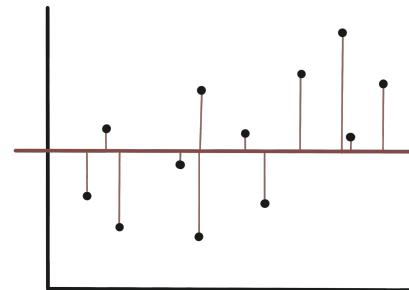
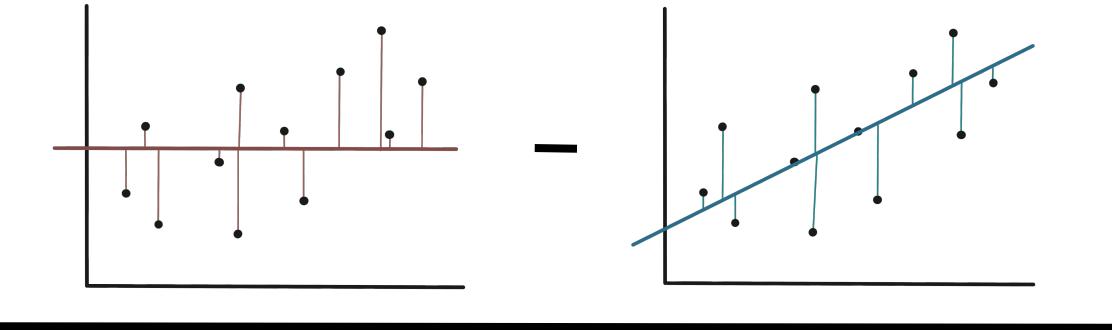
How to Measure Model Success

R^2 :

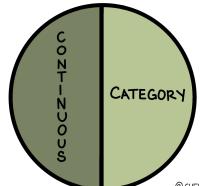


How to Measure Model Success

R^2 :

