

Regression



Simple Linear Regression





Simple Linear Regression

$$\hat{y} = \underline{b_0} + \underline{b_1 X_1}$$

Dependent variable

y-intercept (constant)

Independent variable

Slope coefficient



Simple Linear Regression



~

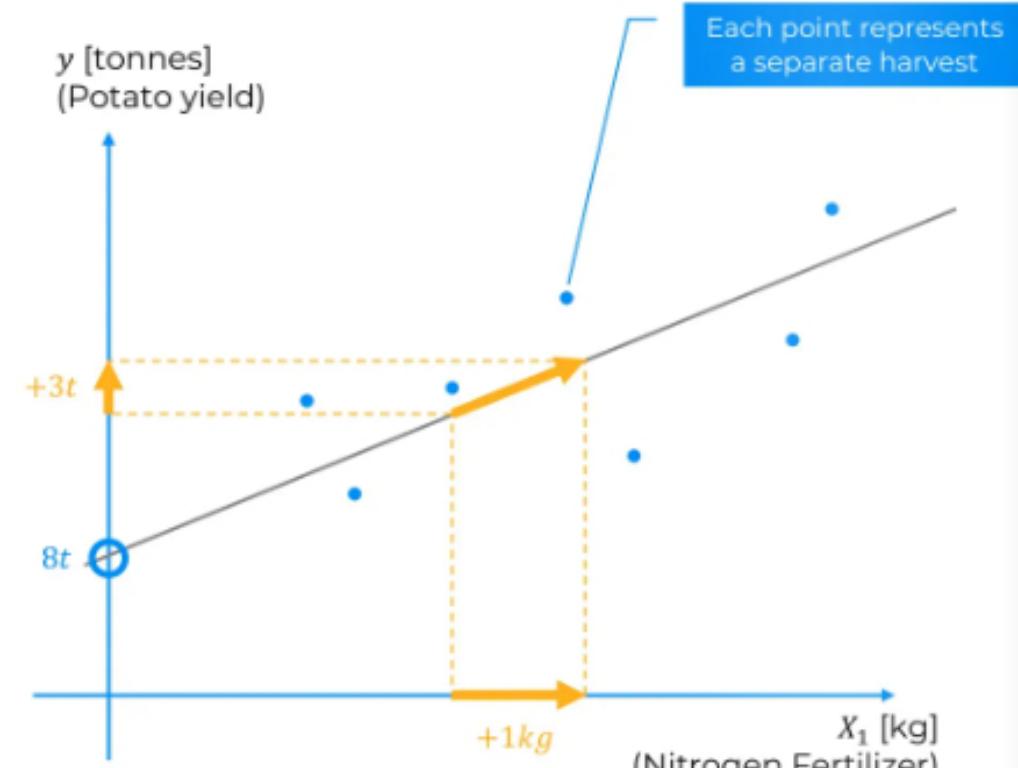


$$\hat{y} = b_0 + b_1 X_1$$

Potatoes[t] = $b_0 + b_1 \times$ Fertilizer[kg]

$$b_0 = 8[t]$$

$$b_1 = 3\left[\frac{t}{kg}\right]$$



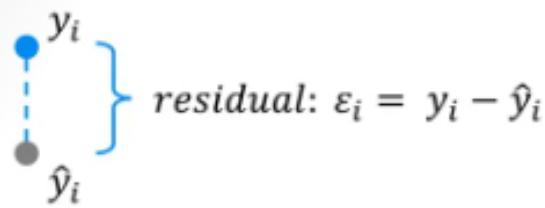
Ordinary Least Squares



Simple Linear Regression



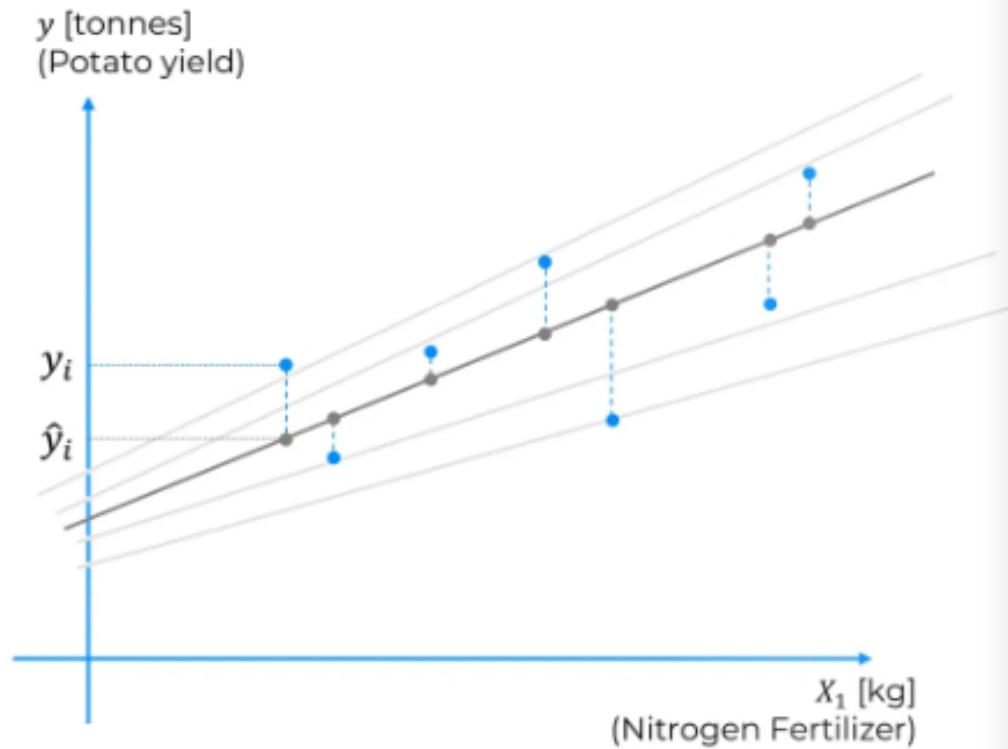
Ordinary Least Squares:



$$\hat{y} = b_0 + b_1 X_1$$

b_0, b_1 such that:

$SUM(y_i - \hat{y}_i)^2$ is minimized



Multiple Linear Regression



Multiple Linear Regression



$$\hat{y} = b_0 + b_1X_1 + b_2X_2 + \cdots + b_nX_n$$

Dependent variable
y-intercept (constant)

Independent variable 1
Slope coefficient 1

Independent variable 2
Slope coefficient 2

Independent variable n
Slope coefficient n



Multiple Linear Regression



~



$$Potatoes[t] = 8t + 3 \frac{t}{kg} \times Fertilizer[kg] - 0.54 \frac{t}{^{\circ}C} \times AvgTemp[{}^{\circ}C] + 0.04 \frac{t}{mm} \times Rain[mm]$$



Additional Reading



The Application of Multiple Linear Regression and Artificial Neural Network Models for Yield Prediction of Very Early Potato Cultivars before Harvest

Magdalena Piekutowska et. al. (2021)

Link:

<https://www.mdpi.com/2073-4395/11/5/885>

Quantitative Yield Forecast		
Models RY1 and NY1	Yield Forecast before Harvest (40 Days from Full Emergence)	Data Range
INSO	insolation sum [h] in the periods: planting—June 20,	275.3–711.7
TEMP	average daily air temperature [$^{\circ}\text{C}$] in the periods: planting—20 June	10.8–15.7
PREC	precipitation [mm] in the periods: planting—20 June	38.7–258.2
NITRO	sum of nitrogen fertilization [kg ha^{-1}] in the periods planting—20 June	80–155
PHOSP	sum of phosphorus fertilization [kg ha^{-1}]	28.2–150
POTAS	sum of potassium fertilization [kg ha^{-1}]	80–306.5
PLANT	planting date [number of days since the beginning of the year]	107–127
EMERG	date of emergence [number of days since the beginning of the year], yield forecast 20th of June	130–151
DENST	densification [plants/plot], yield forecast June 20	35–60
PH	Soil pH [in 1 mol KCl]	5.8–7
SFERTP	soil fertility in phosphorus [mg $\text{P}_2\text{O}_5 \cdot 100 \text{ g}^{-1}$ soil]	14–26.2
SPERTK	soil fertility in potassium [mg $\text{K}_2\text{O} \cdot 100 \text{ g}^{-1}$ soil]	11.7–19.2
SFERTM	soil fertility in magnesium [mg $\text{Mg} \cdot 100 \text{ g}^{-1}$ soil]	3–9.1
YIELDPI	tuber yield [$\text{t} \cdot \text{ha}^{-1}$], harvest 40 days from full emergence	11.6–41.3

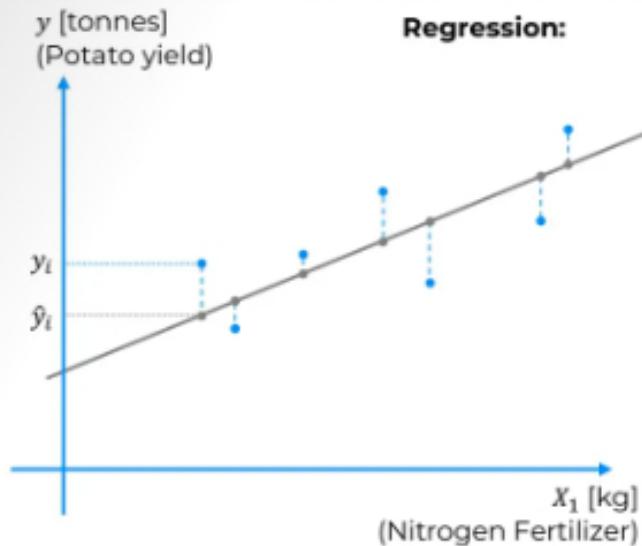


R Squared

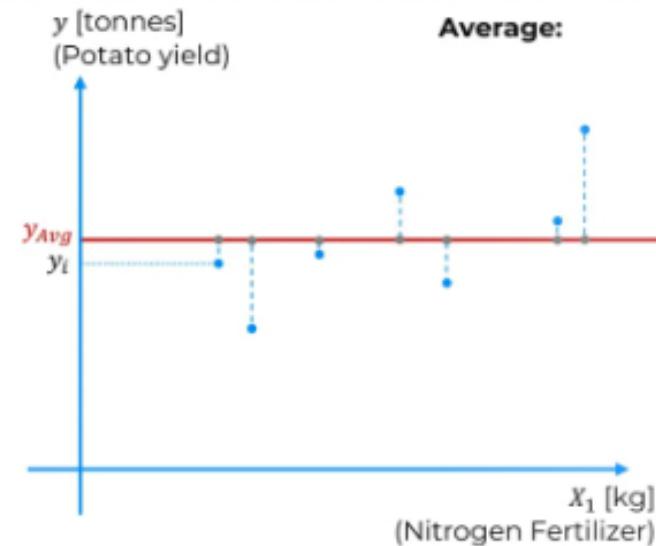




R Squared



$$SS_{res} = \text{SUM}(y_i - \hat{y}_i)^2$$



$$SS_{tot} = \text{SUM}(y_i - y_{avg})^2$$

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

Rule of thumb (for our tutorials)*:

- 1.0 = Perfect fit (suspicious)
- ~0.9 = Very good
- <0.7 = Not great
- <0.4 = Terrible
- <0 = Model makes no sense for this data

*This is highly dependent on the context



Adjusted R Squared





Adjusted R Squared

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

R² – Goodness of fit
(greater is better)

Problem:

$$\hat{y} = b_0 + b_1X_1 + b_2X_2 + b_3X_3$$

$$SS_{res} = \text{SUM}(y_i - \hat{y}_i)^2$$

SS_{tot} doesn't change

SS_{res} will decrease or stay the same

(This is because of Ordinary Least Squares: SS_{res} → Min)

Solution:

$$Adj\ R^2 = 1 - (1 - R^2) \times \frac{n - 1}{n - k - 1}$$

k – number of independent variables

n – sample size

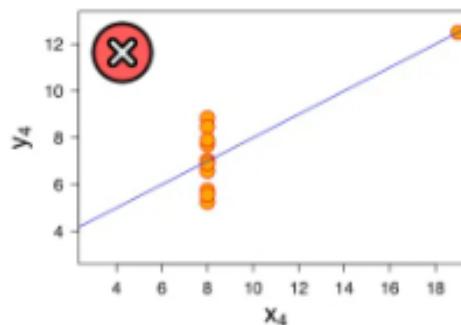
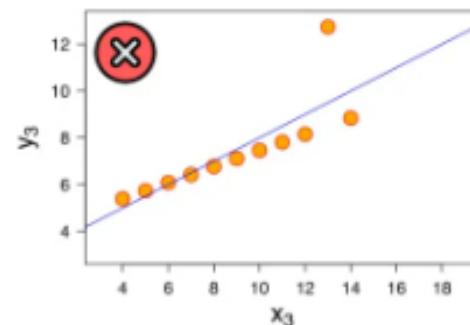
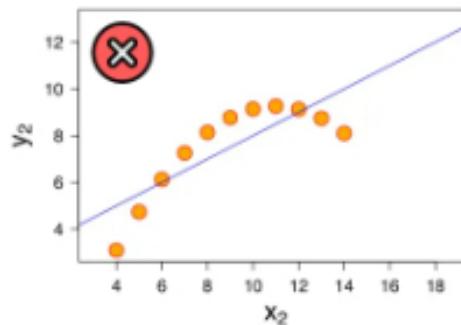
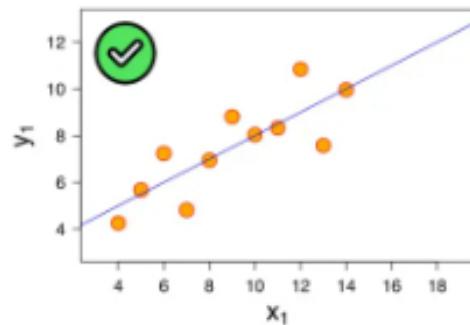


