

Data Preprocessing



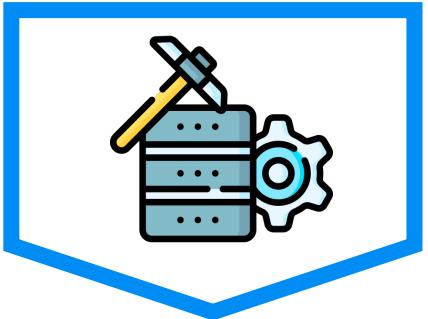
The Machine Learning Process

The Machine Learning Process



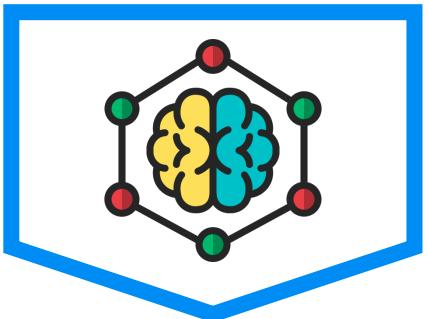
DISTRIBUTION © SUPERDATASCIENCE

www.superdatascience.com



Data Pre-Processing

- Import the data
- Clean the data
- Split into training & test sets
- Feature Scaling



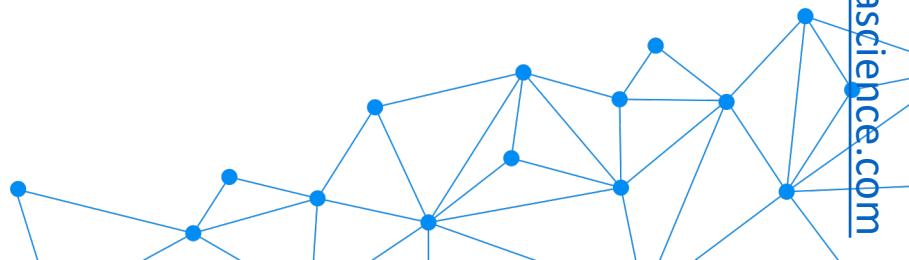
Modelling

- Build the model
- Train the model
- Make predictions



Evaluation

- Calculate performance metrics
- Make a verdict



Training Set & Test Set

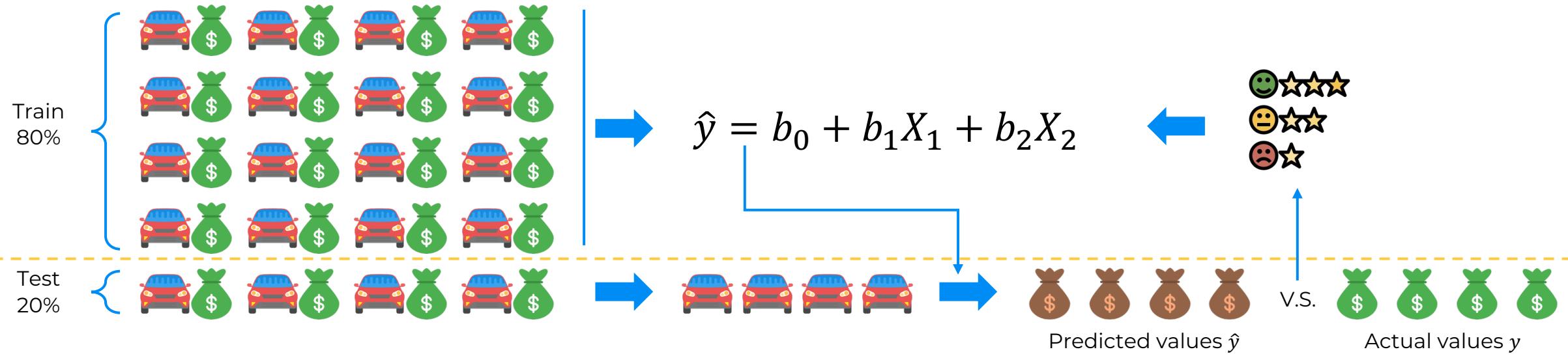




Training Set & Test Set



~





Feature Scaling



Feature Scaling

NOT FOR
DISTRIBUTION © SUPERDATASCIENCE www.superdatascience.com

X1	X2	X3	X4
\$ 179.43	56.784	34.6181	3.55
\$ 641.87	62.054	47.7306	1.692
\$ 556.30	64.13	55.596	1.559
\$ 578.47	63.377	52.7121	1.679
\$ 591.16	61.553	46.1315	1.984
\$ 242.03	58.29	39.2952	2.942
\$ 364.66	59.93	42.4628	2.494
\$ 190.68	57.271	36.2725	3.419
\$ 547.23	63.763	54.1971	1.634
\$ 359.69	59.375	41.5105	2.128
\$ 438.08	60.484	43.493	2.47
\$ 637.17	62.525	49.428	1.725



Feature Scaling

Normalization

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

[0 ; 1]

Standardization

$$X' = \frac{X - \mu}{\sigma}$$

[-3 ; +3]

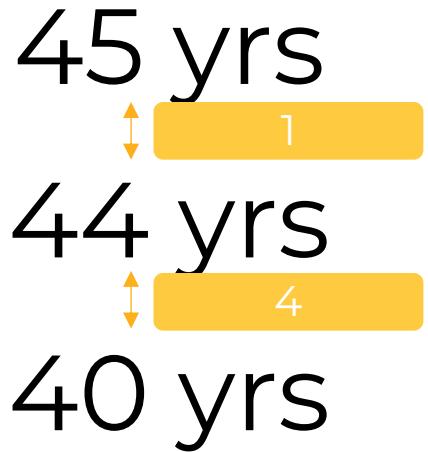
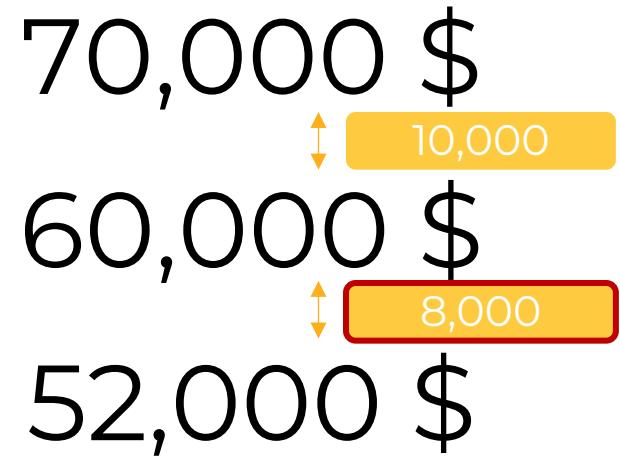
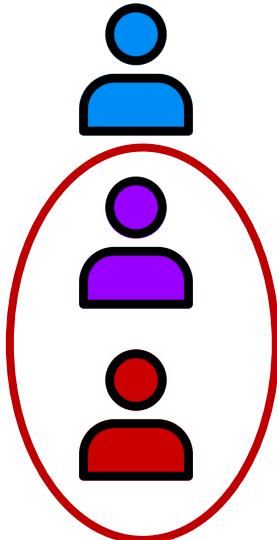


Feature Scaling



NOT FOR

DISTRIBUTION © SUPERDATASCIENCE

www.superdatascience.com

Feature Scaling

Normalization

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

[0 ; 1]



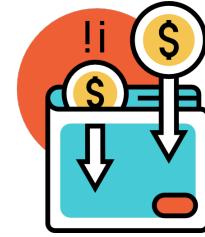
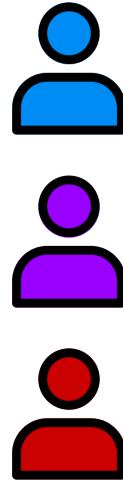
Feature Scaling