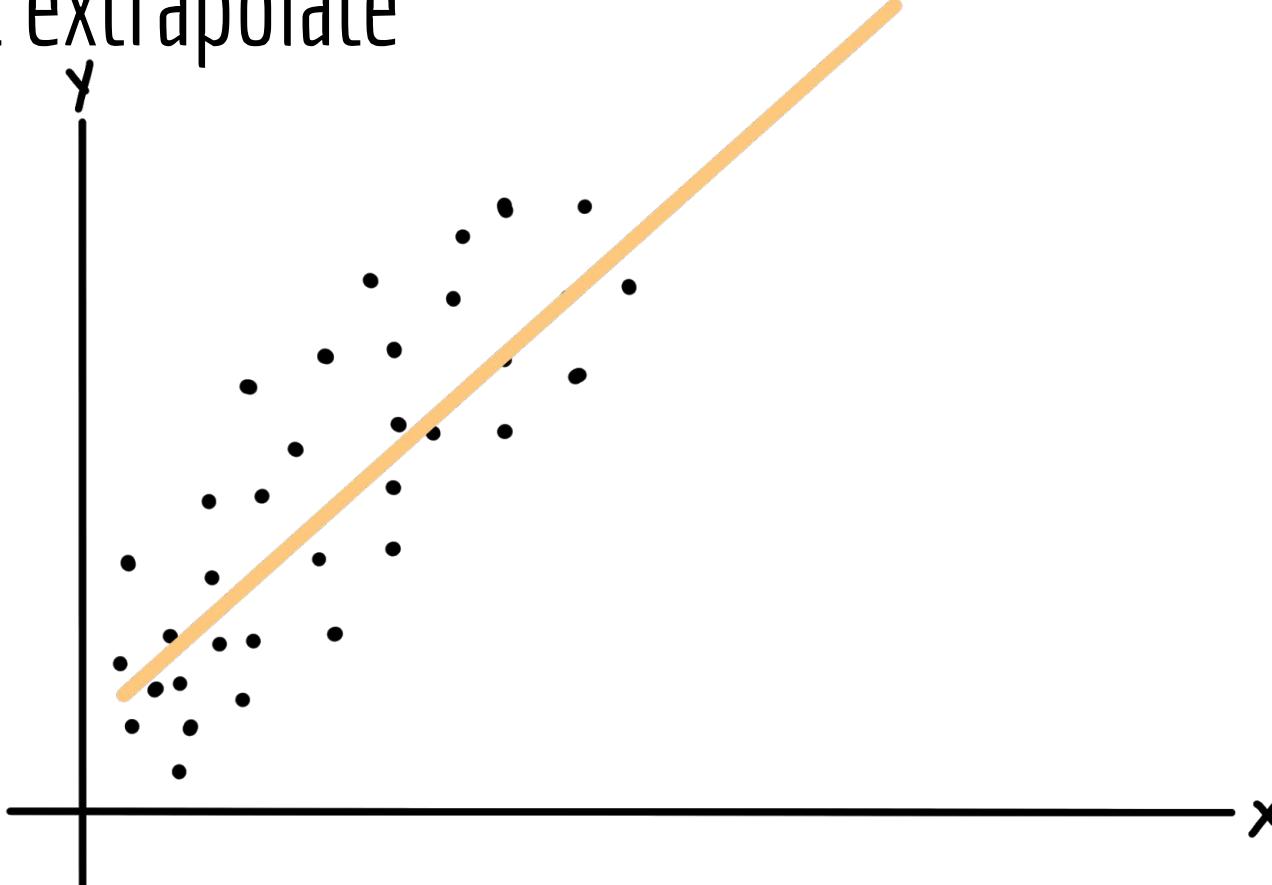


PREDICT

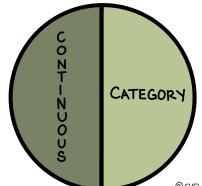


@CHELSEAPARLETT

Don't extrapolate

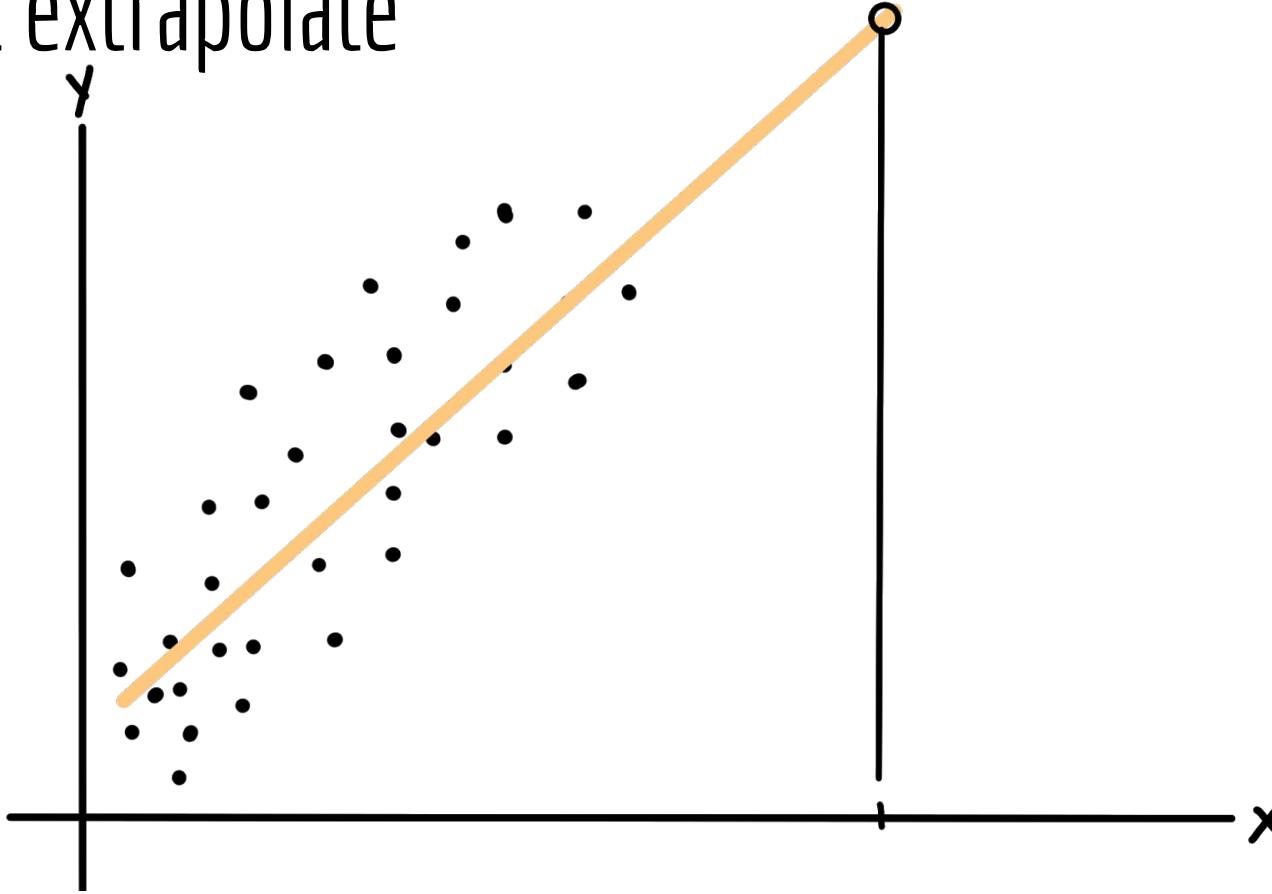