Assumptions Of Linear Regression

Assumptions of Linear Regression



Anscombe's quartet (1973):

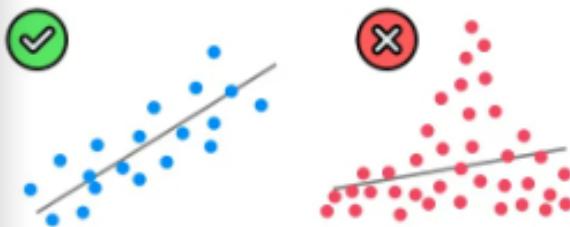




Assumptions of Linear Regression

1. Linearity

(Linear relationship between Y and each X)



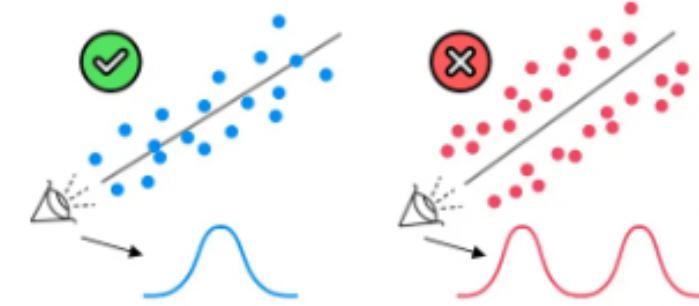
2. Homoscedasticity

(Equal variance)



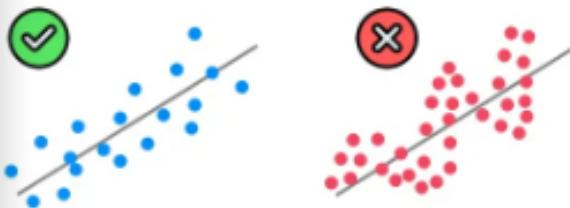
3. Multivariate Normality

(Normality of error distribution)



4. Independence

(of observations. Includes "no autocorrelation")



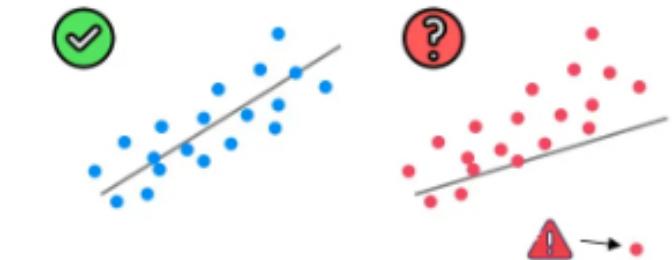
5. Lack of Multicollinearity

(Predictors are not correlated with each other)

$$\checkmark X_1 \not\sim X_2 \quad \times X_1 \sim X_2$$

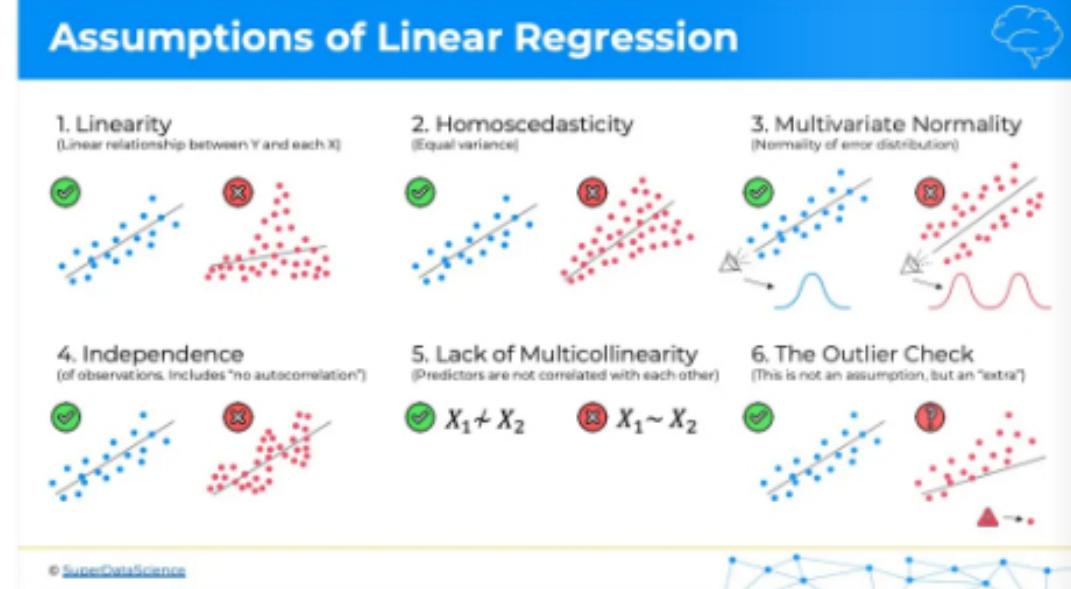
6. The Outlier Check

(This is not an assumption, but an "extra")



Bonus

Download the Assumptions poster at:
superdatascience.com/assumptions





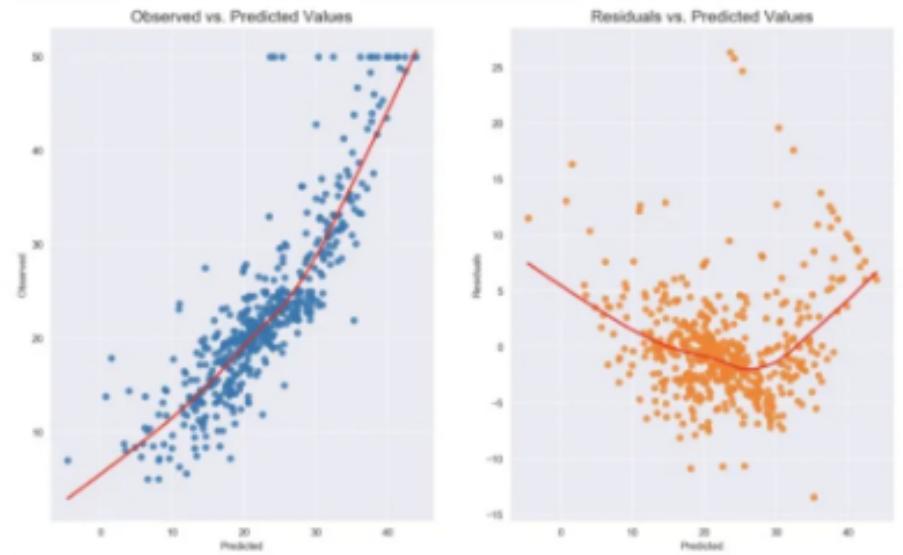
Additional Reading

*Verifying the Assumptions of Linear Regression
in Python and R*

Eryk Lewinson (2019)

Link:

towardsdatascience.com/verifying-the-assumptions-of-linear-regression-in-python-and-r-f4cd2907d4c0



Dummy Variables

Dummy Variables

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + ???$$

Dummy Variables

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1$$



Dummy Var. Trap

Dummy Variable Trap

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1$$

Dummy Variable Trap

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70			California
191,050.39	153,441.51			California
182,901.99	144,372.41			New York
166,187.94	142,107.34			California

Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$D_2 = 1 - D_1$$

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1 + \underline{b_5 * D_2}$$



Dummy Variable Trap

Profit	R&D Spend	Admin	Marketing	State	Dummy Variables	
					New York	California
192,261.83	165,349.20	136,897.80	471,784.10	New York	1	0
191,792.06	162,597.70	151,377.59	443,898.53	California	0	1
191,050.39	153,441.51	101,145.55	407,934.54	California	0	1
182,901.99	144,372.41	118,671.85	383,199.62	New York	1	0
166,187.94	142,107.34	91,391.77	366,168.42	California	0	1

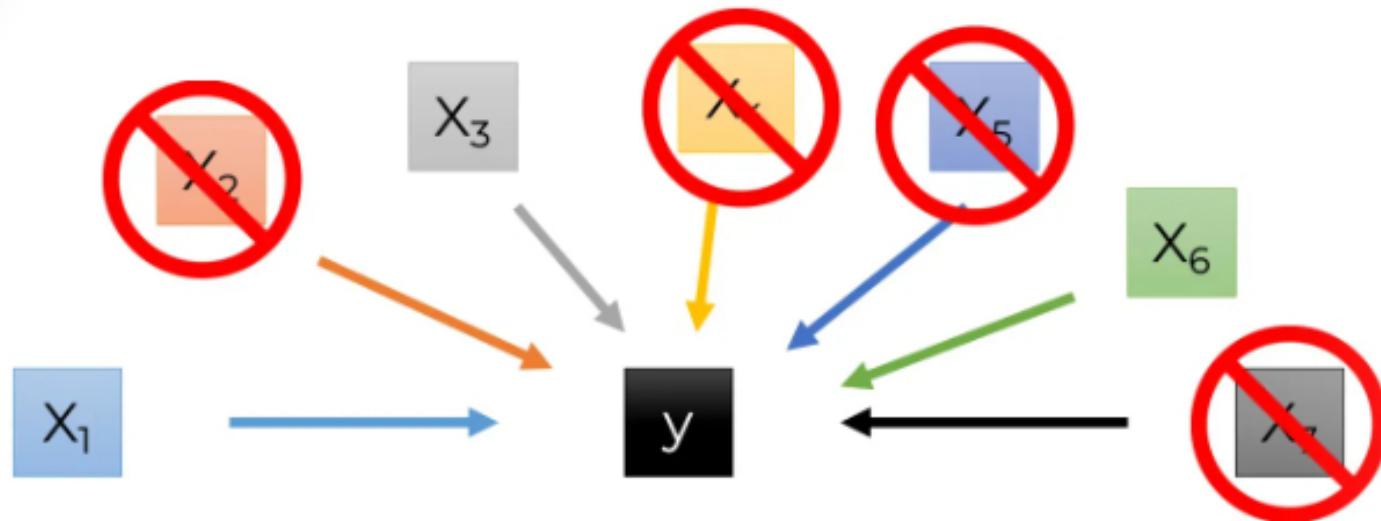
$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1 + b_5 * D_2$$



Always omit one dummy variable

Building A Model (Step-By-Step)

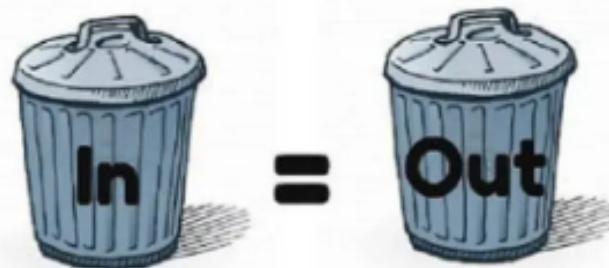
Building A Model



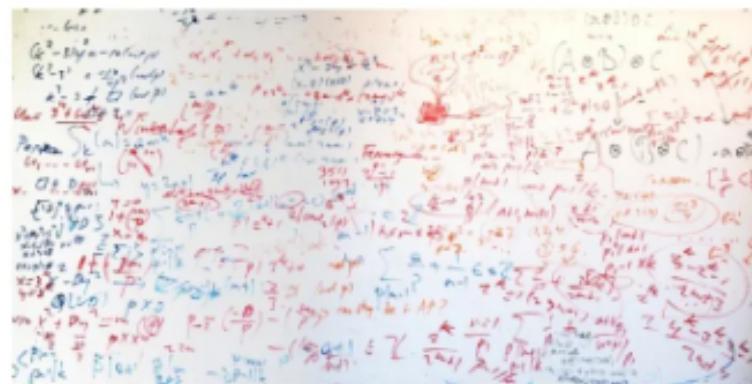
Why?