NOT FOR
DISTRIBUTION

DISTRIBUTION © SUPERDATASCIENCE

www.superdatascience.com



70,000 \$
60,000 \$
52,000 \$



45 yrs
44 yrs
40 yrs



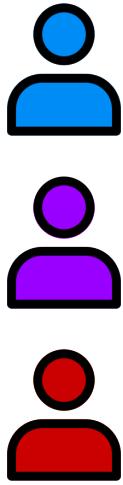
Feature Scaling



NOT FOR DISTRIBUTION

DISTRIBUTION © SUPERDATASCIENCE

www.superdatascience.com



1
0.444
0



45 yrs
44 yrs
40 yrs



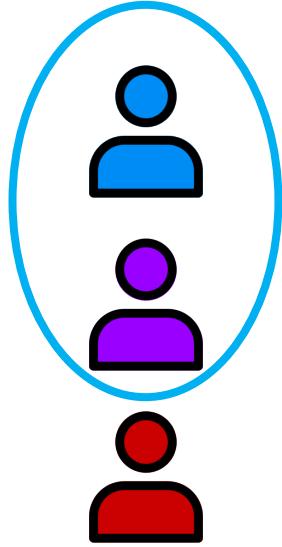
Feature Scaling



NOT FOR
DISTRIBUTION

DISTRIBUTION © SUPERDATASCIENCE

www.superdatascience.com



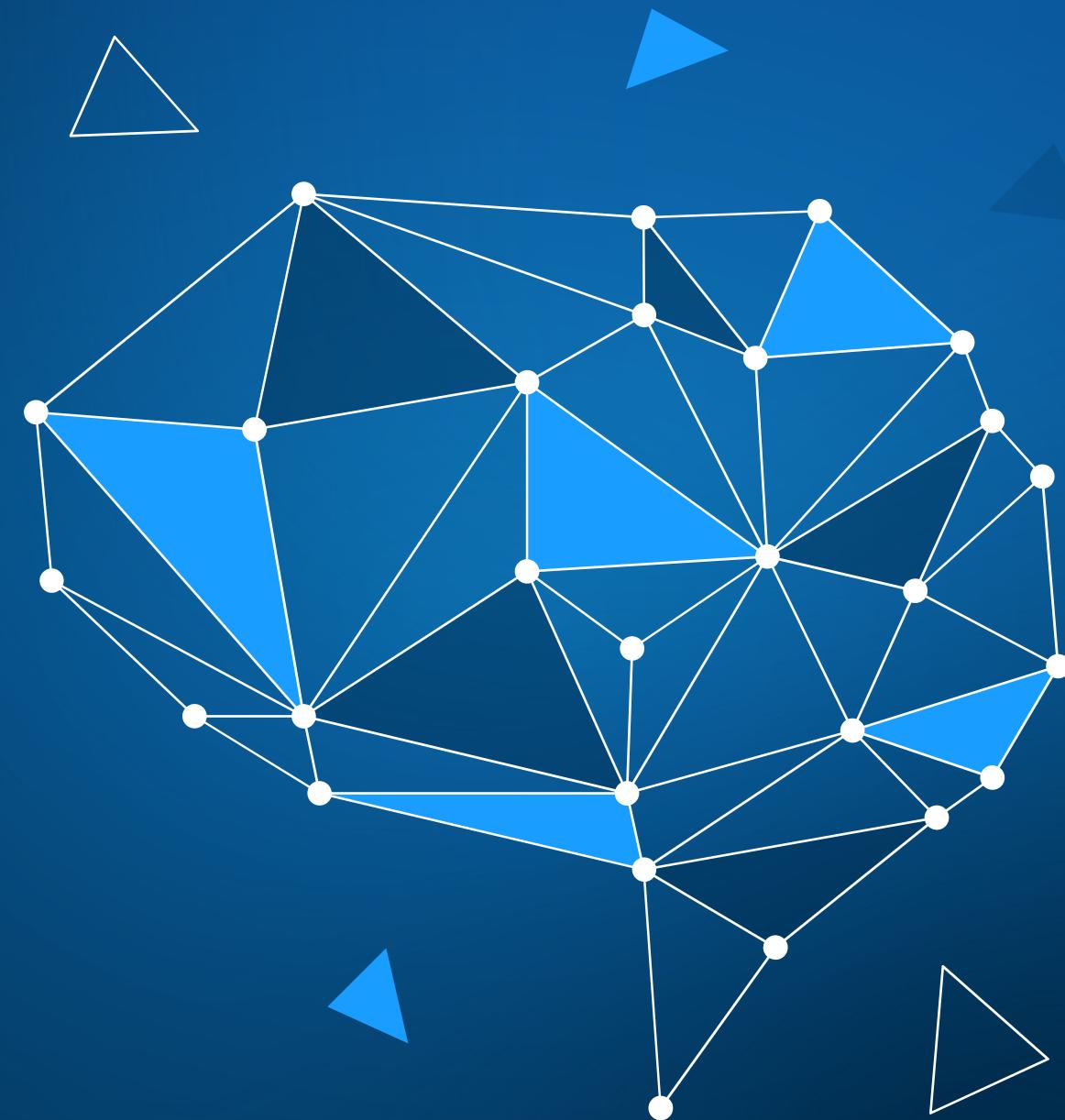
1
0.444
0



1
0.75
0



Regression



Simple Linear Regression





Simple Linear Regression

NOT FOR DISTRIBUTION © SUPERDATASCIENCE

www.superdatascience.com

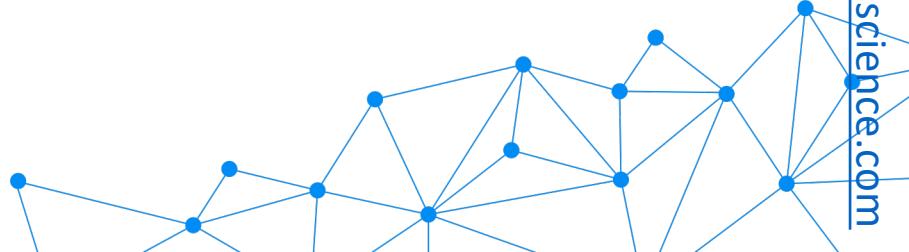
$$\hat{y} = b_0 + b_1 X_1$$

Dependent variable

y-intercept (constant)