PREDICT

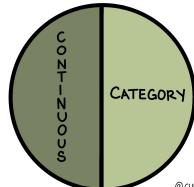


@CHELSEAPARLETT

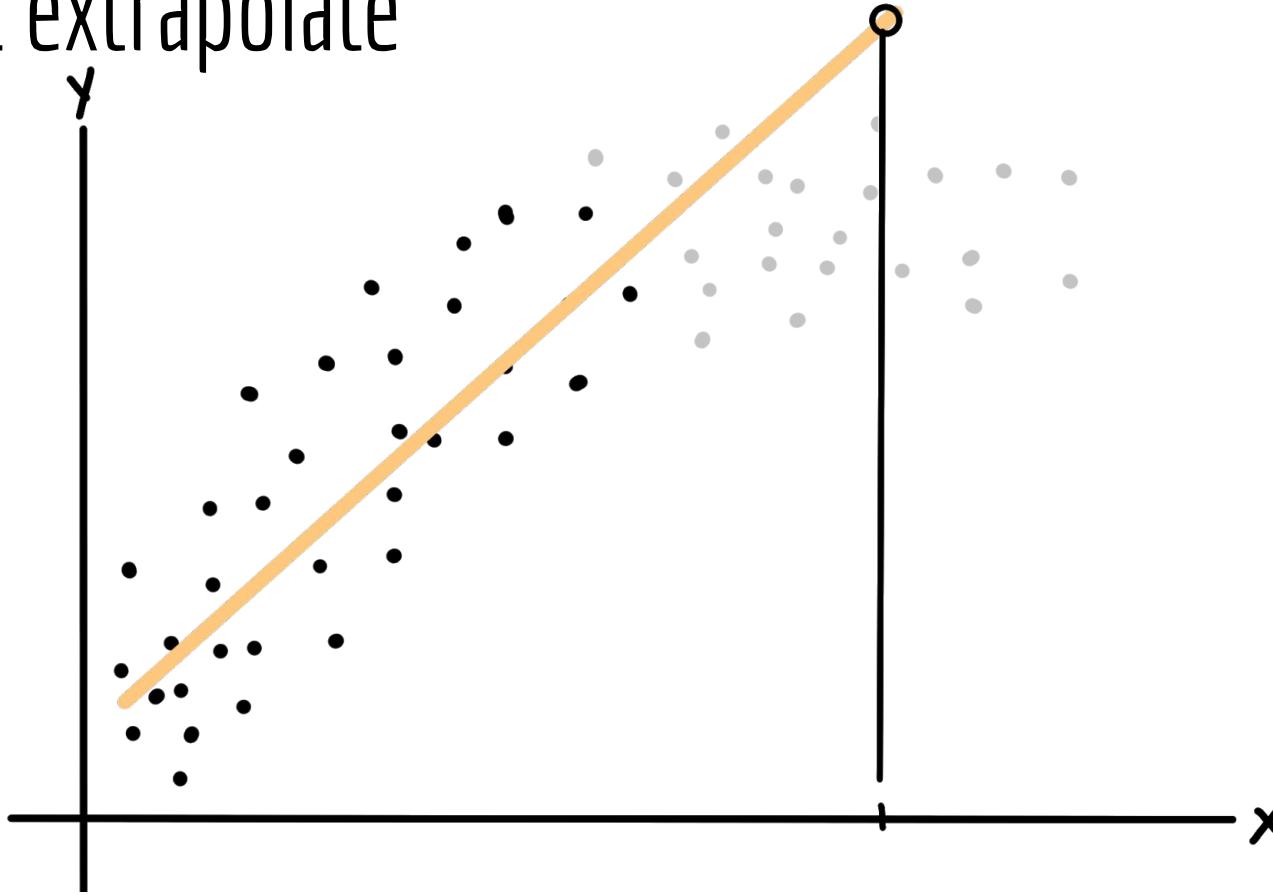
Don't extrapolate



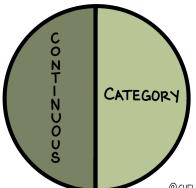