Building A Model

1)



2)



Building A Model

5 methods of building models:

1. All-in
2. Backward Elimination
3. Forward Selection
4. Bidirectional Elimination
5. Score Comparison

} Stepwise
Regression

Building A Model

“All-in” – cases:

- Prior knowledge; OR
- You have to; OR
- Preparing for Backward Elimination



Building A Model

Backward Elimination

STEP 1: Select a significance level to stay in the model (e.g. SL = 0.05)



STEP 2: Fit the full model with all possible predictors



STEP 3: Consider the predictor with the highest P-value. If $P > SL$, go to STEP 4, otherwise go to FIN



STEP 4: Remove the predictor



STEP 5: Fit model without this variable*



FIN: Your Model Is Ready

Building A Model

Forward Selection

STEP 1: Select a significance level to enter the model (e.g. SL = 0.05)



STEP 2: Fit all simple regression models $y \sim x_n$. Select the one with the lowest P-value



STEP 3: Keep this variable and fit all possible models with one extra predictor added to the one(s) you already have



STEP 4: Consider the predictor with the lowest P-value. If $P < SL$, go to STEP 3, otherwise go to FIN

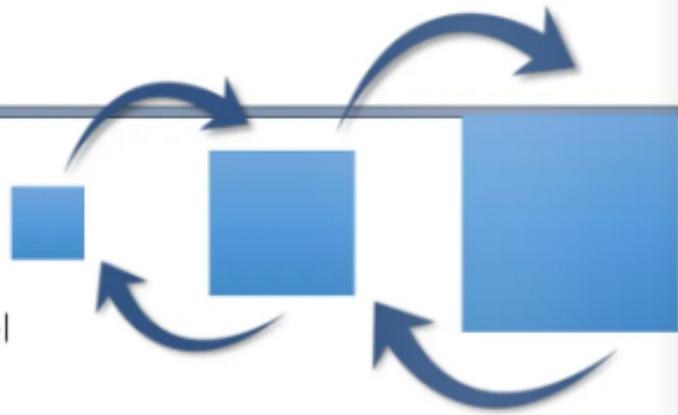


FIN: Keep the previous model

Building A Model

Bidirectional Elimination

STEP 1: Select a significance level to enter and to stay in the model
e.g.: SLENTER = 0.05, SLSTAY = 0.05



STEP 2: Perform the next step of Forward Selection (new variables must have: $P < \text{SLENTER}$ to enter)

STEP 3: Perform ALL steps of Backward Elimination (old variables must have $P < \text{SLSTAY}$ to stay)

STEP 4: No new variables can enter and no old variables can exit



FIN: Your Model Is Ready

Building A Model

All Possible Models

STEP 1: Select a criterion of goodness of fit (e.g. Akaike criterion)



STEP 2: Construct All Possible Regression Models: $2^N - 1$ total combinations



STEP 3: Select the one with the best criterion



FIN: Your Model Is Ready



Example:
10 columns means
1,023 models

Building A Model

5 methods of building models:

1. All-in
2. Backward Elimination
3. Forward Selection
4. Bidirectional Elimination
5. Score Comparison

Section Recap

Section Recap

In this section we learned:

1. How to create dummies for categorical IVs
2. How to avoid the dummy variable trap
3. Backward, Forward, Bidirectional, All Possible
4. We actually built a model. Step-By-Step!!
5. How to use adjusted R-squared in modelling
6. How to interpret coefficients of a MLR

Polynomial Regression

Regressions

Simple
Linear
Regression

$$y = b_0 + b_1 x_1$$

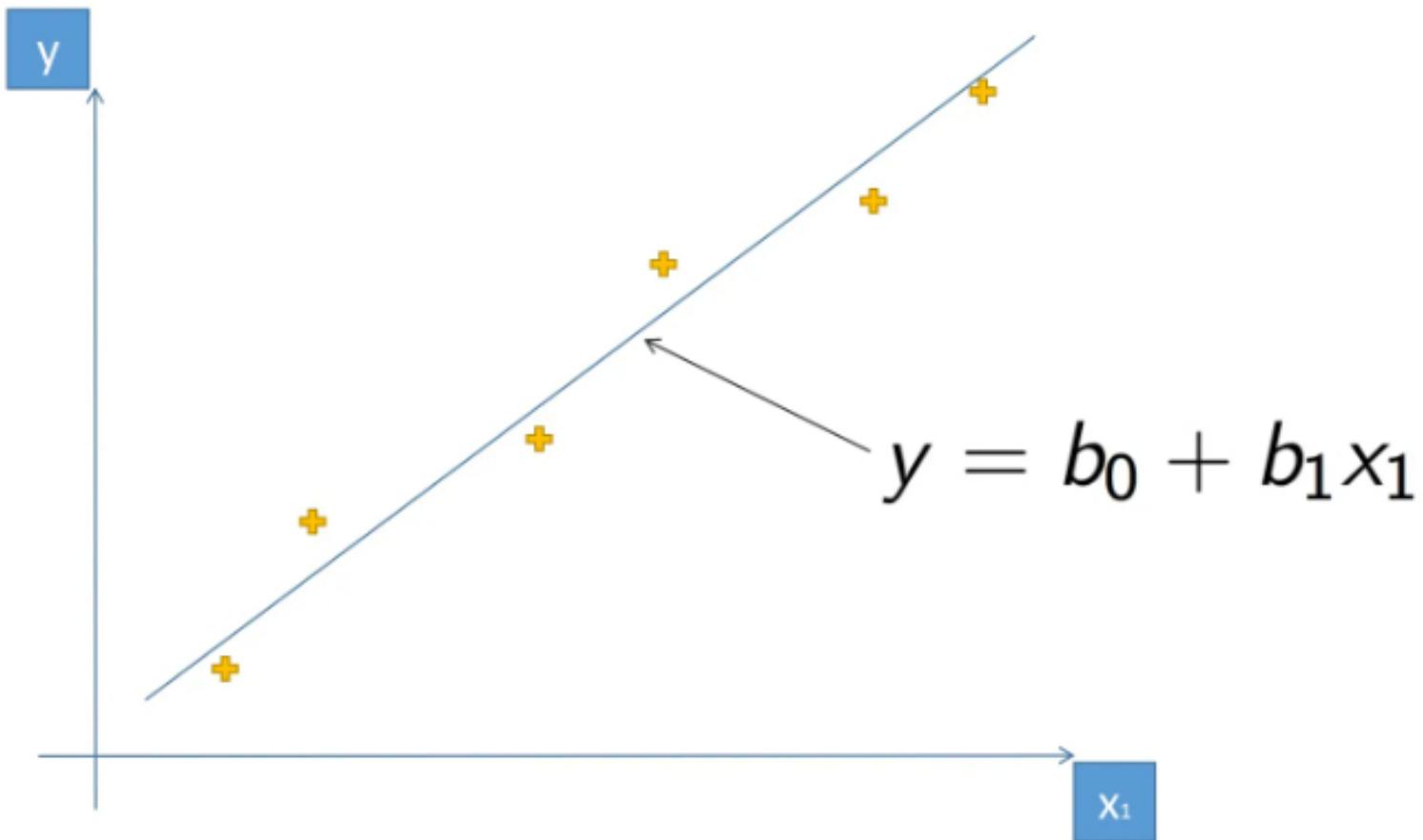
Multiple
Linear
Regression

$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

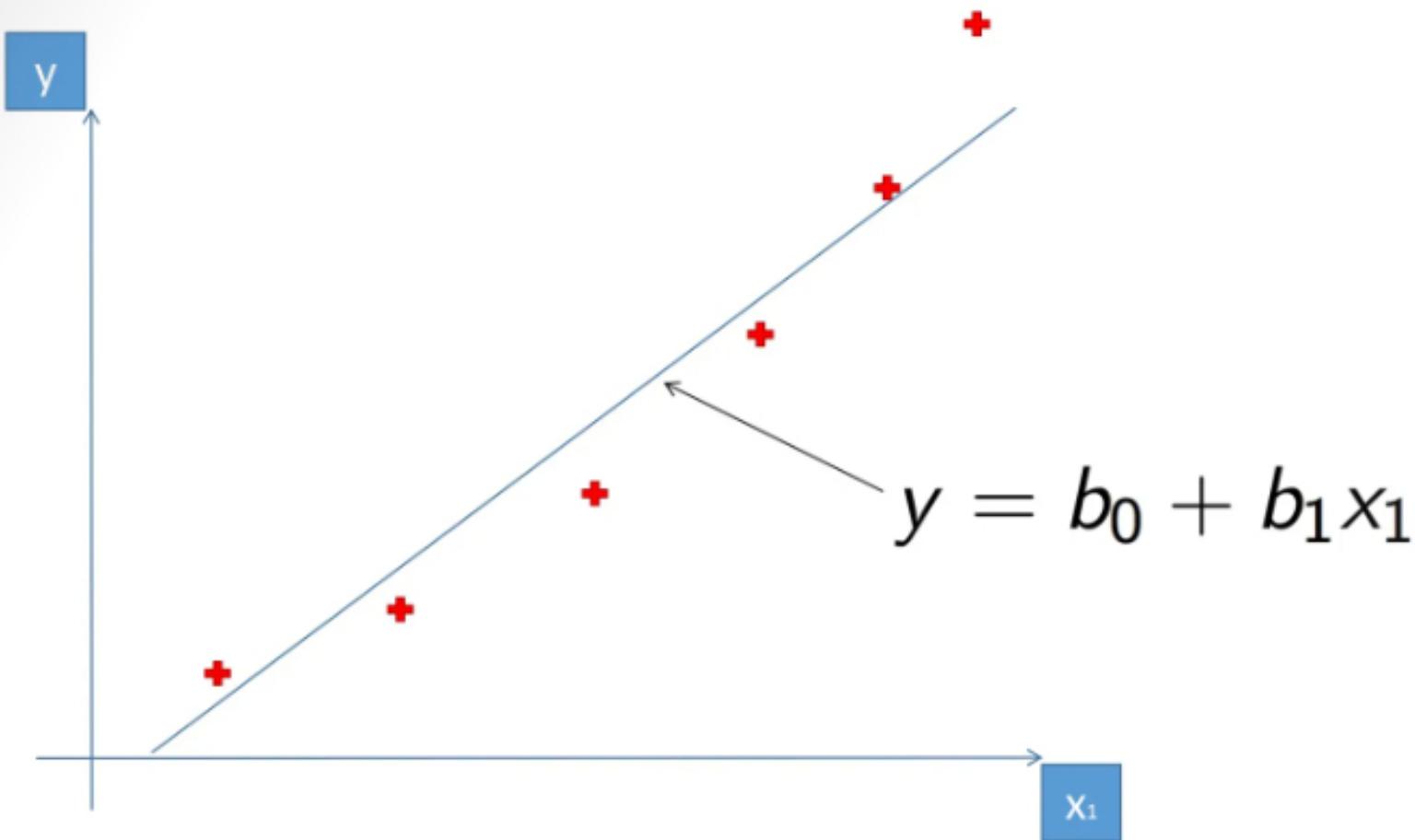
Polynomial
Linear
Regression

$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \dots + b_n x_1^n$$

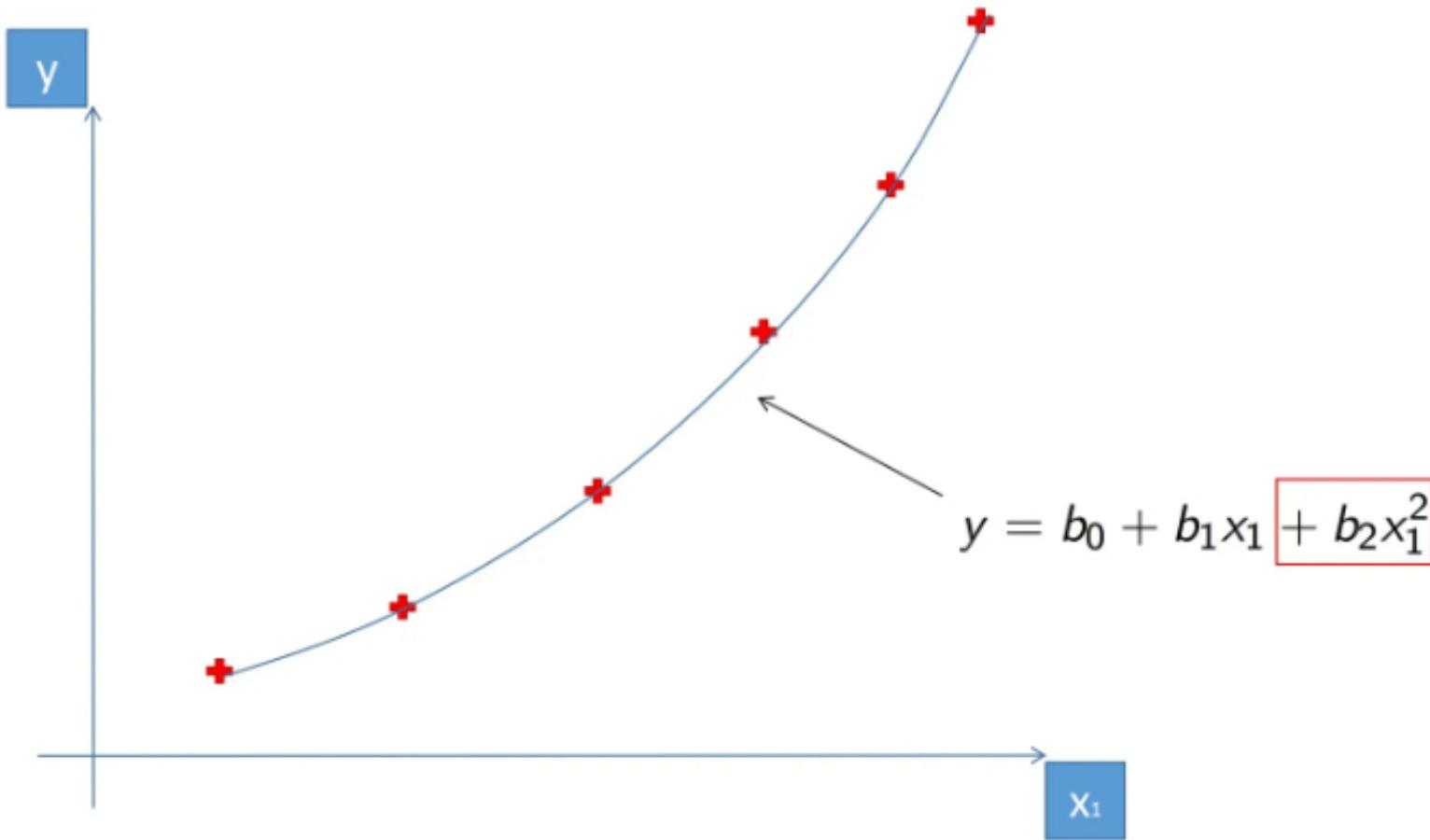
Simple Linear Regression



Simple Linear Regression



Polynomial Regression



Polynomial Regression

One Question: Why “Linear”?