Slope coefficient

Independent variable





Simple Linear Regression



~

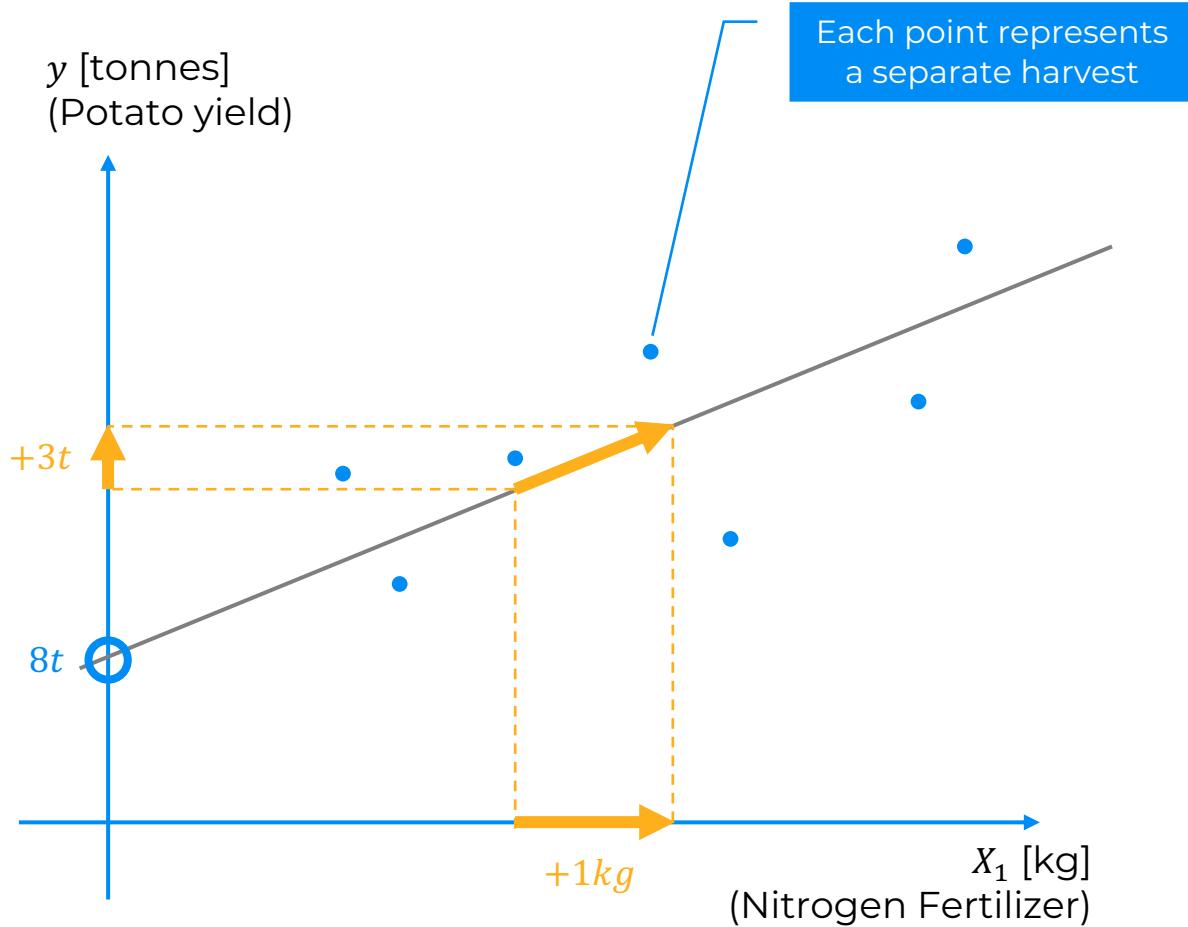


$$\hat{y} = b_0 + b_1 X_1$$

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

$$b_0 = 8[t]$$

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



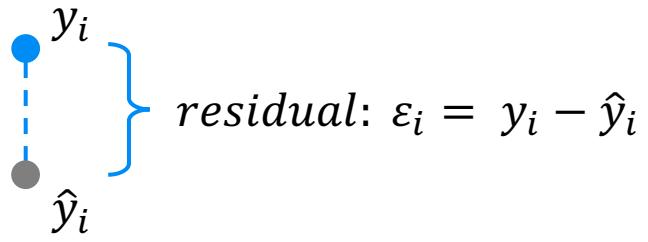
Ordinary Least Squares





Simple Linear Regression

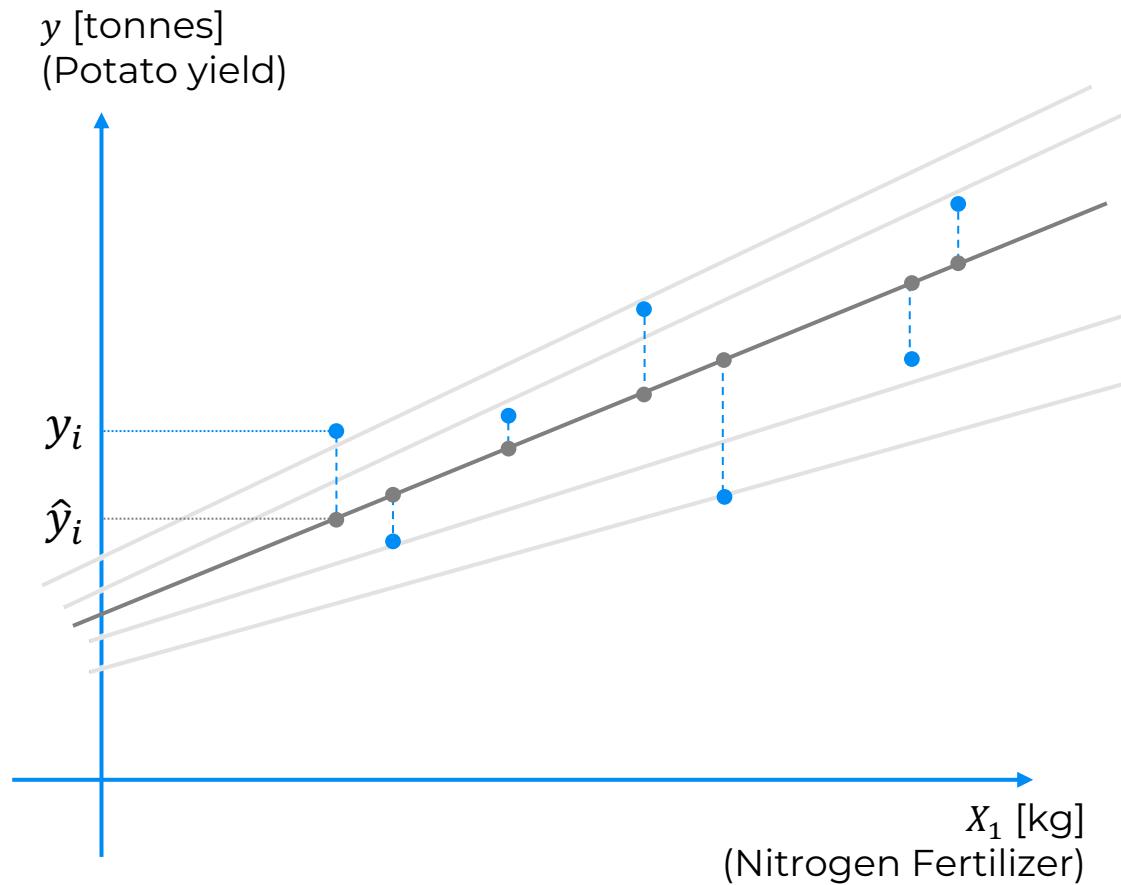
Ordinary Least Squares:



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

b_0, b_1 such that:

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



Multiple Linear Regression





Multiple Linear Regression

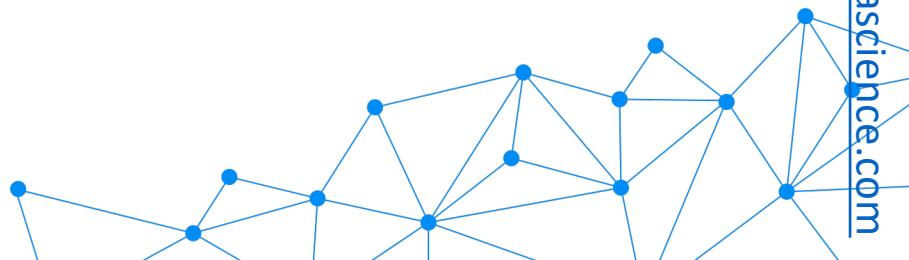
$$\hat{y} = b_0 + b_1X_1 + b_2X_2 + \cdots + b_nX_n$$

Dependent variable
y-intercept (constant)

Independent variable 1
Slope coefficient 1

Independent variable 2
Slope coefficient 2

Independent variable n
Slope coefficient n

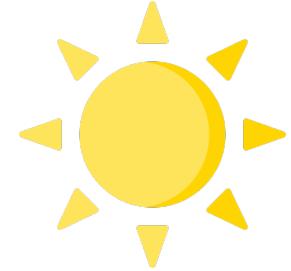




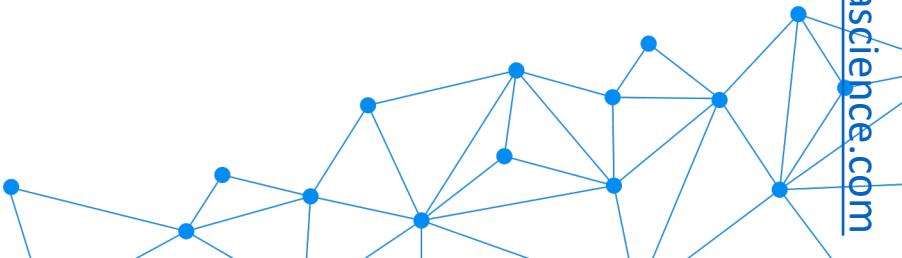
Multiple Linear Regression



~



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



Additional Reading

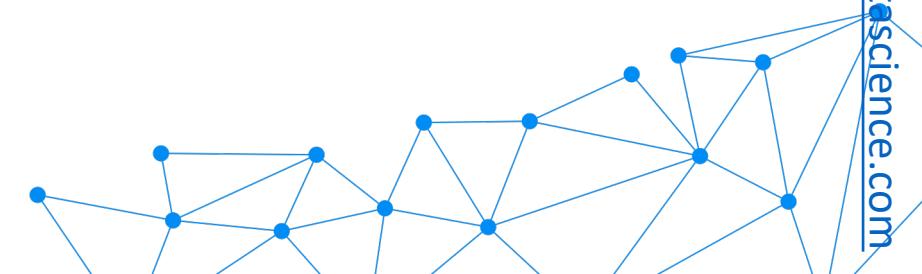
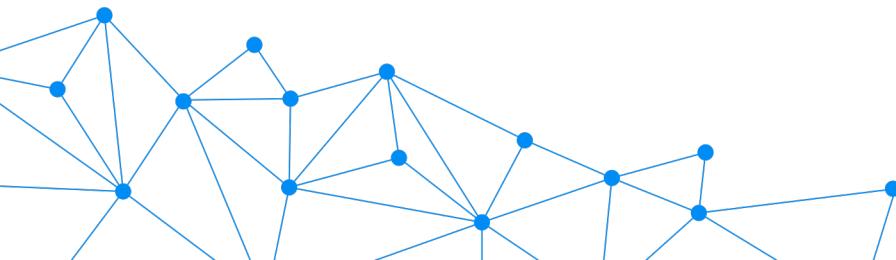
The Application of Multiple Linear Regression and Artificial Neural Network Models for Yield Prediction of Very Early Potato Cultivars before Harvest

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 [$\text{kg}\cdot\text{ha}^{-1}$] in the periods: planting—20 June	80–155
PHOSP	sum of phosphorus fertilization [$\text{kg}\cdot\text{ha}^{-1}$]	28.2–150
POTAS	sum of potassium fertilization [$\text{kg}\cdot\text{ha}^{-1}$]	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 [$\text{mg P}_2\text{O}_5\cdot 100 \text{ g}^{-1}$ soil]	14–26.2
SFERTK	soil fertility in potassium [$\text{mg K}_2\text{O}\cdot 100 \text{ g}^{-1}$ soil]	11.7–19.2
SFERTM	soil fertility in magnesium [$\text{mg Mg}\cdot 100 \text{ g}^{-1}$ soil]	3–9.1
YIELDP1	tuber yield [$\text{t}\cdot\text{ha}^{-1}$], harvest 40 days from full emergence	11.6–41.3