PREDICT



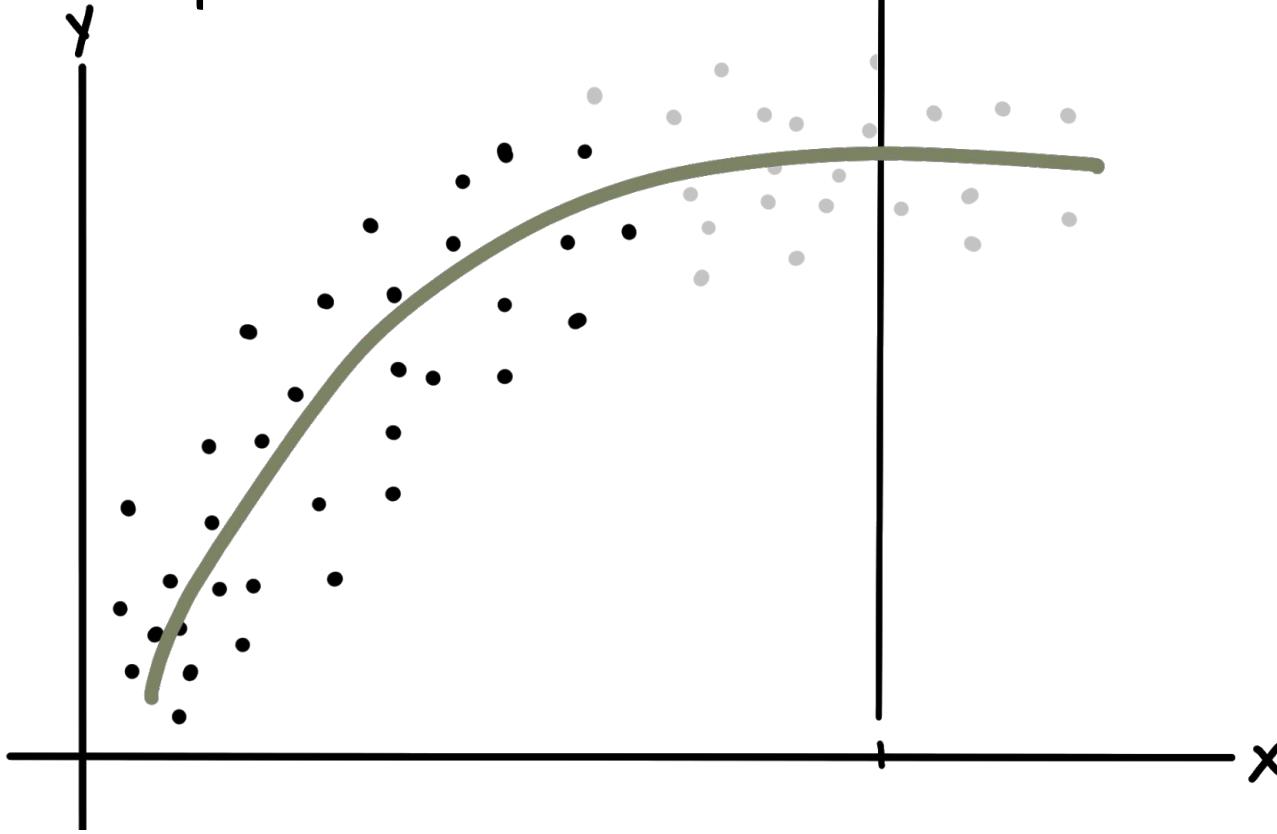
Don't extrapolate



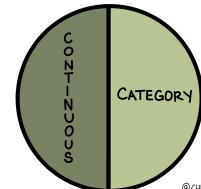
PREDICT



Don't extrapolate



PREDICT