Polynomial Regression

Polynomial
Linear
Regression

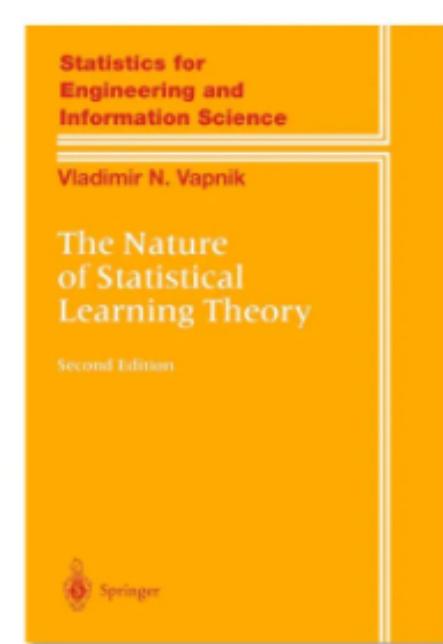
$$y = b_0 + b_1x_1 + b_2x_1^2 + \dots + b_nx_1^n$$

SVR Intuition

SVR Intuition



Vladimir Vapnik



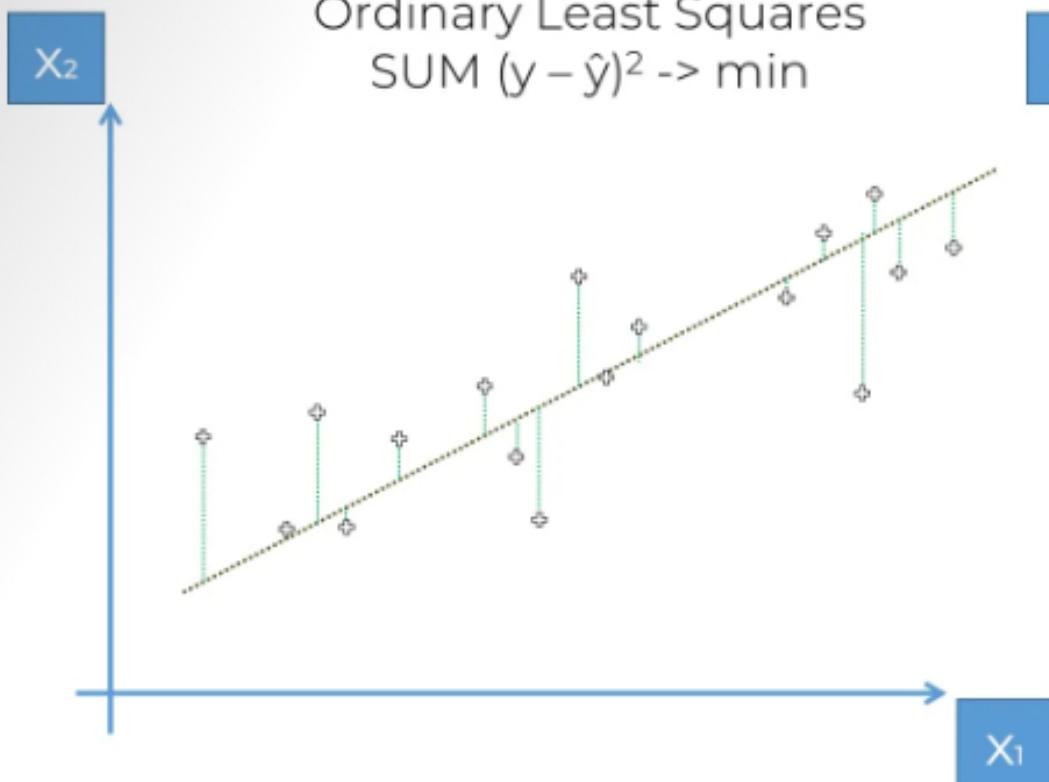
1992

SVR Intuition

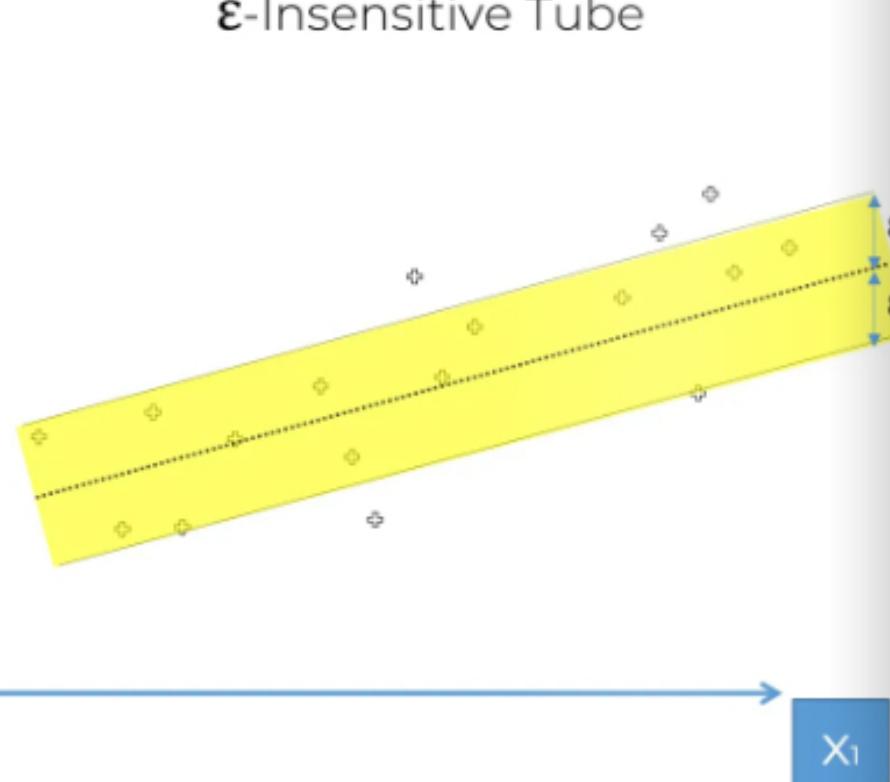
NOT FOR
DISTRIBUTION © SUPERDATASCIENCE

www.superdatascience.com

Ordinary Least Squares
 $\text{SUM } (y - \hat{y})^2 \rightarrow \min$



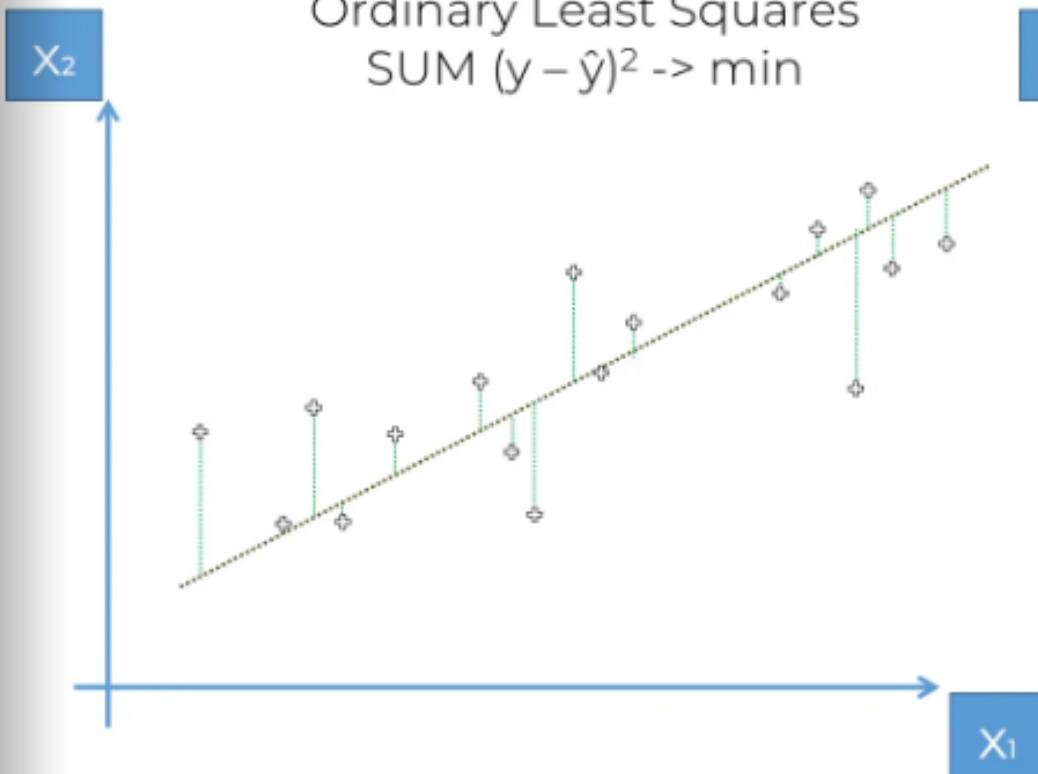
ϵ -Insensitive Tube



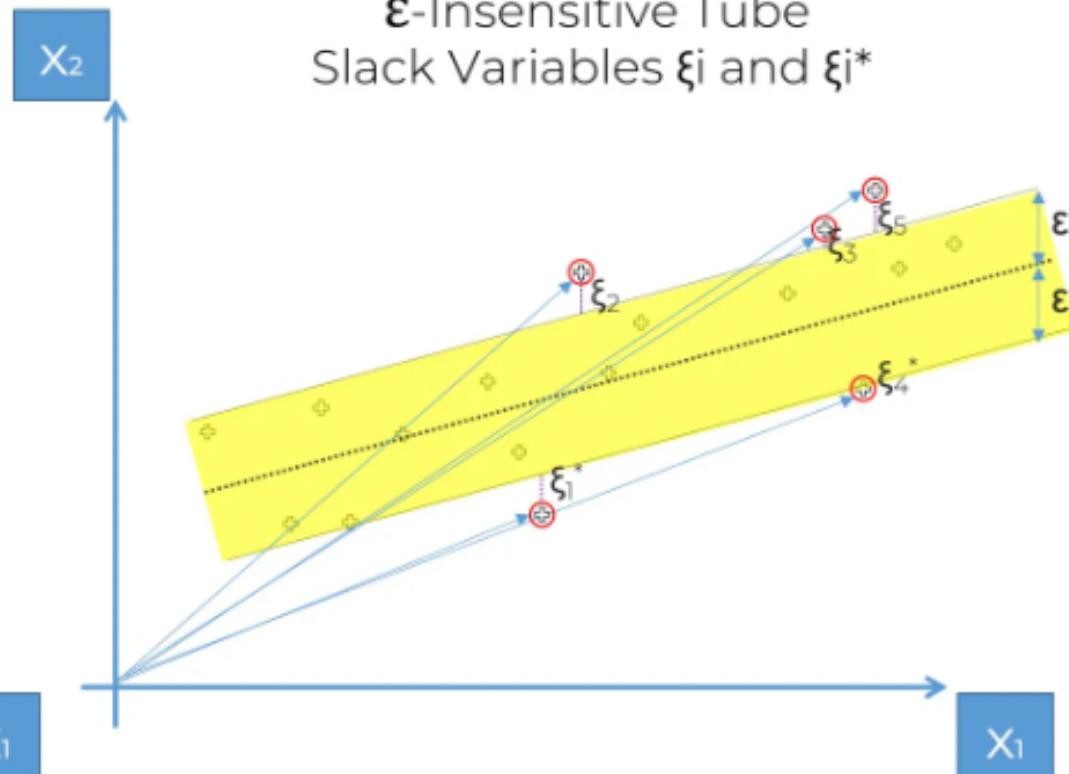
SVR Intuition

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \rightarrow \min$$

Ordinary Least Squares
SUM $(y - \hat{y})^2 \rightarrow \min$



ϵ -Insensitive Tube
Slack Variables ξ_i and ξ_i^*



SVR Intuition

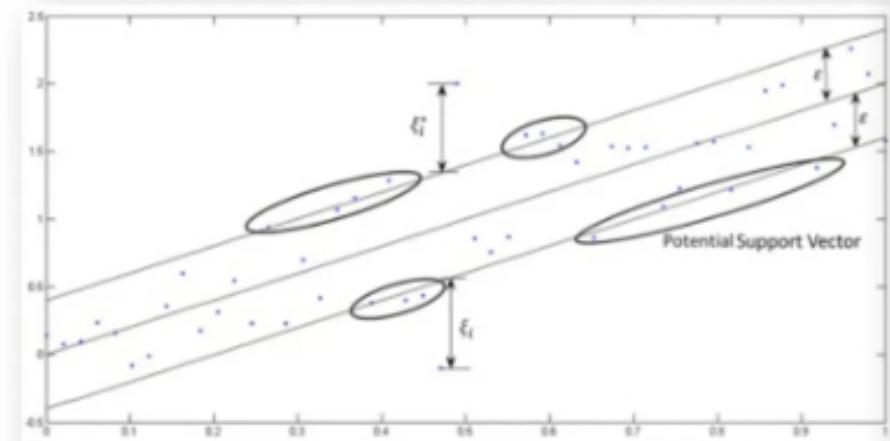
Additional Reading:

*Chapter 4 – Support Vector Regression
(from: Efficient Learning Machines:
Theories, Concepts, and Applications for
Engineers and System Designers)*

By Mariette Awad & Rahul Khanna (2015)