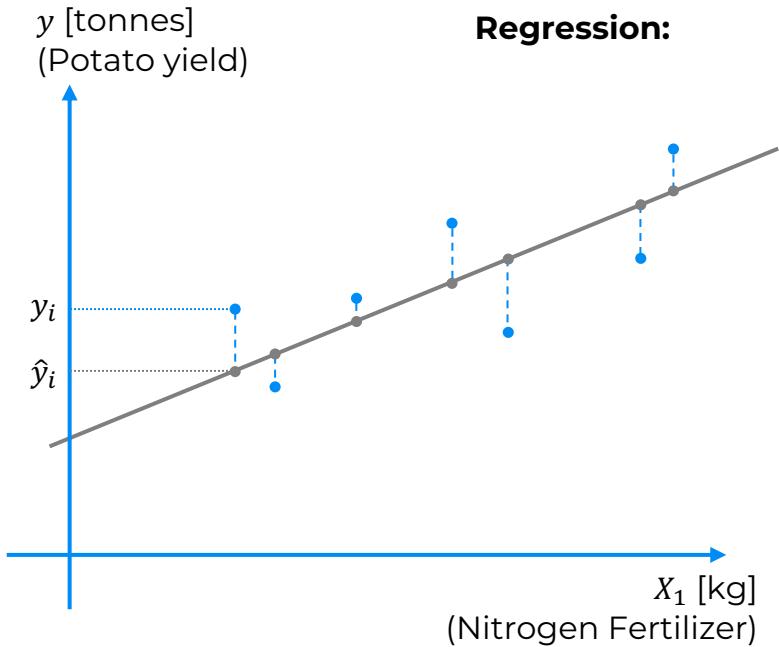


R Squared

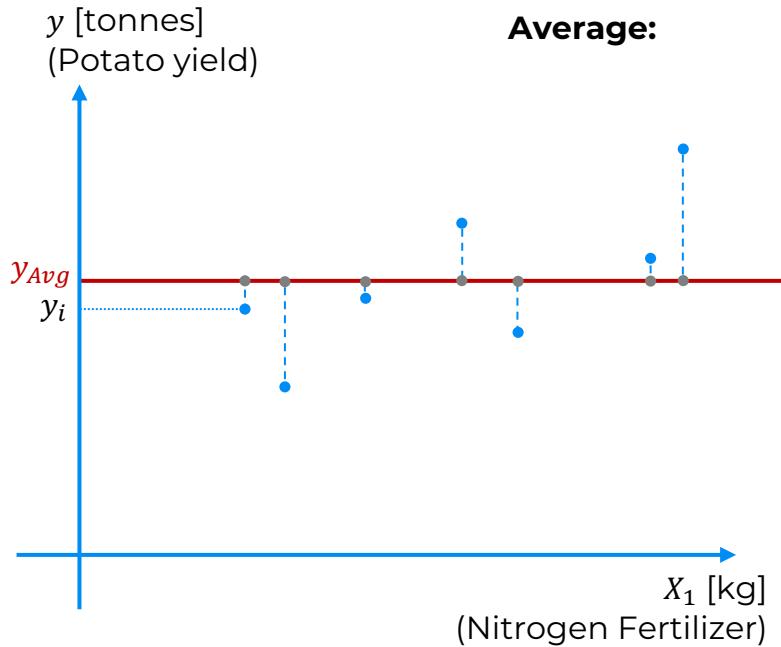




R Squared



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



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

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

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_{tot} doesn't change

SS_{res} will decrease or stay the same

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

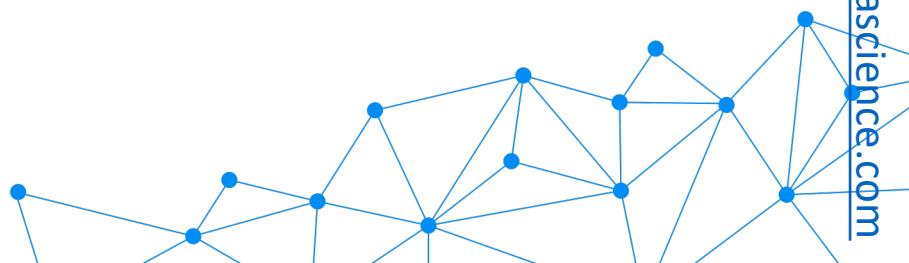
(This is because of Ordinary Least Squares: $SS_{res} \rightarrow \text{Min}$)

Solution:

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

k – number of independent variables

n – sample size



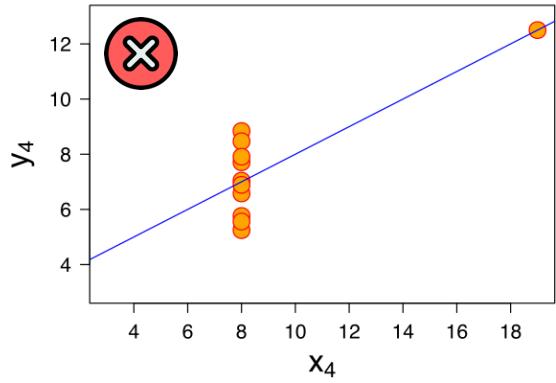
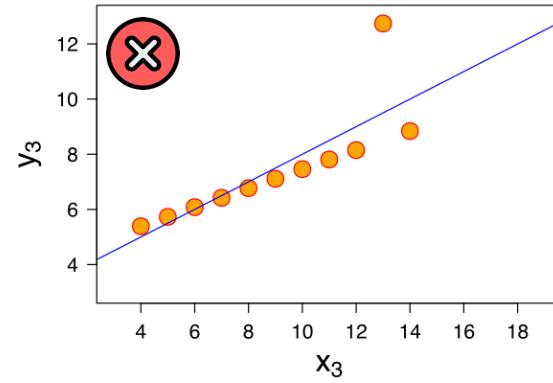
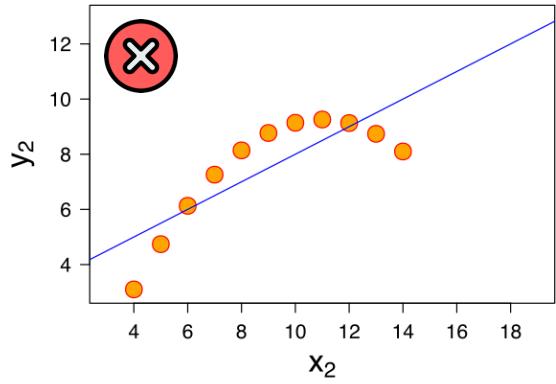
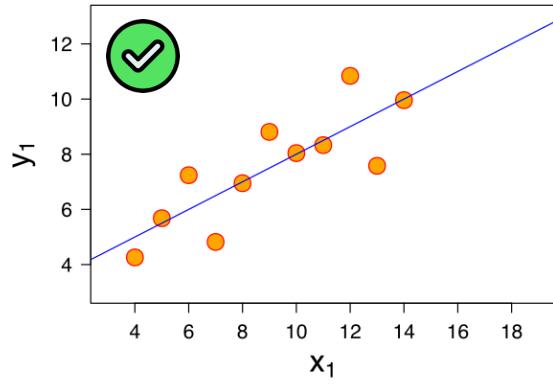
Assumptions Of Linear Regression





Assumptions of Linear Regression

Anscombe's quartet (1973):

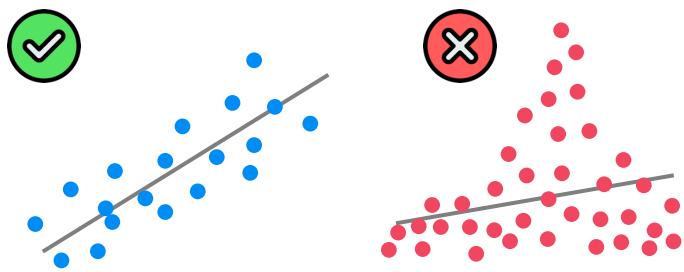




Assumptions of Linear Regression

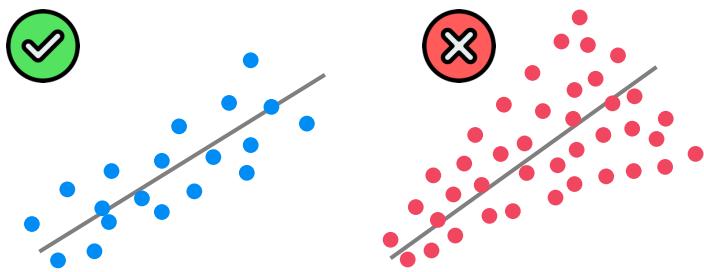
1. Linearity

(Linear relationship between Y and each X)



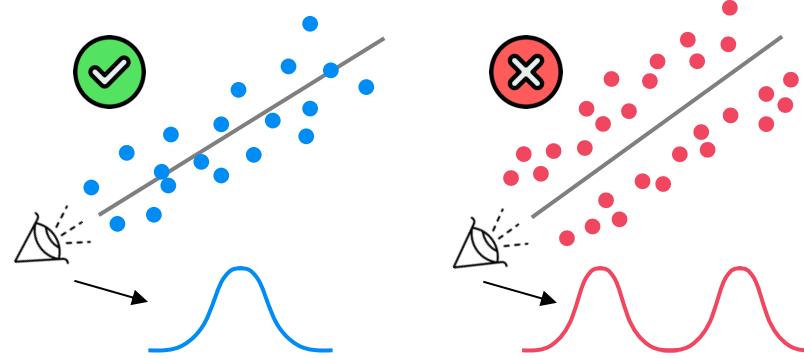
2. Homoscedasticity

(Equal variance)



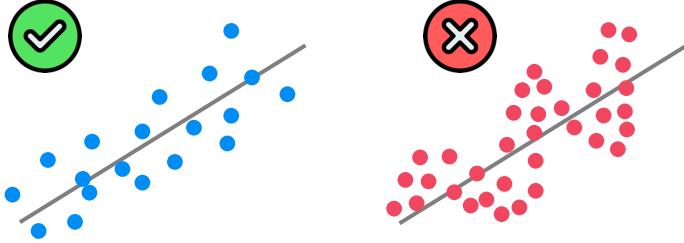
3. Multivariate Normality

(Normality of error distribution)



4. Independence

(of observations. Includes “no autocorrelation”)