Doing LR for prediction

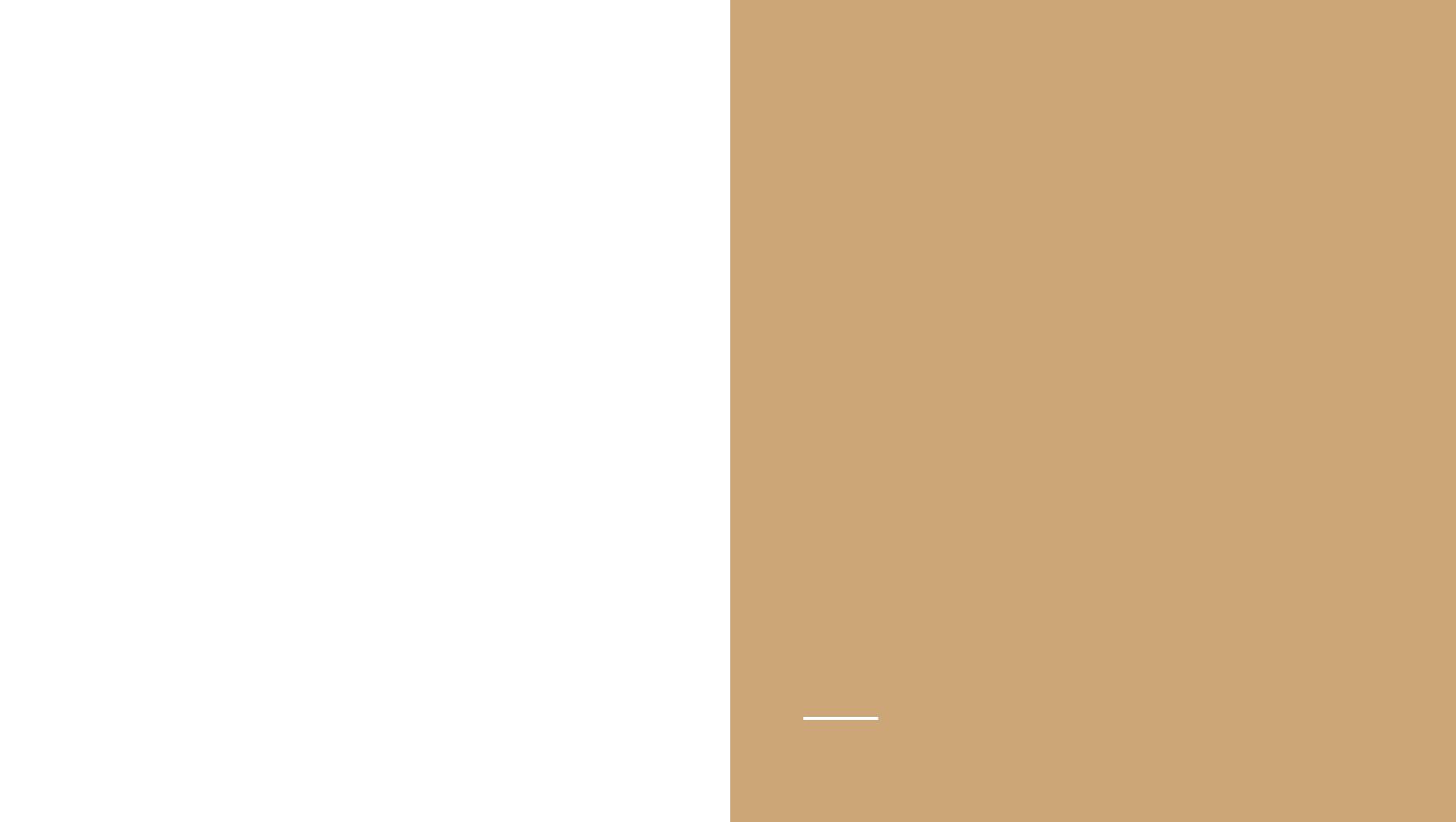
sklearn:

```
x = Model()
```

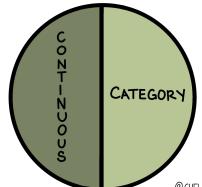
```
x.fit(X, y)
```

```
x.predict(Xnew)
```