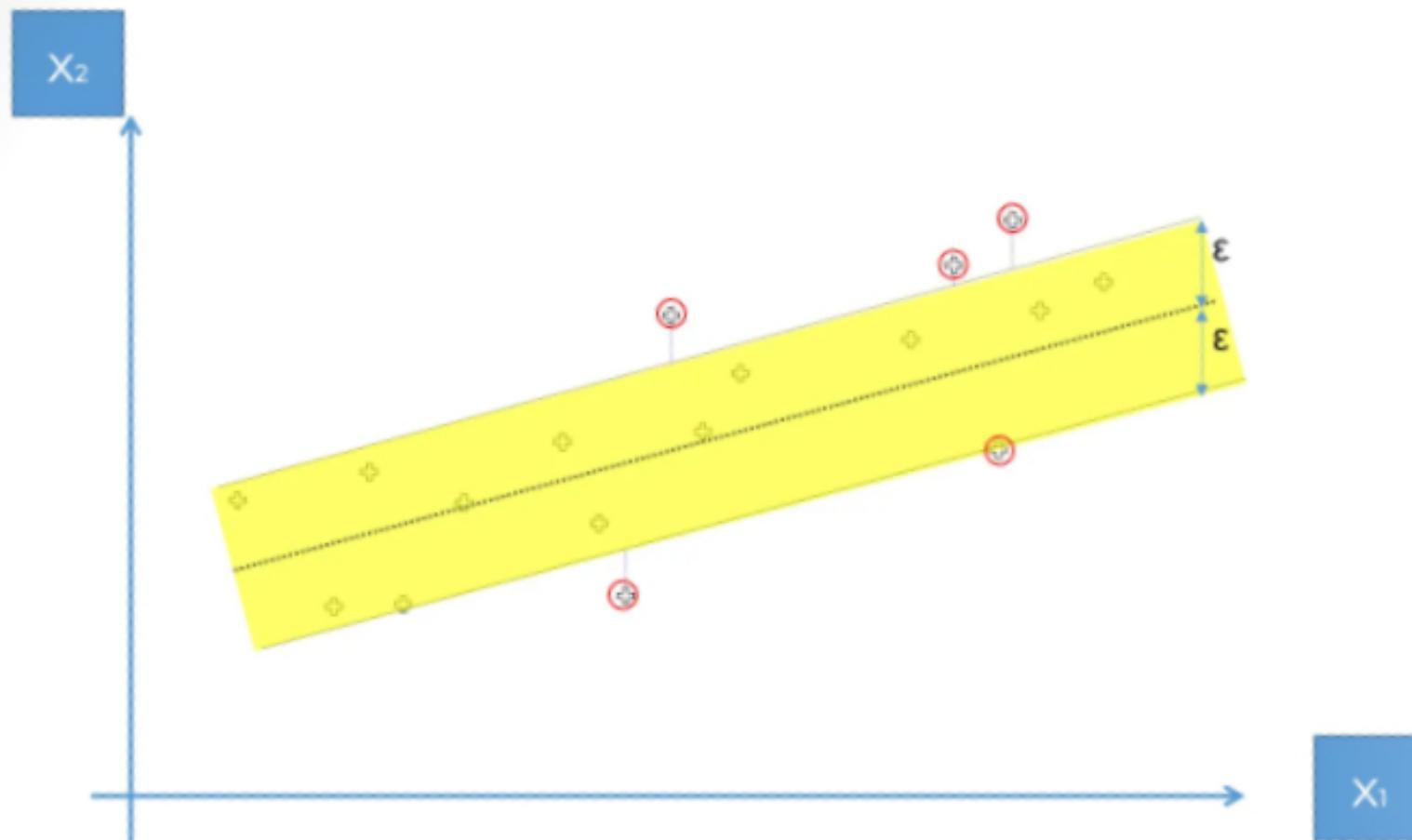
Link:

<https://core.ac.uk/download/pdf/81523322.pdf>



Heads-up about Non-Linear SVR

SVR Intuition



Python - Google Drive

Copy of support_vector_regression.ipynb

File Edit View Insert Runtime Tools Help All changes saved

RAM Disk Editing

Position_Salaries.csv

1 to 10 of 10 entries Filter

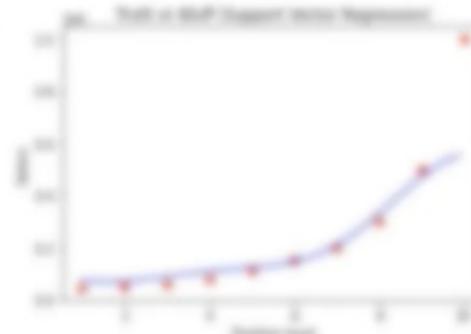
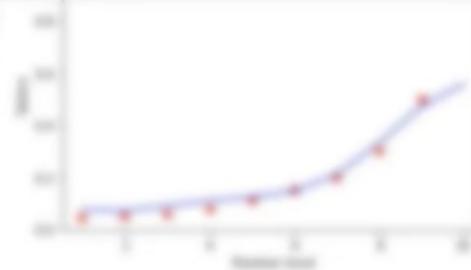
Position	Level	Salary
Business Analyst	1	45000
Junior Consultant	2	50000
Senior Consultant	3	60000
Manager	4	80000
Country Manager	5	110000
Region Manager	6	150000
Partner	7	200000
Senior Partner	8	300000
C-level	9	500000
CEO	10	1000000

Show 10 per page

Visualising the SVR results (for higher resolution and smoother curve)

```
1 X_grid = np.arange(min(sc_X.inverse_transform(X)), max(sc_X.inverse_transform(X)), 0.1)
2 X_grid = X_grid.reshape((len(X_grid), 1))
3 plt.scatter(sc_X.inverse_transform(X), sc_y.inverse_transform(y), color = 'red')
4 plt.plot(X_grid, sc_y.inverse_transform(regressor.predict(sc_X.transform(X_grid))), color = 'blue')
5 plt.title('Truth or Bluff (Support Vector Regression)')
6 plt.xlabel('Position level')
7 plt.ylabel('Salary')
8 plt.show()
```

Disk 76.87 GB available



Heads-up about Non-Linear SVR

Section on SVM:

- SVM Intuition

Section on Kernel SVM:

- Kernel SVM Intuition
- Mapping to a higher dimension
- The Kernel Trick
- Types of Kernel Functions
- Non-linear Kernel SVR

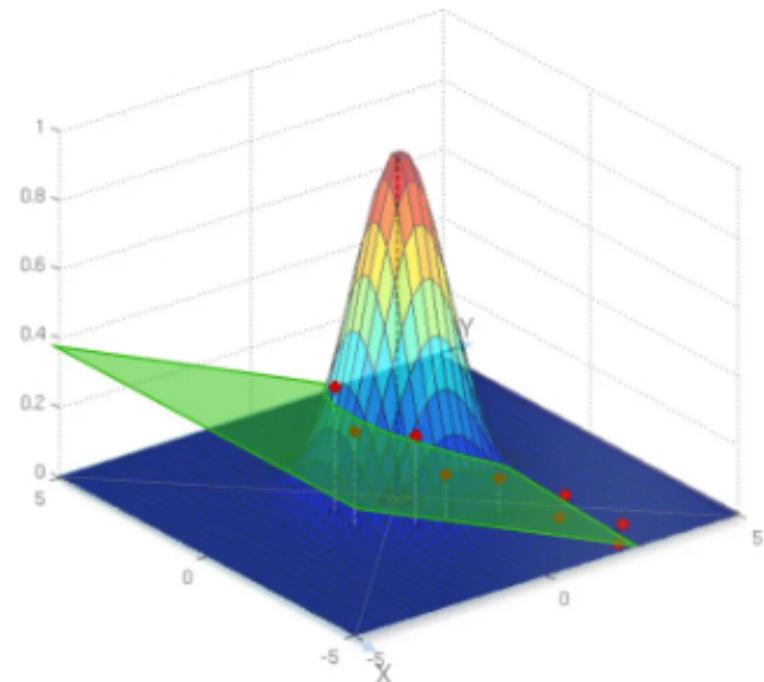


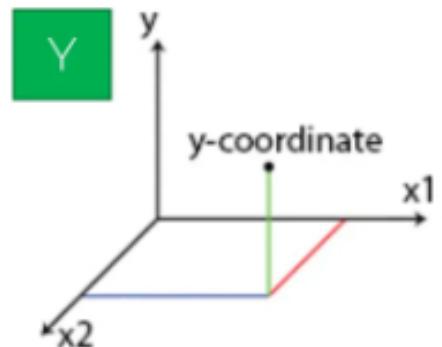
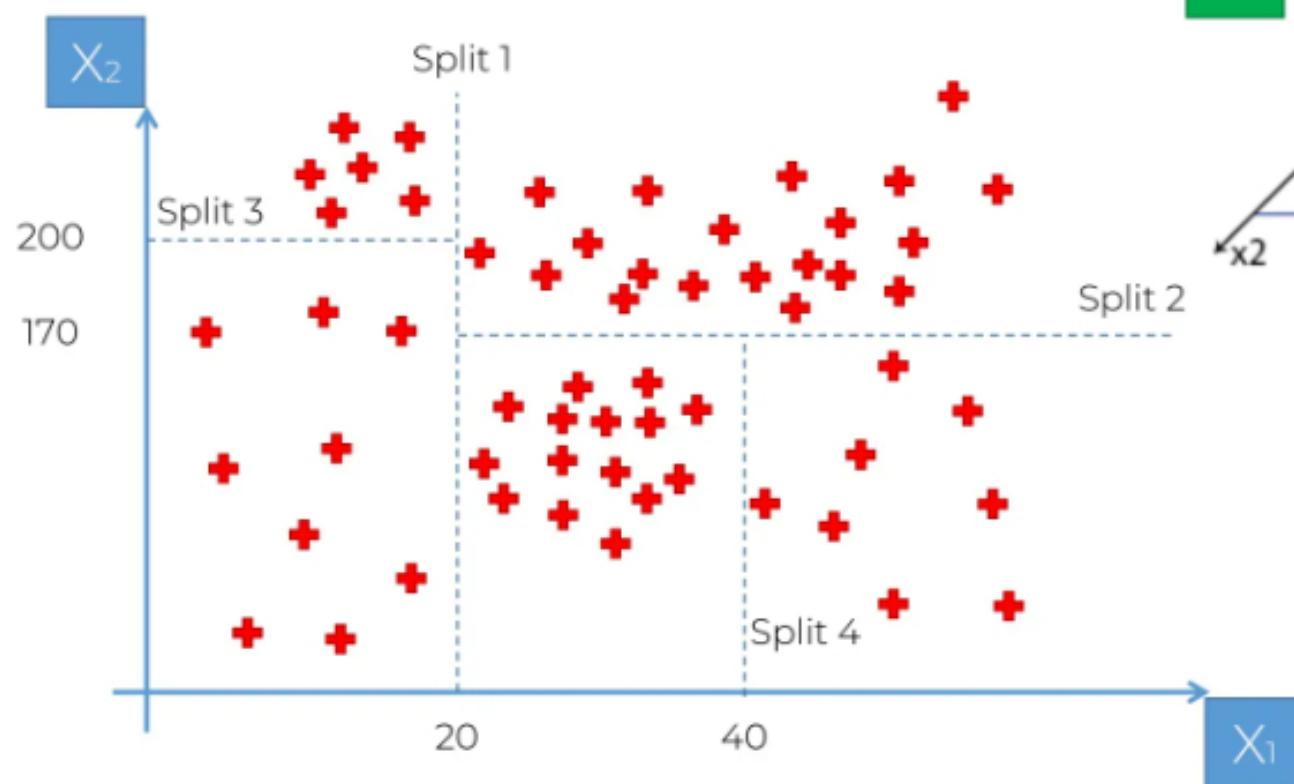
Image source: <http://www.cs.toronto.edu/~duvenaud/cookbook/index.html>

Decision Tree Intuition

Decision Tree Intuition



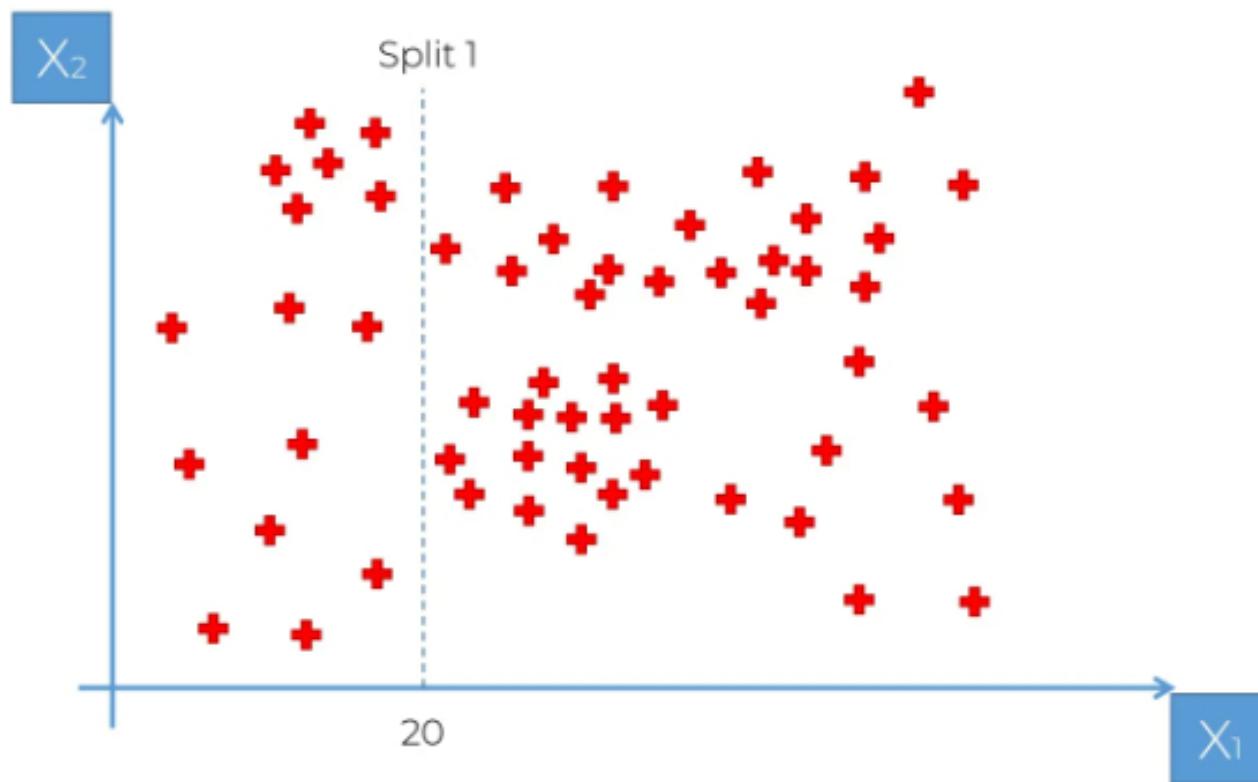
Decision Tree Intuition



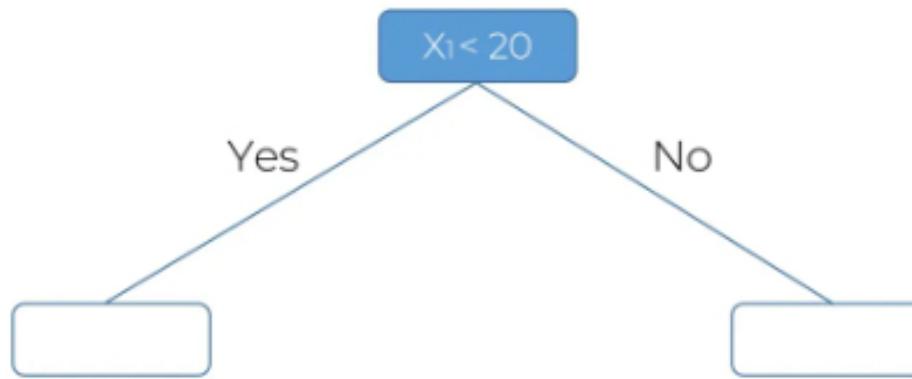
Decision Tree Intuition

Rewind...