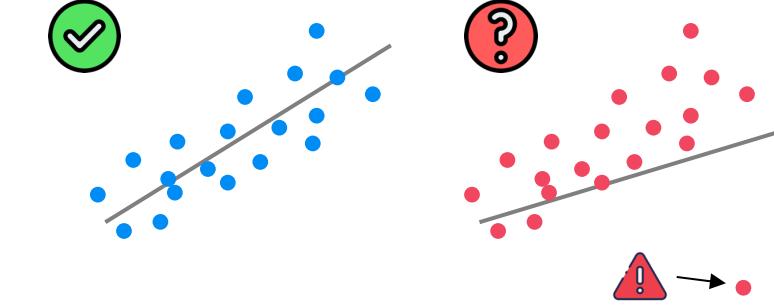
5. Lack of Multicollinearity

(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

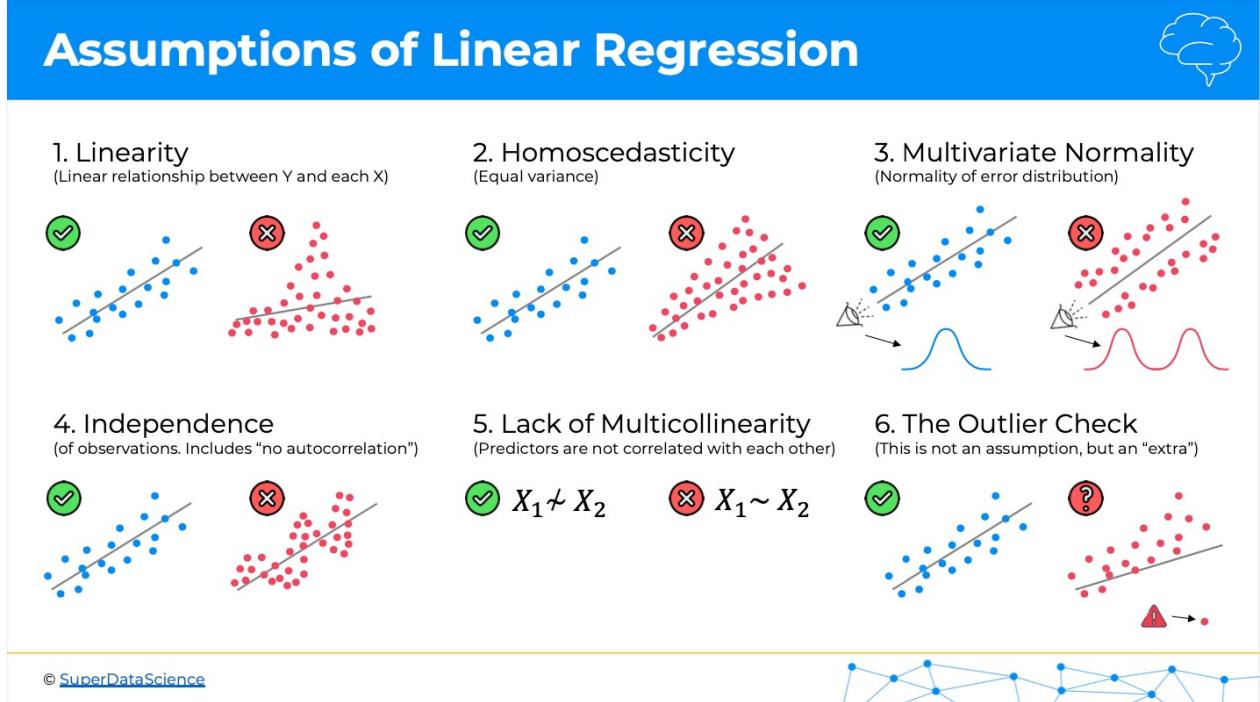
(This is not an assumption, but an “extra”)



Bonus



Download the Assumptions poster at:
superdatascience.com/assumptions



Additional Reading

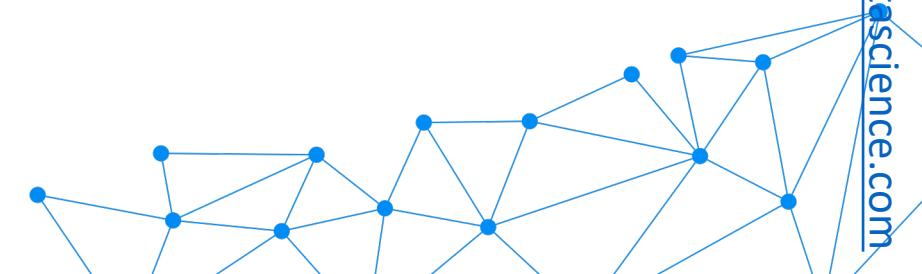
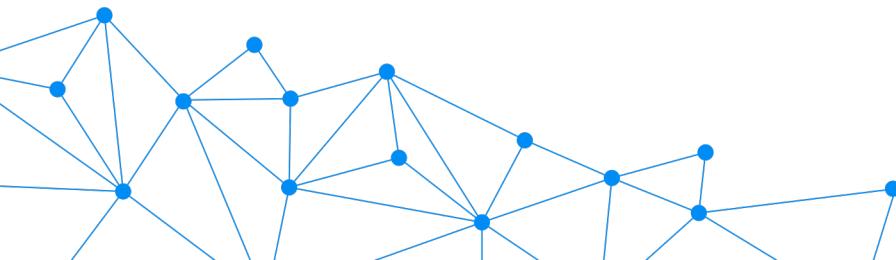


*Verifying the Assumptions of Linear Regression
in Python and R*

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

$$D_2 = 1 - D_1$$

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



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

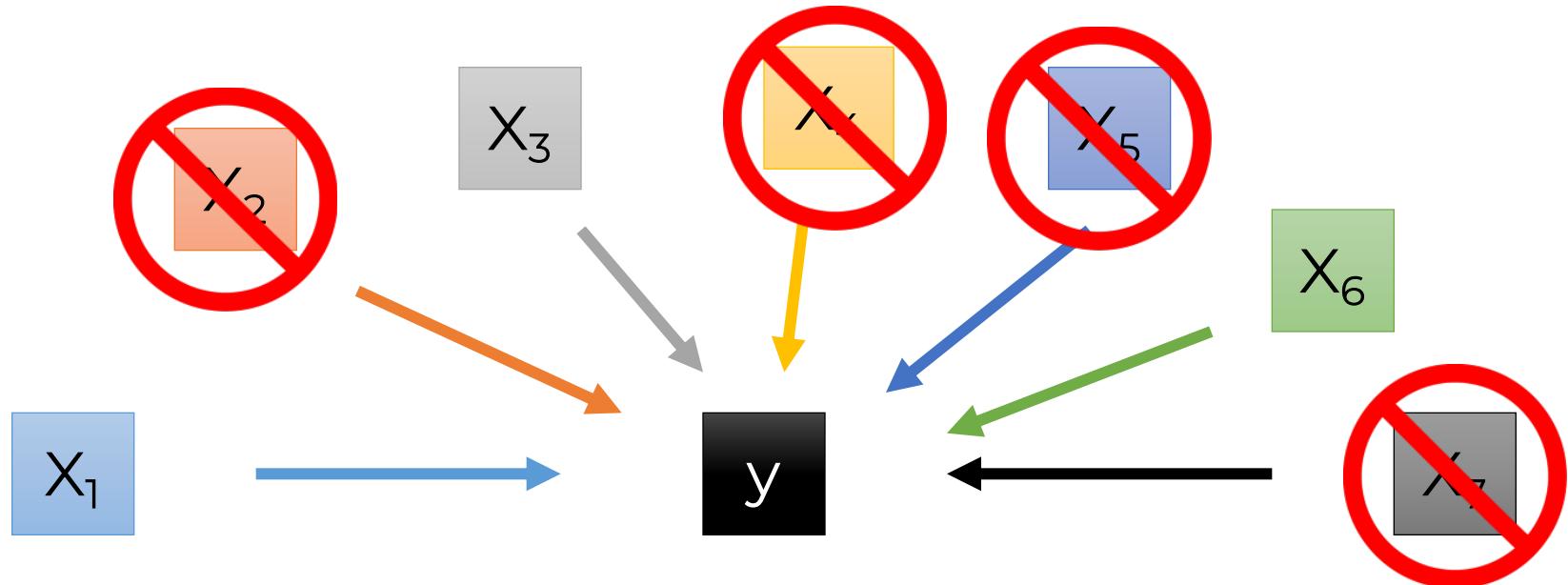
New York	California
1	0
0	1
0	1
1	0
0	1

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

Always omit one dummy variable

Building A Model (Step-By-Step)

Building A Model

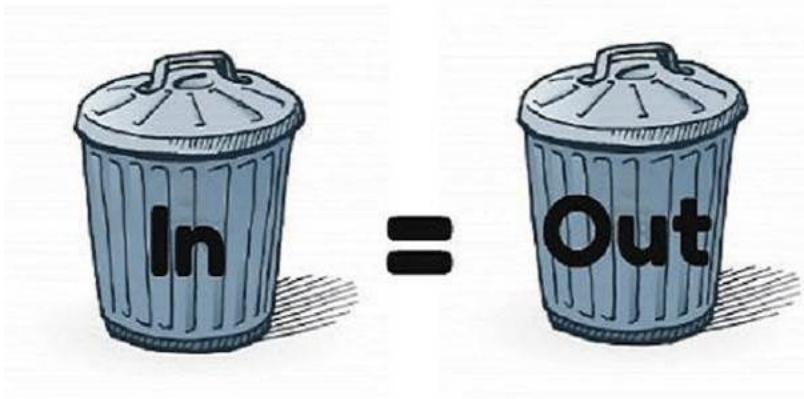


Why?

Building A Model

7)

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



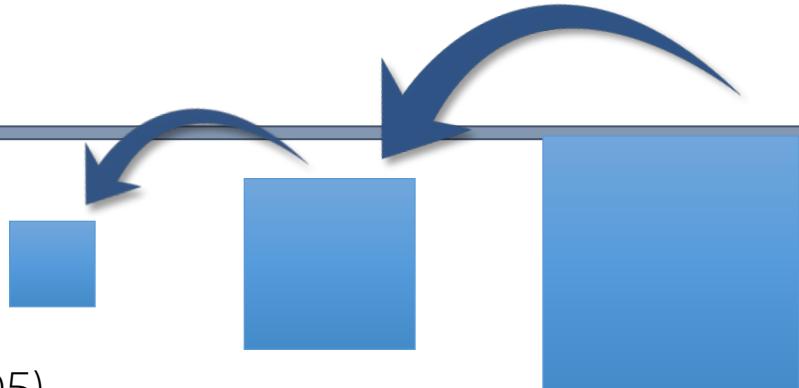
STEP 3: Consider the predictor with the highest P-value. If $P > SL$, go to STEP 4, otherwise go to FIN



STEP 4: Remove the predictor



STEP 5: Fit model without this variable*



FIN: Your Model Is Ready

Building A Model

Forward Selection

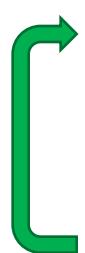
STEP 1: Select a significance level to enter the model (e.g. SL = 0.05)