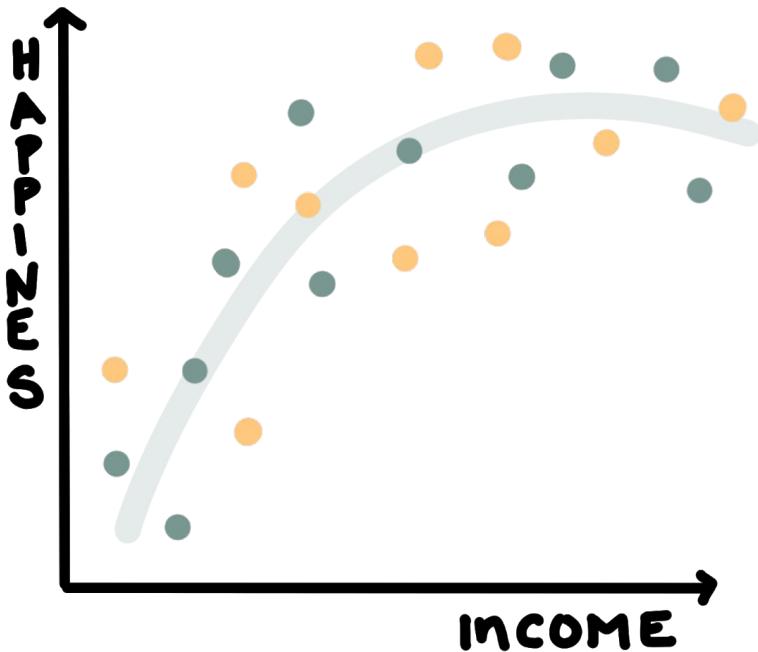
```
x.score(X, y)
```



PREDICT



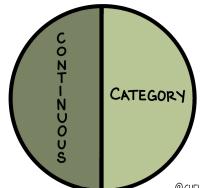
Bias and Variance



- data we have now
- future data

@CHELSEAPARLETT

PREDICT



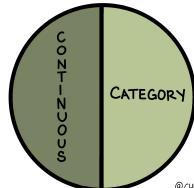
Bias and Variance



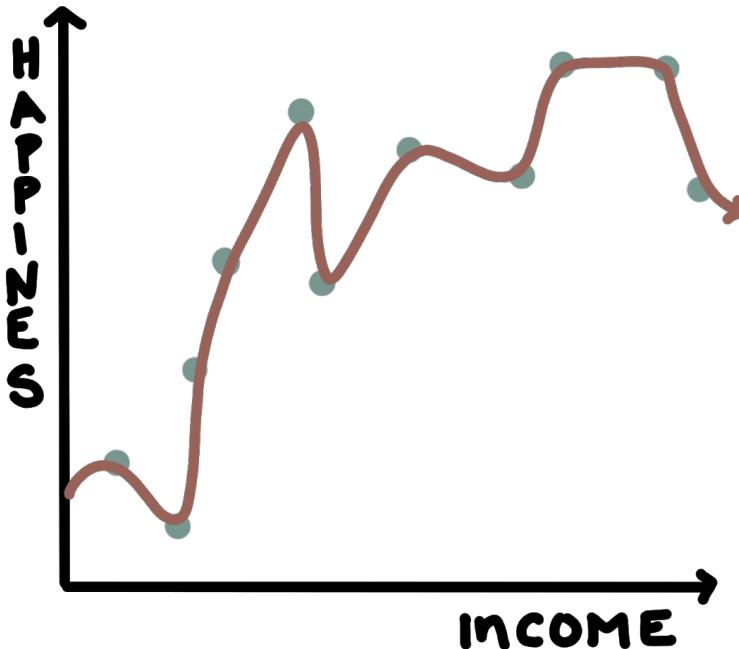
- data we have now
- future data

@CHELSEAPARLETT

PREDICT



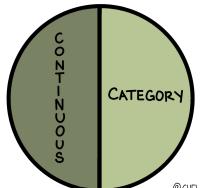
Bias and Variance



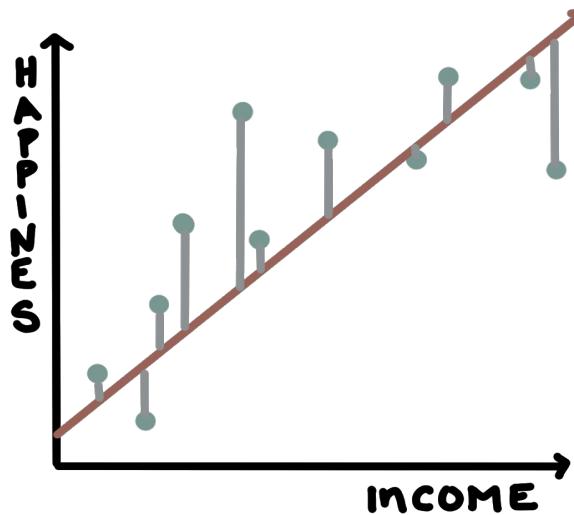
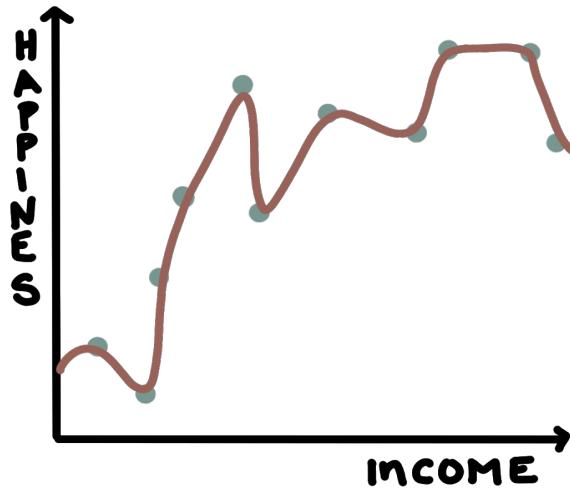
- data we have now
- future data