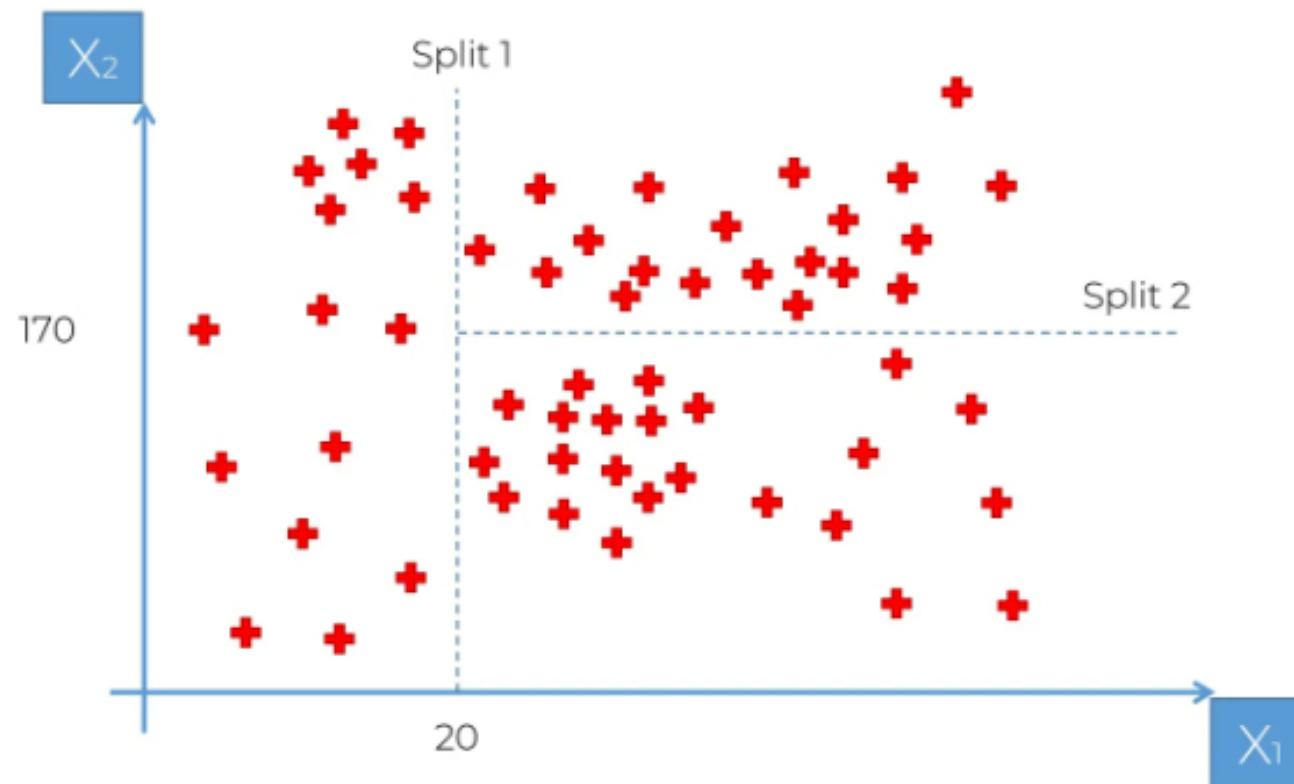
Decision Tree Intuition



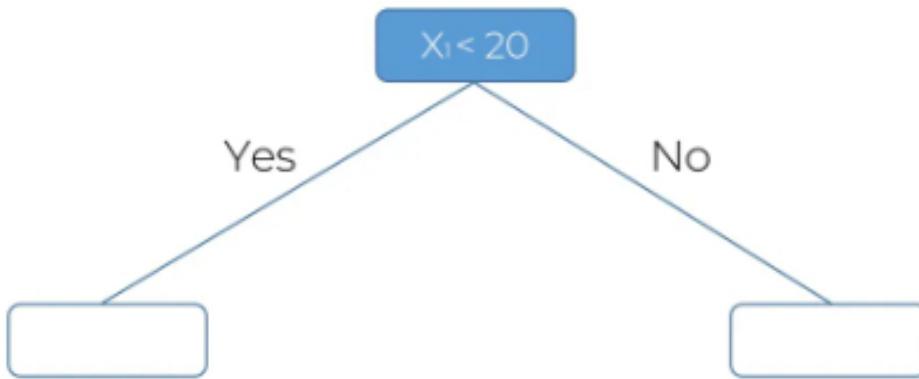
Decision Tree Intuition



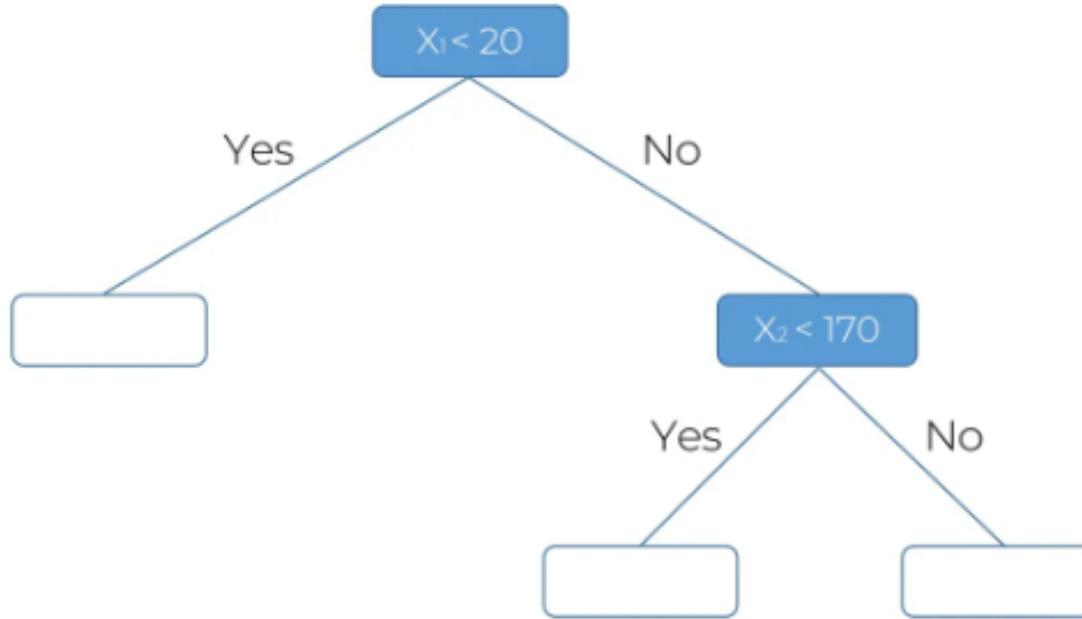
Decision Tree Intuition



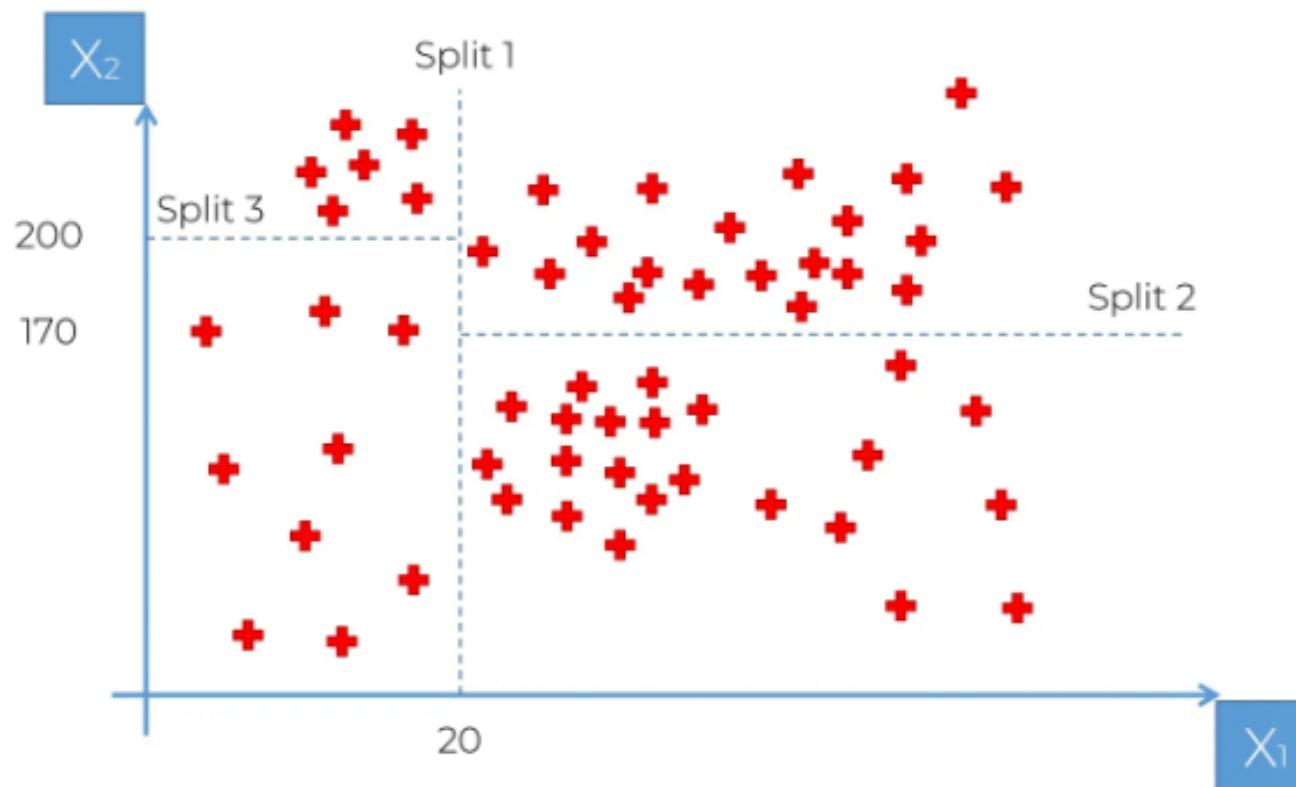
Decision Tree Intuition



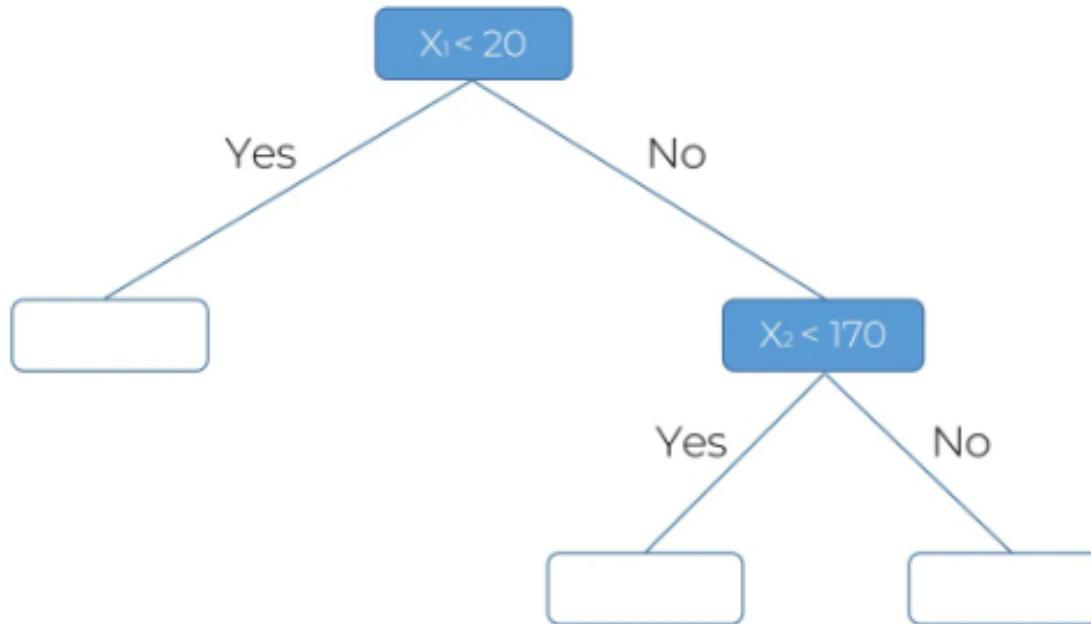
Decision Tree Intuition



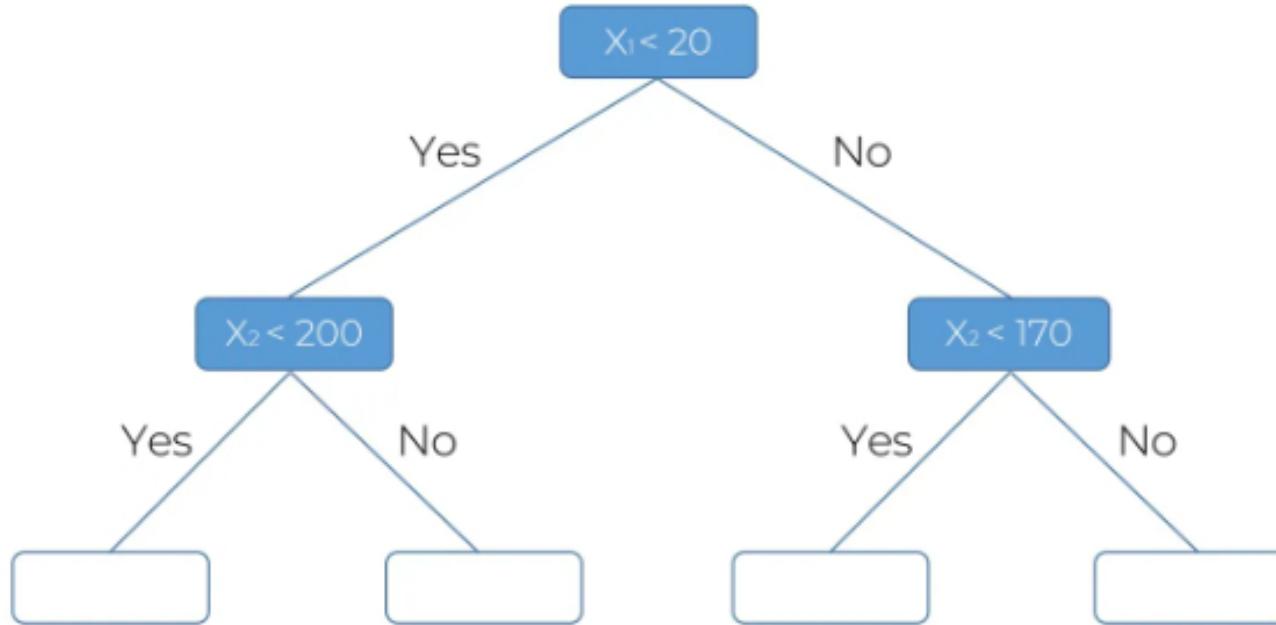
Decision Tree Intuition



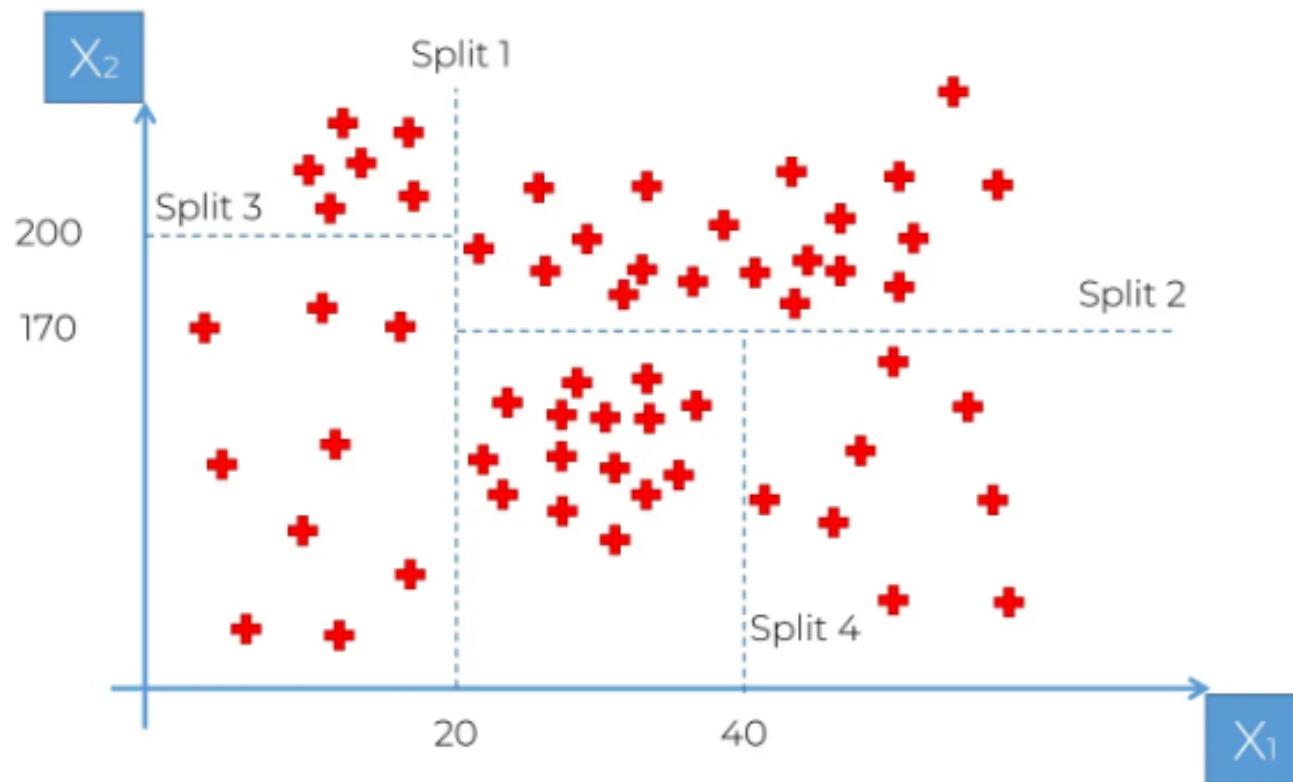
Decision Tree Intuition



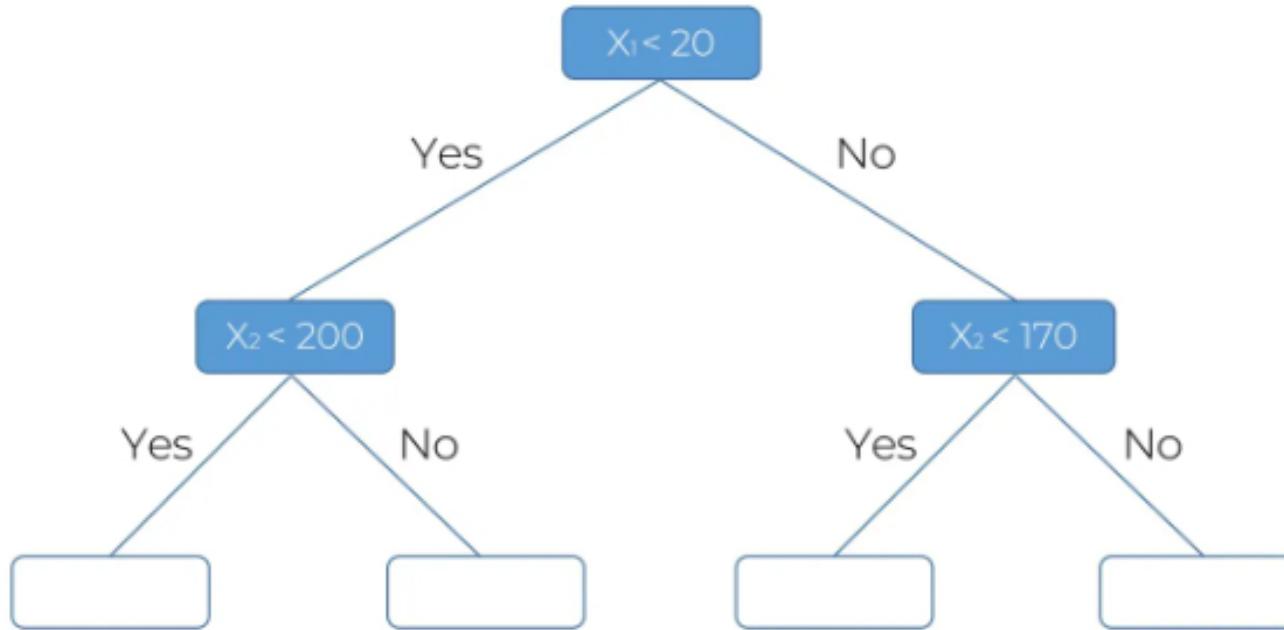
Decision Tree Intuition



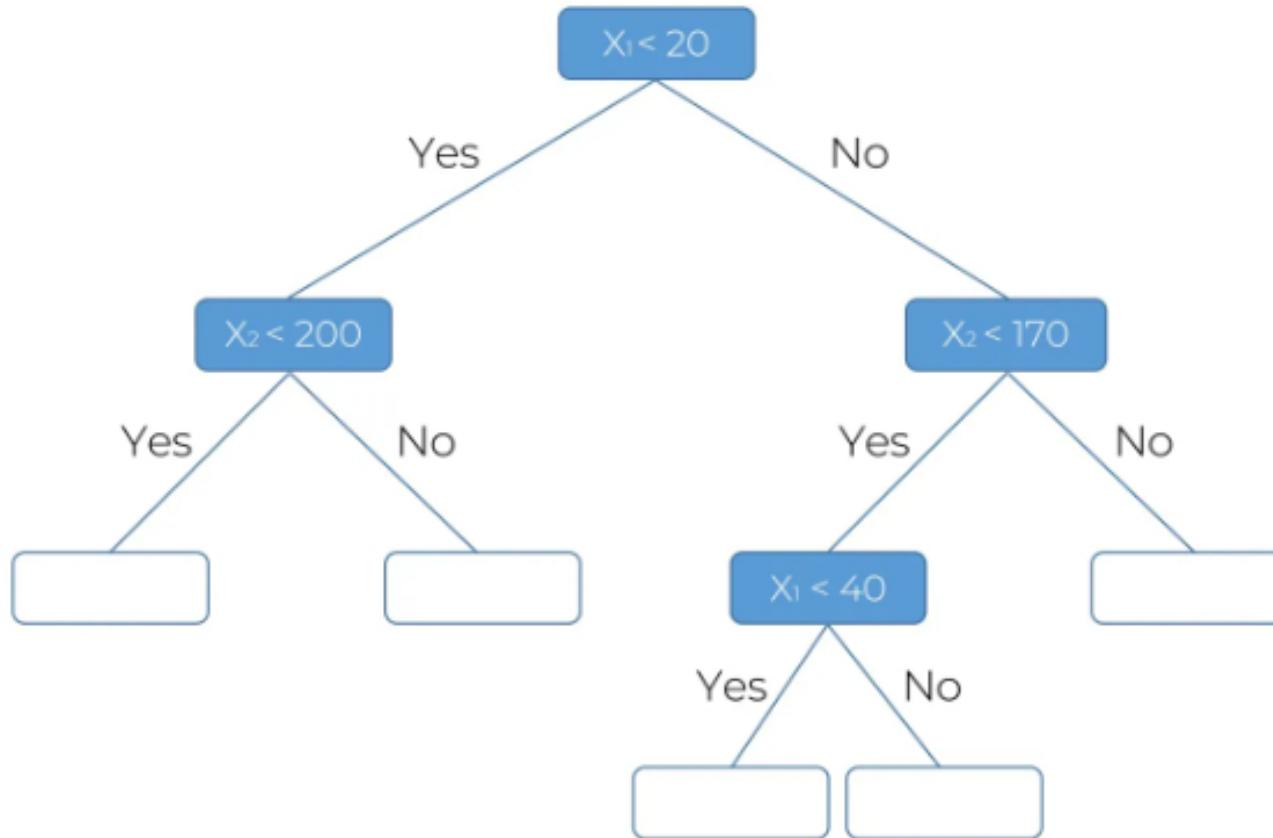
Decision Tree Intuition



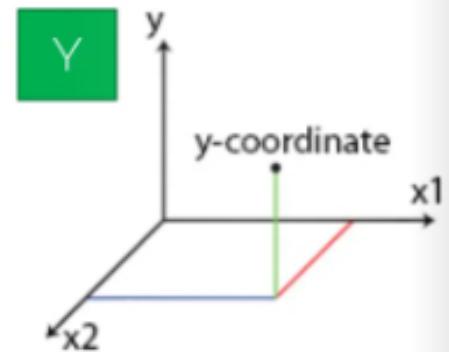
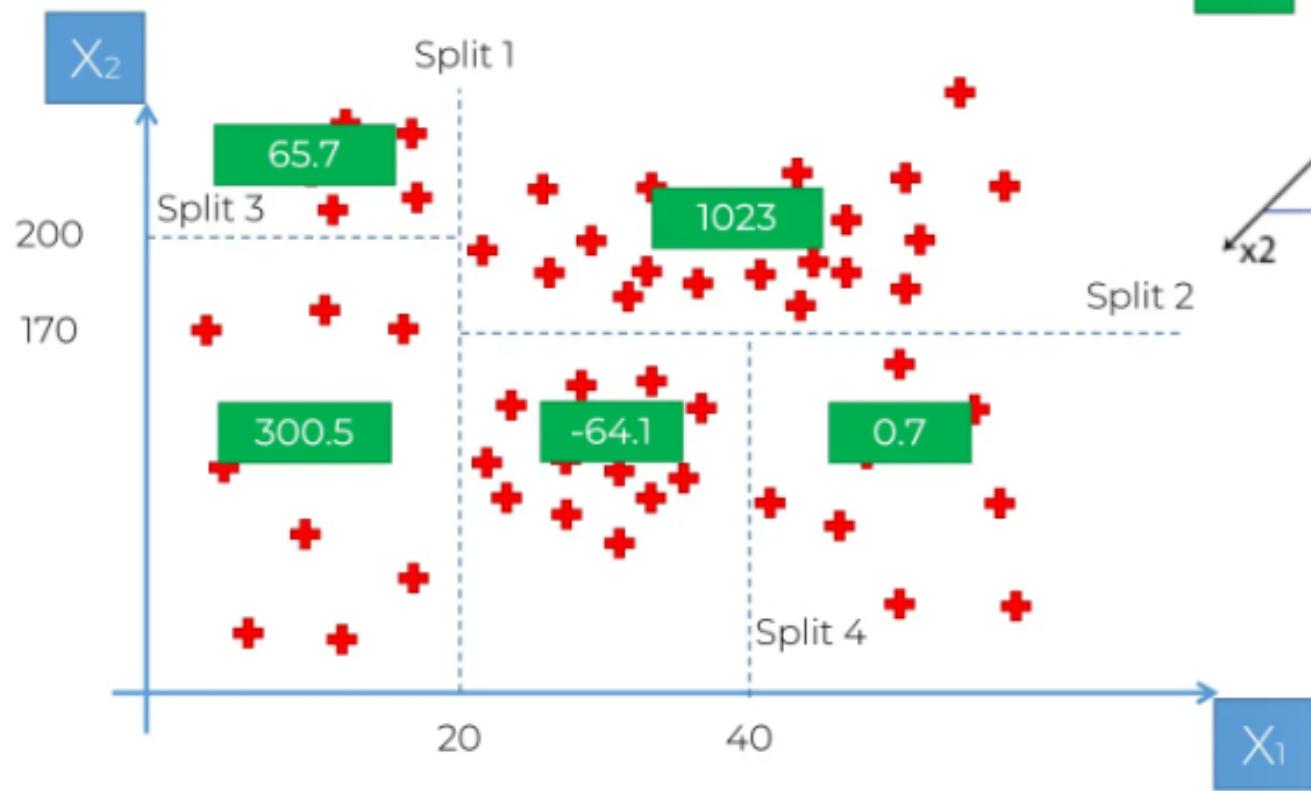
Decision Tree Intuition



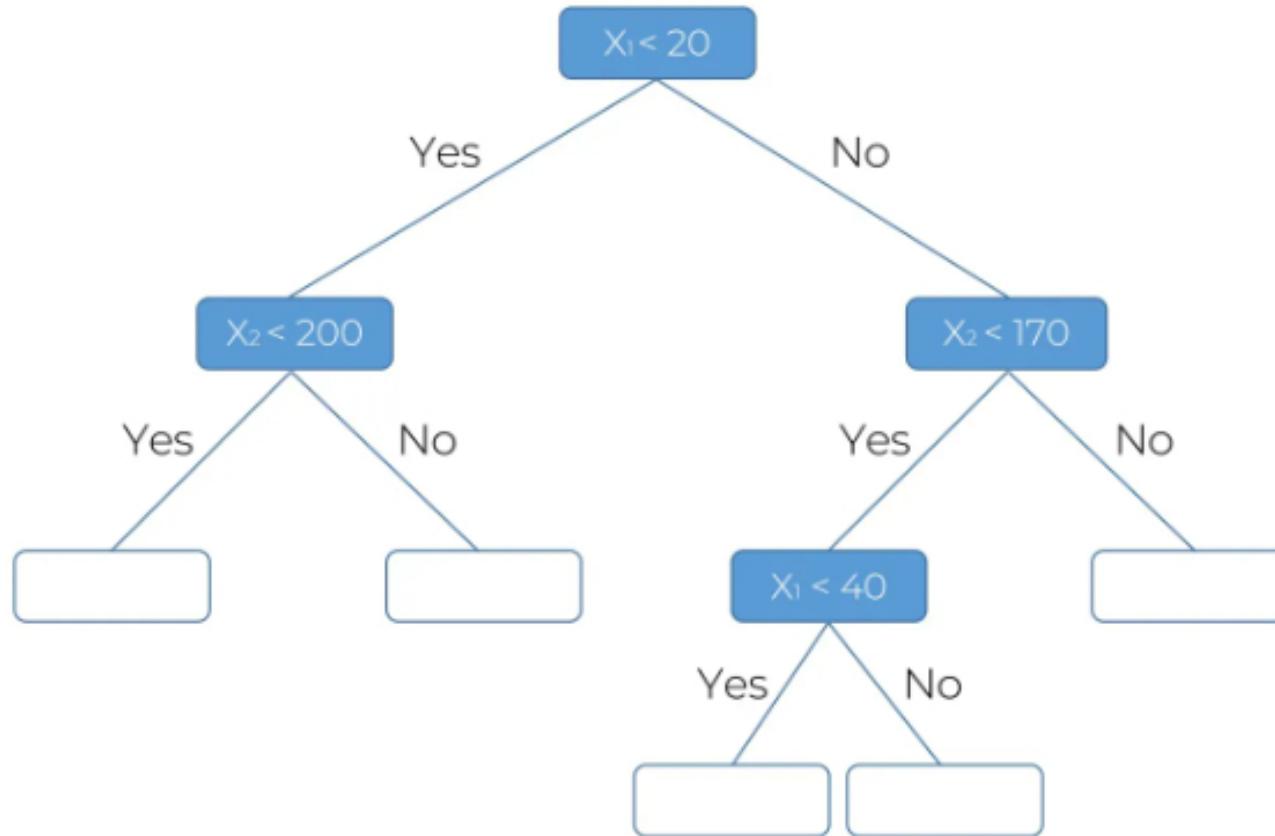
Decision Tree Intuition