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

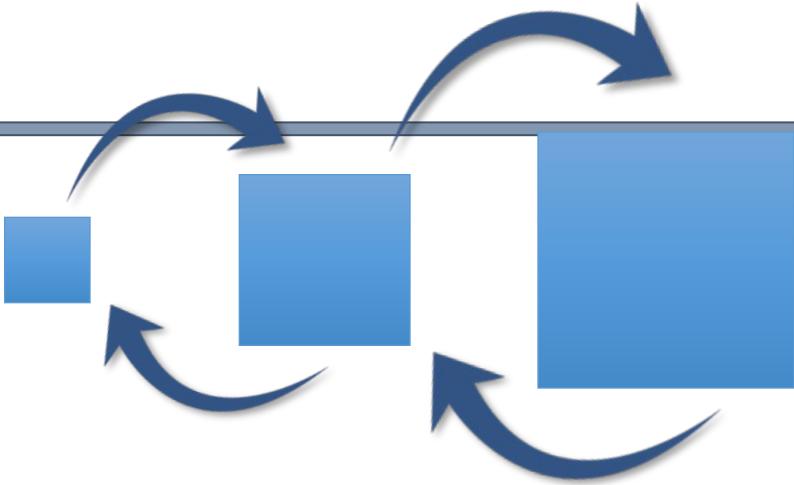


FIN: Keep the previous model

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)



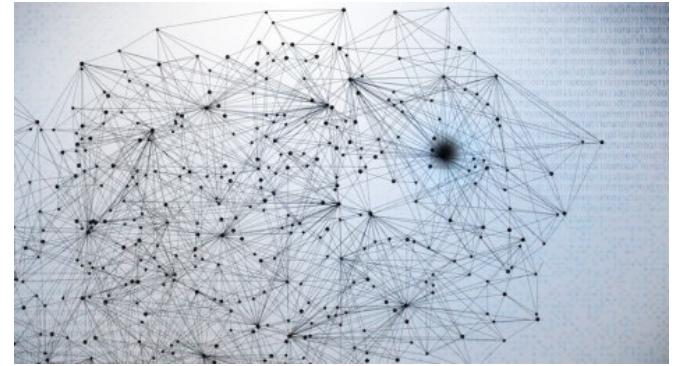
STEP 2: Construct All Possible Regression Models: $2^N - 1$ total combinations