@CHELSEAPARLETT

PREDICT

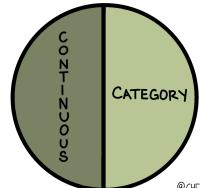


Bias and Variance

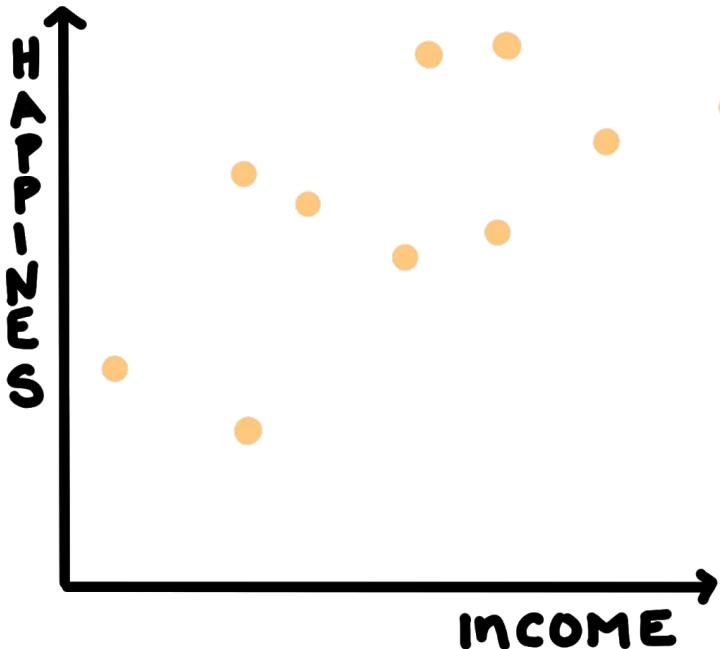


- data we have now
- future data

PREDICT



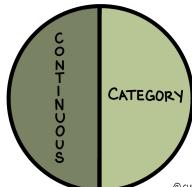
Bias and Variance



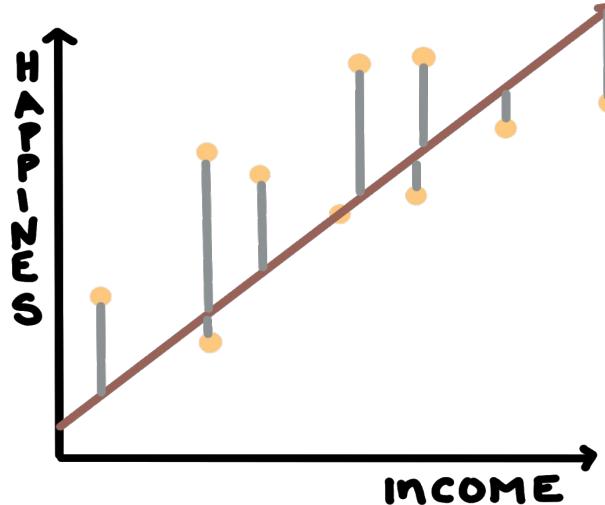
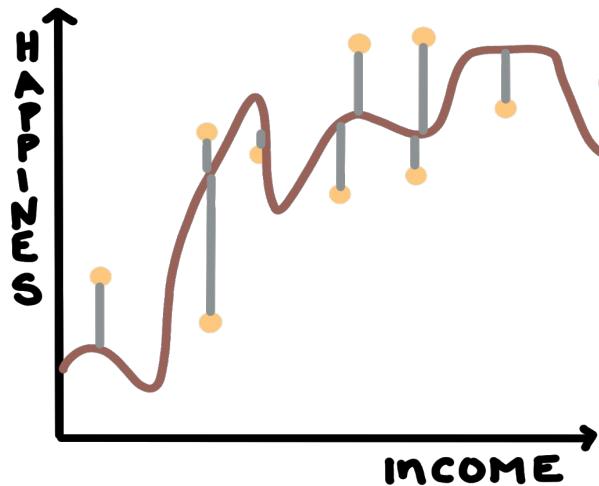
- data we have now
- future data