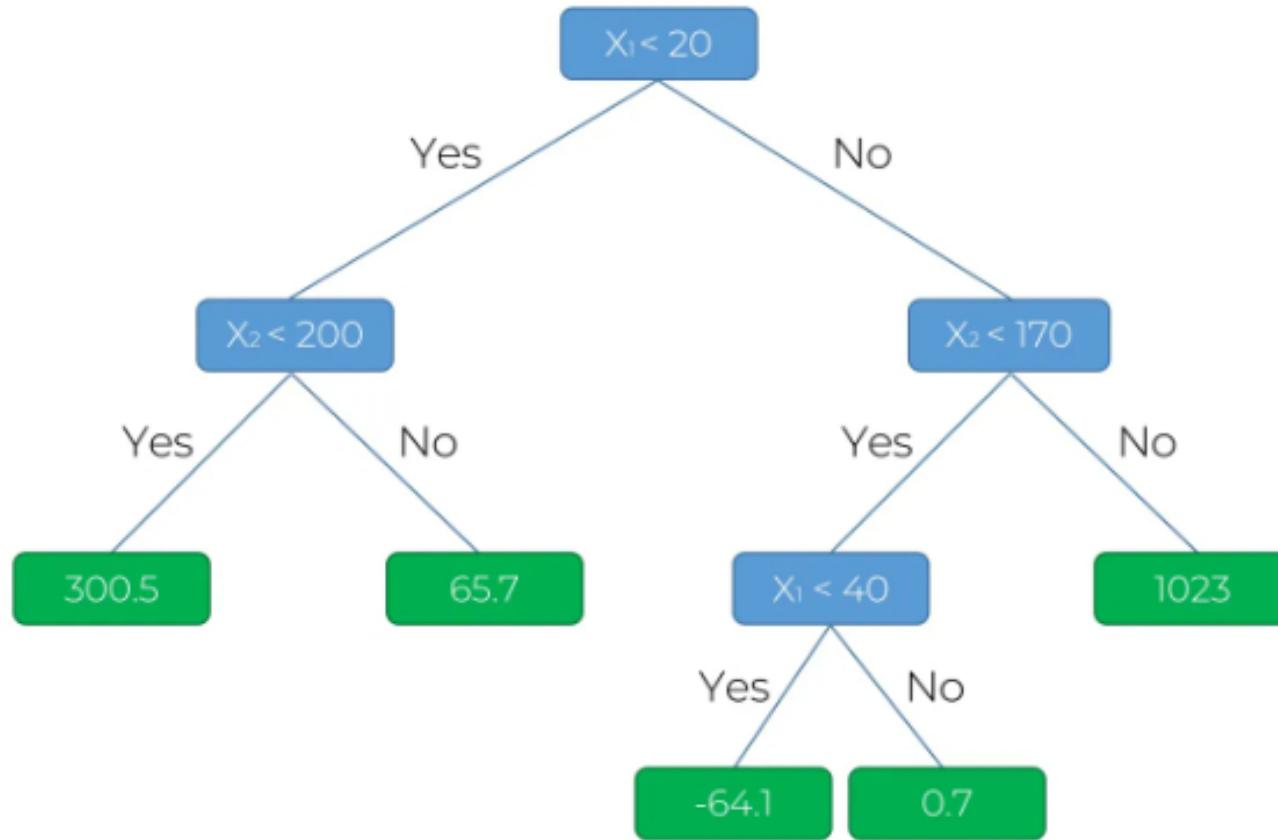
Decision Tree Intuition



Decision Tree Intuition



Decision Tree Intuition



Random Forest Intuition

Random Forest Intuition

Ensemble Learning

Random Forest Intuition

STEP 1: Pick at random K data points from the Training set.



STEP 2: Build the Decision Tree associated to these K data points.



STEP 3: Choose the number Ntree of trees you want to build and repeat STEPS 1 & 2



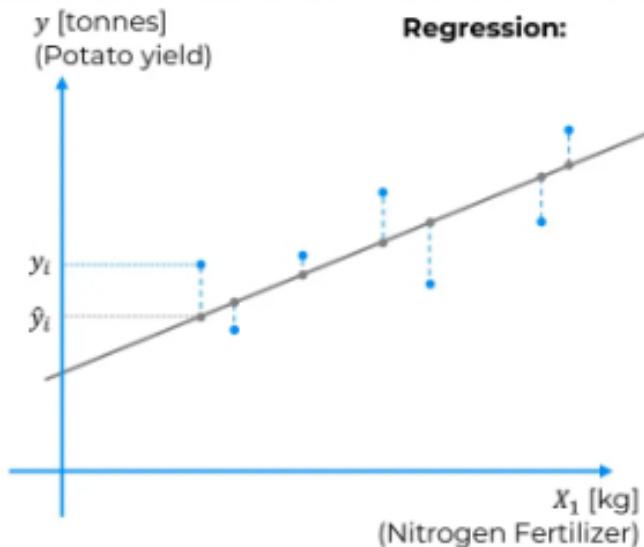
STEP 4: For a new data point, make each one of your Ntree trees predict the value of Y to for the data point in question, and assign the new data point the average across all of the predicted Y values.

R Squared

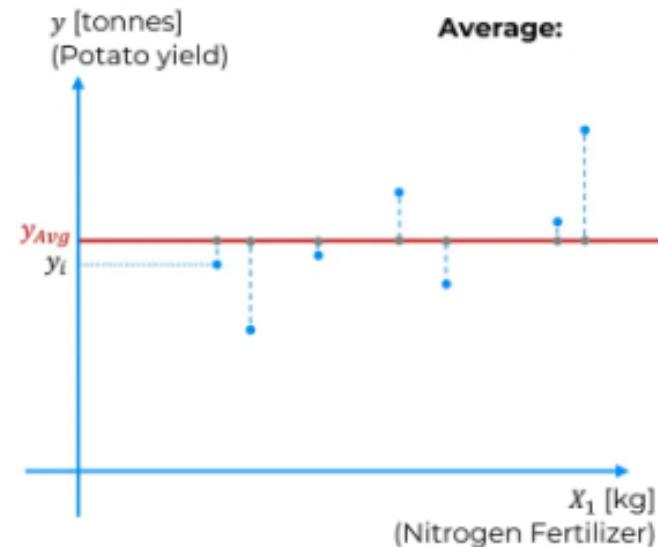
© SuperDataScience



R Squared



$$SS_{res} = \text{SUM}(y_i - \hat{y}_i)^2$$



$$SS_{tot} = \text{SUM}(y_i - y_{avg})^2$$

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

Rule of thumb (for our tutorials)*:

- 1.0 = Perfect fit (suspicious)
- ~0.9 = Very good
- <0.7 = Not great
- <0.4 = Terrible
- <0 = Model makes no sense for this data

*This is highly dependent on the context



Adjusted R Squared



Adjusted R Squared



$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

R² – Goodness of fit
(greater is better)

Problem:

$$\hat{y} = b_0 + b_1X_1 + b_2X_2 + b_3X_3$$

$$SS_{res} = \text{SUM}(y_i - \hat{y}_i)^2$$

SS_{tot} doesn't change

SS_{res} will decrease or stay the same

(This is because of Ordinary Least Squares: SS_{res} → Min)

Solution:

$$\text{Adj } R^2 = 1 - (1 - R^2) \times \frac{n - 1}{n - k - 1}$$

k – number of independent variables

n – sample size