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

$$y = b_0 + b_1 x_1$$

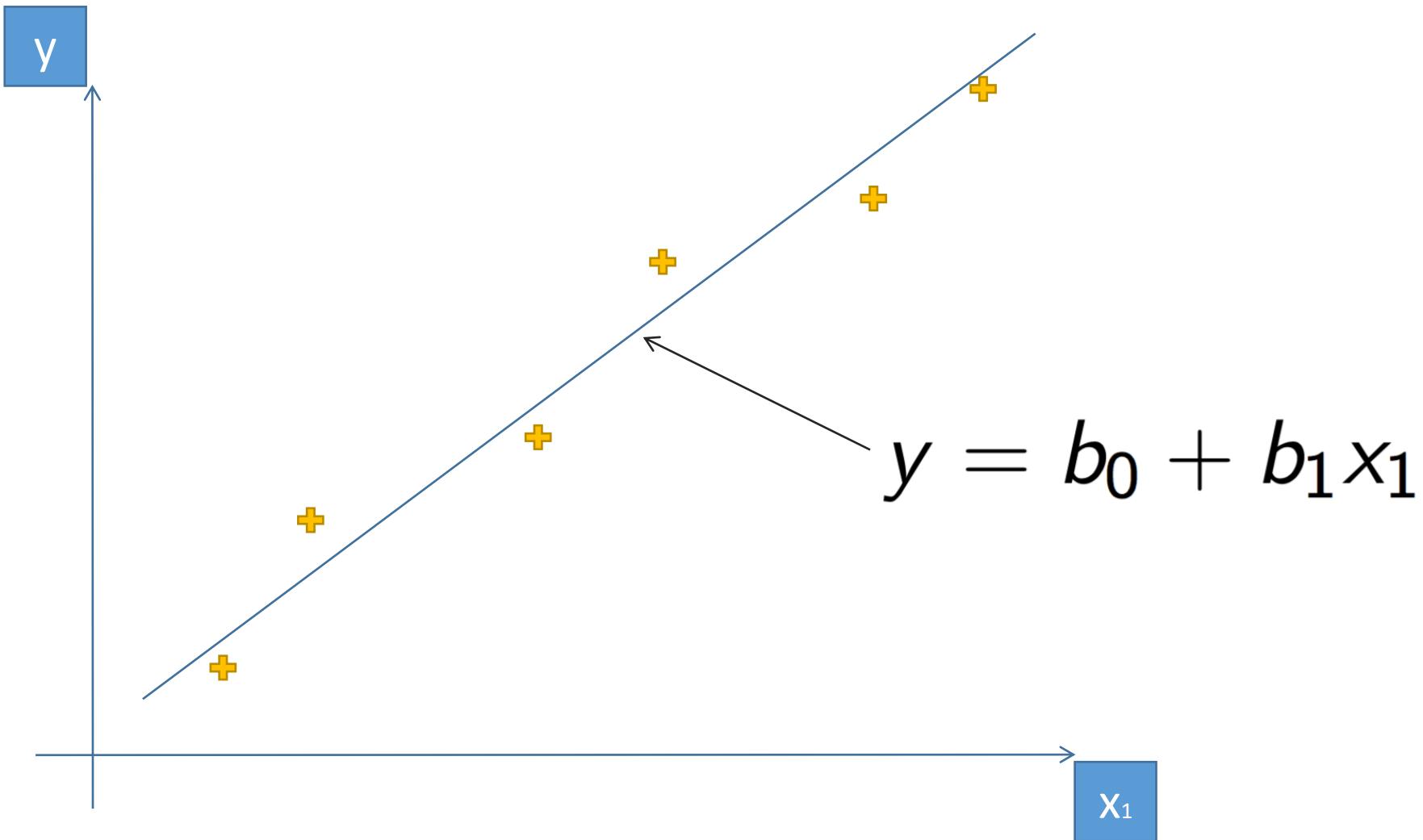
Multiple
Linear
Regression

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

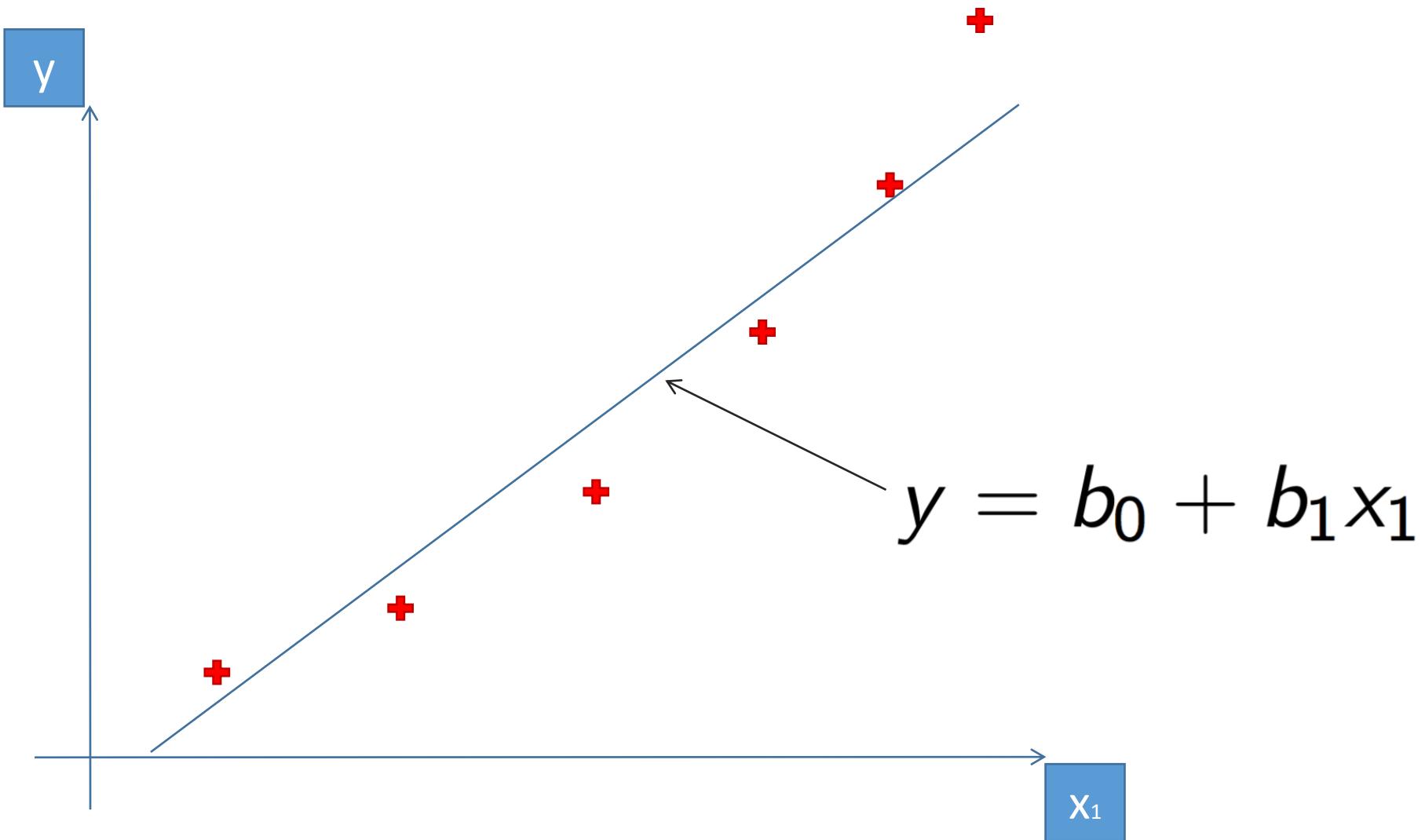
Polynomial
Linear
Regression

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

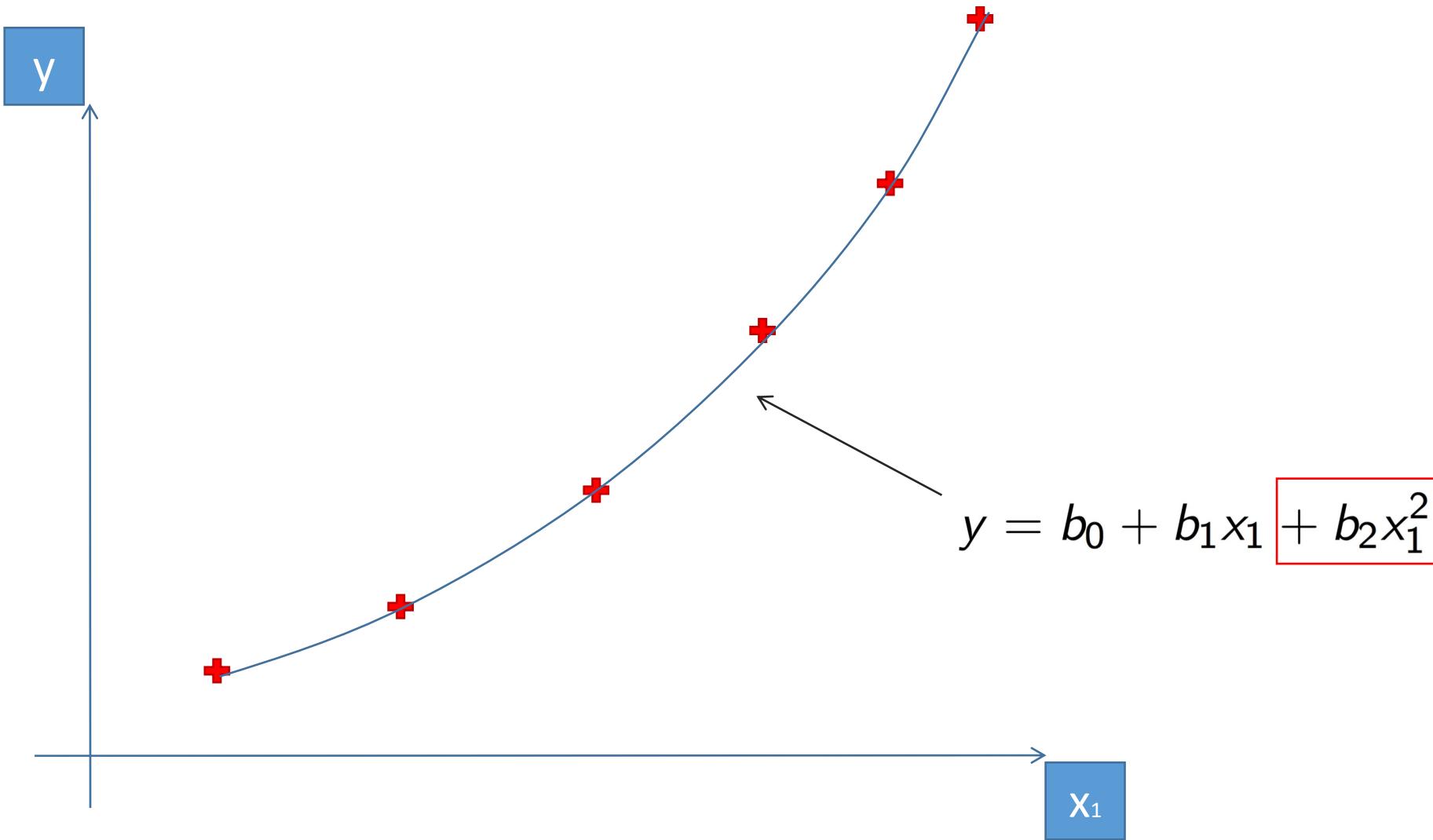
Simple Linear Regression



Simple Linear Regression



Polynomial Regression



Polynomial Regression

One Question: Why “Linear”?

Polynomial Regression

Polynomial
Linear
Regression

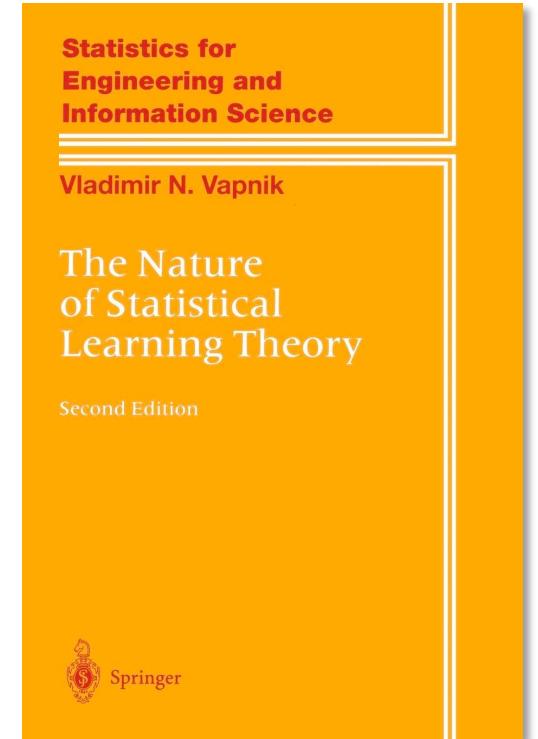
$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \dots + b_n x_1^n$$