@CHELSEAPARLETT

PREDICT

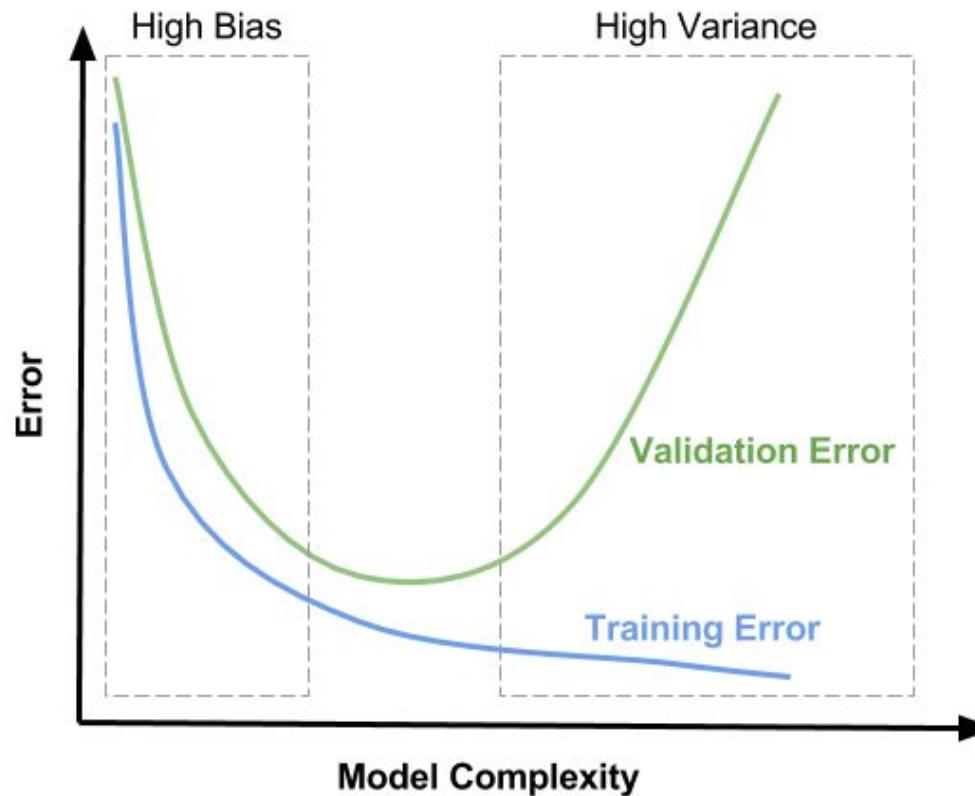


Bias and Variance

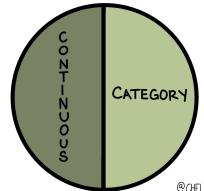


■ data we have now
■ future data

Bias and Variance



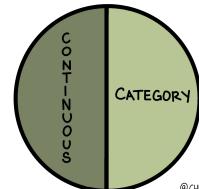
PREDICT



Validation

1. **Split data** into train and test, build model on train, test on test
2. Use performance metrics on test data (or average over multiple) to **evaluate model performance**

PREDICT

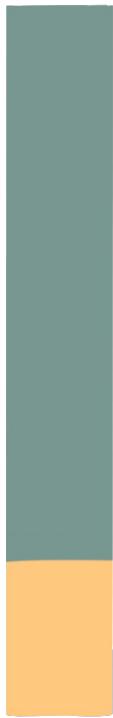


Validation

DATA

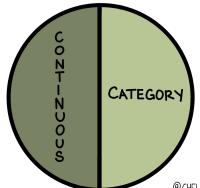
train

test



@CHESEAPARLETT

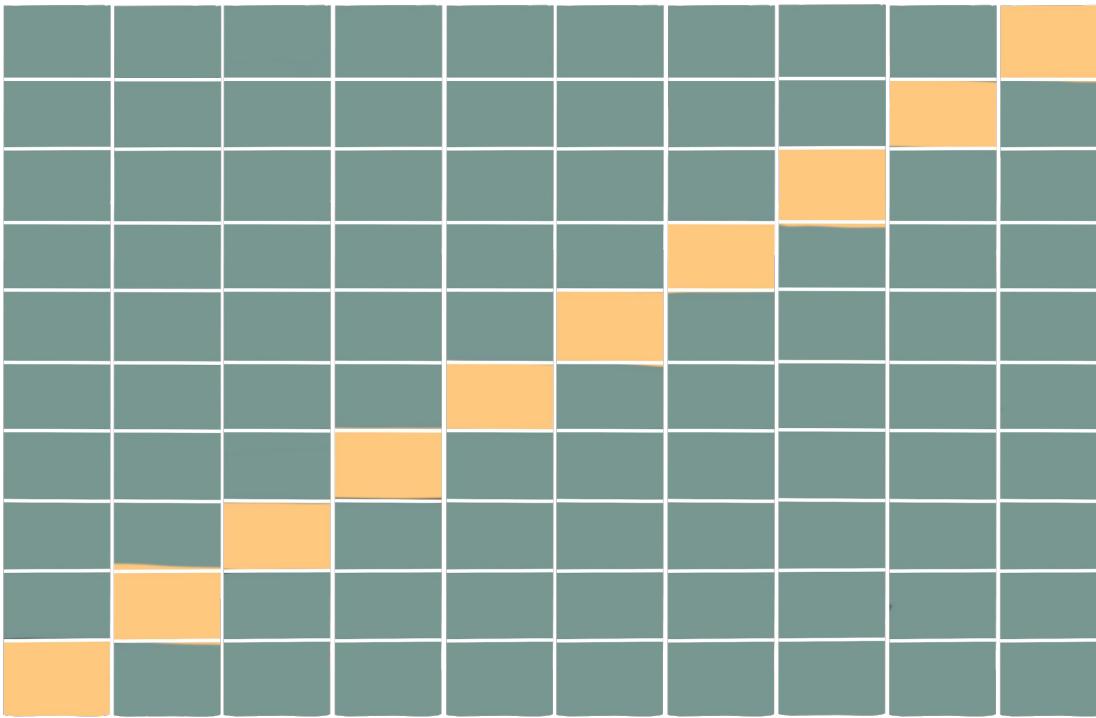
PREDICT



Cross Validation

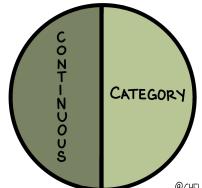
DATA

■ train
■ test



@CHESEAPARLETT

PREDICT



Cross Validation

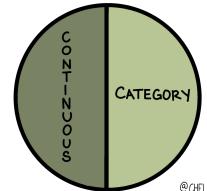
DATA

train
test



@CHESEAPARLETT

PREDICT



Cross Validation

How to decide (Things to Think About)

- Size of your Dataset (rows AND columns)
- Computational Expense