SVR Intuition

SVR Intuition



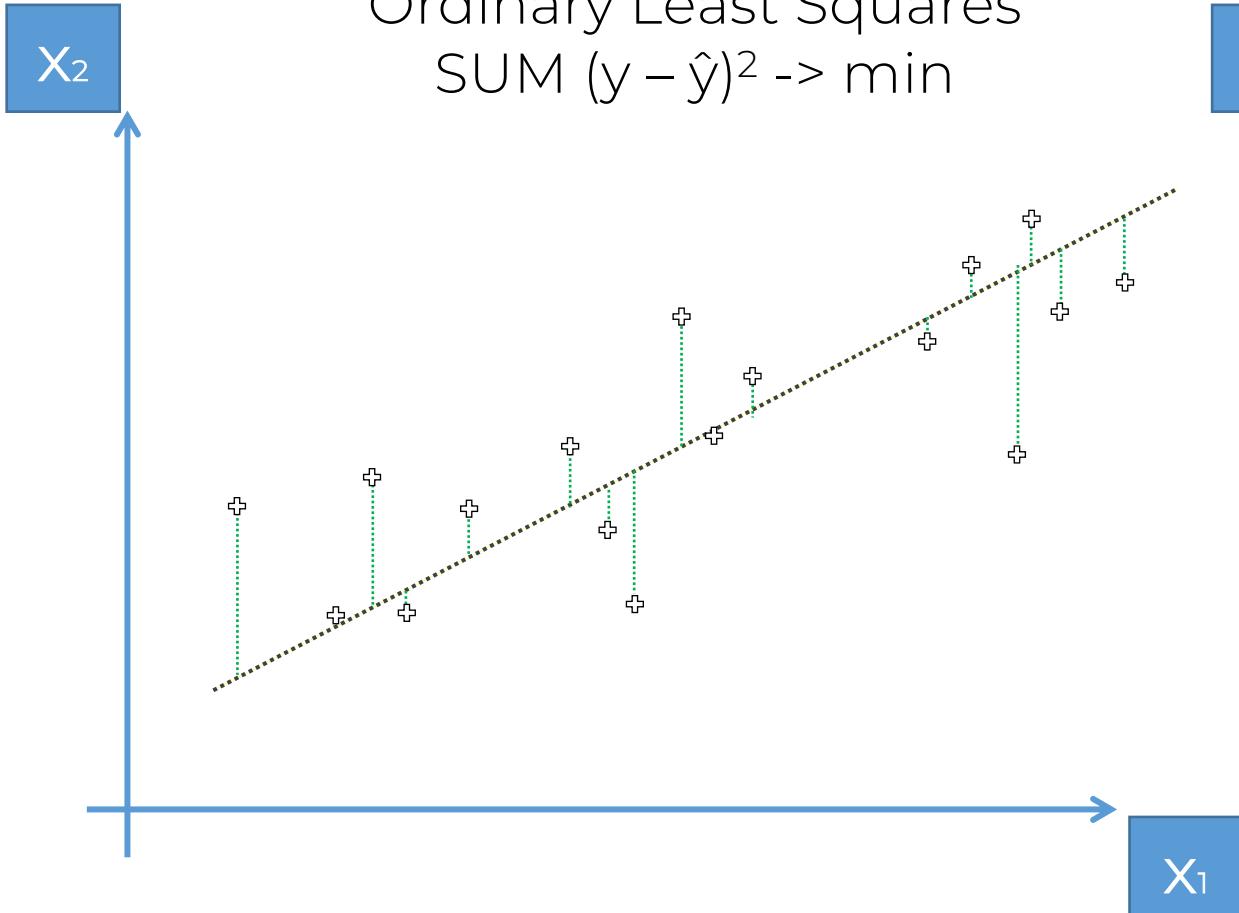
Vladimir Vapnik



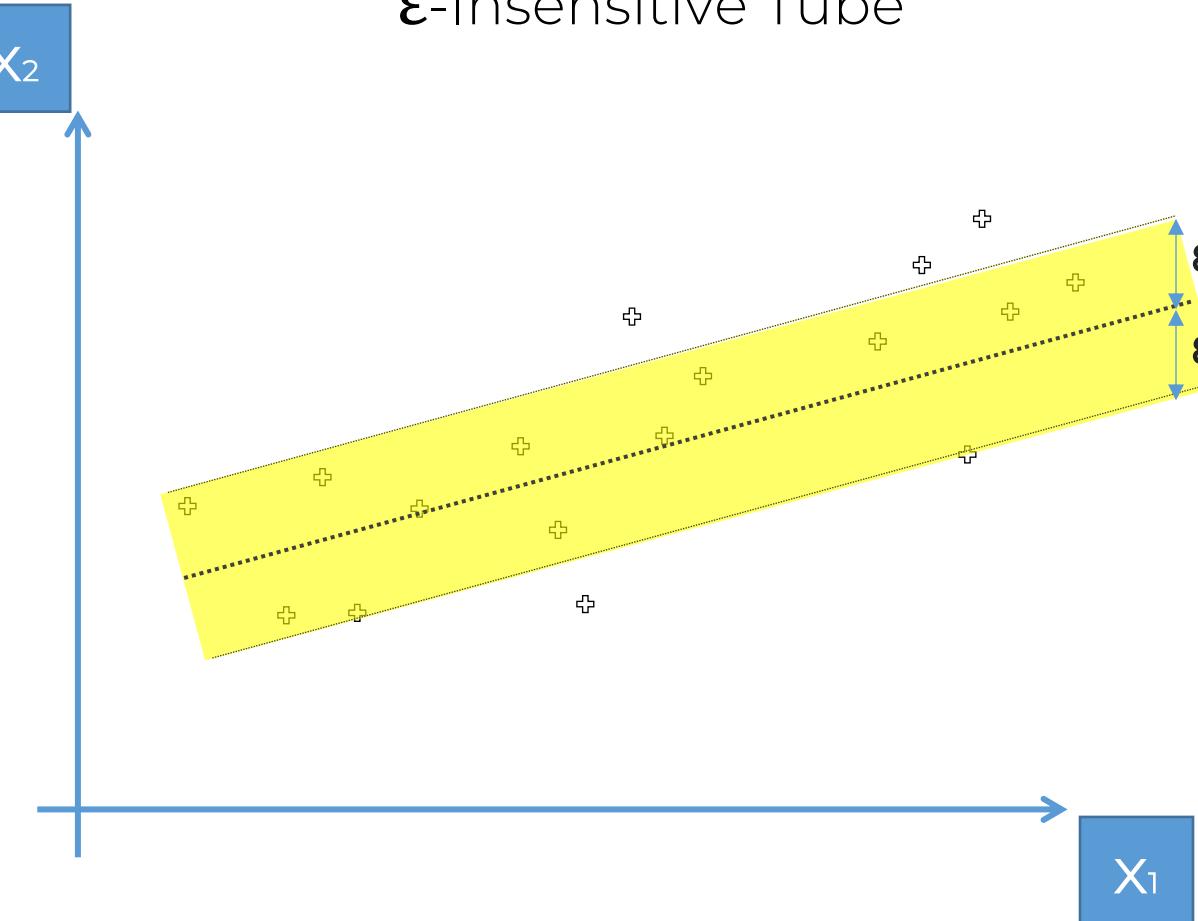
1992

SVR Intuition

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



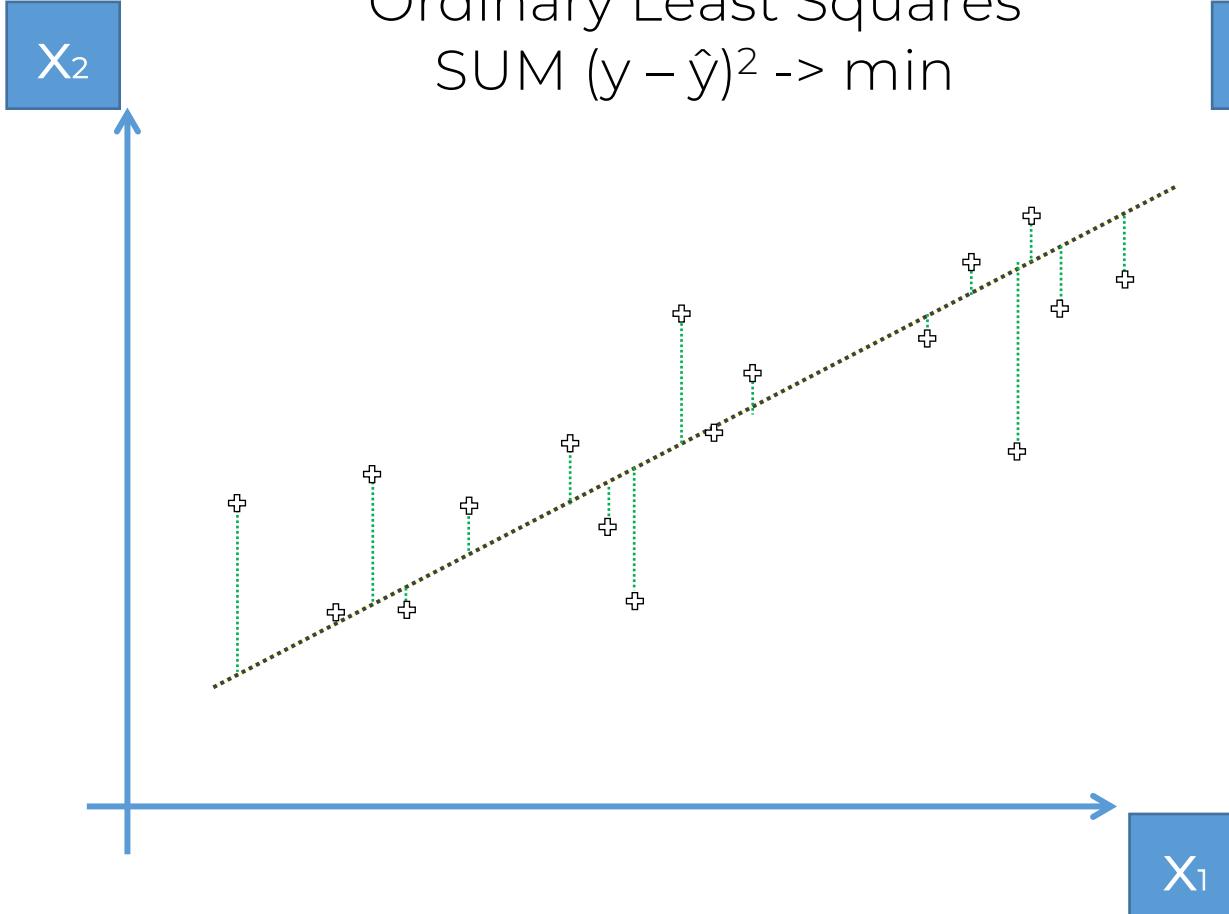
ϵ -Insensitive Tube



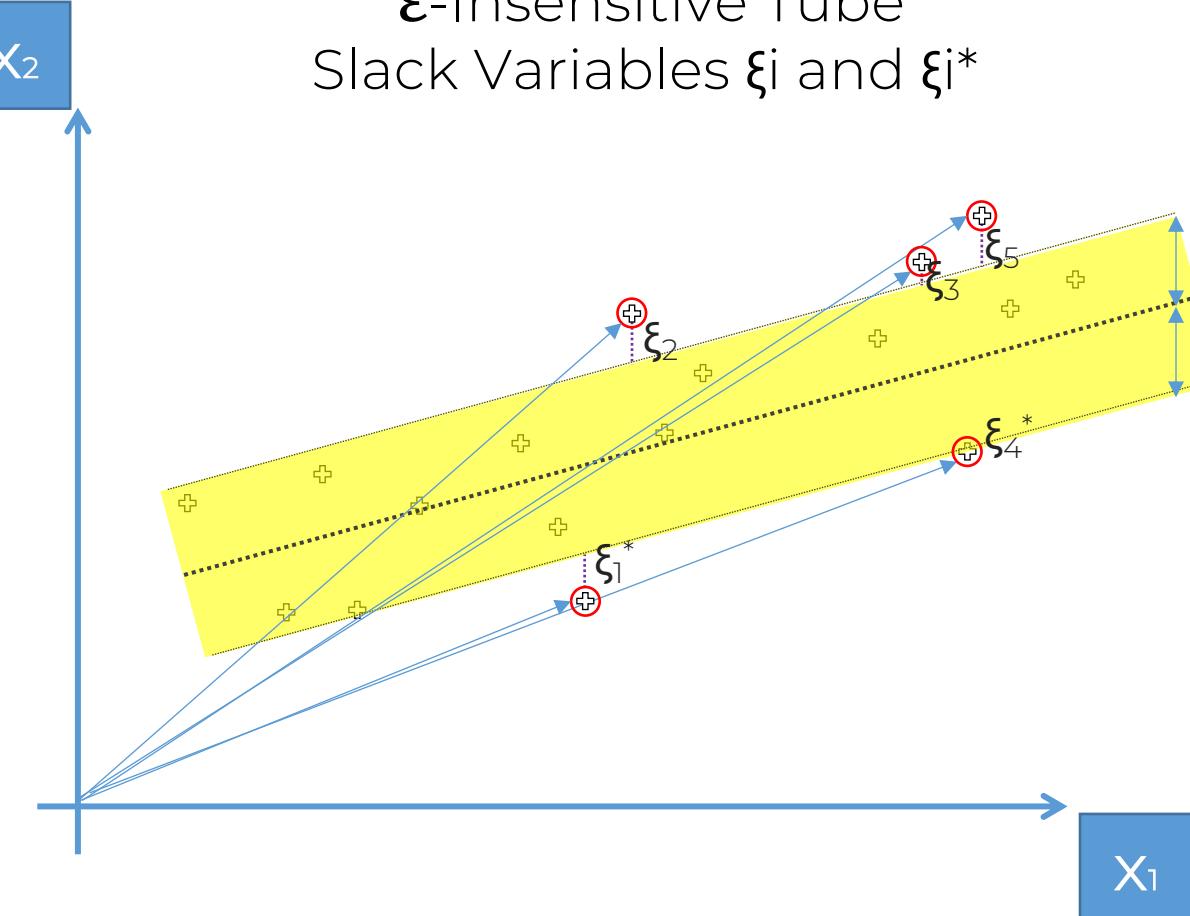
SVR Intuition

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

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



ϵ -Insensitive Tube
Slack Variables ξ_i and ξ_i^*

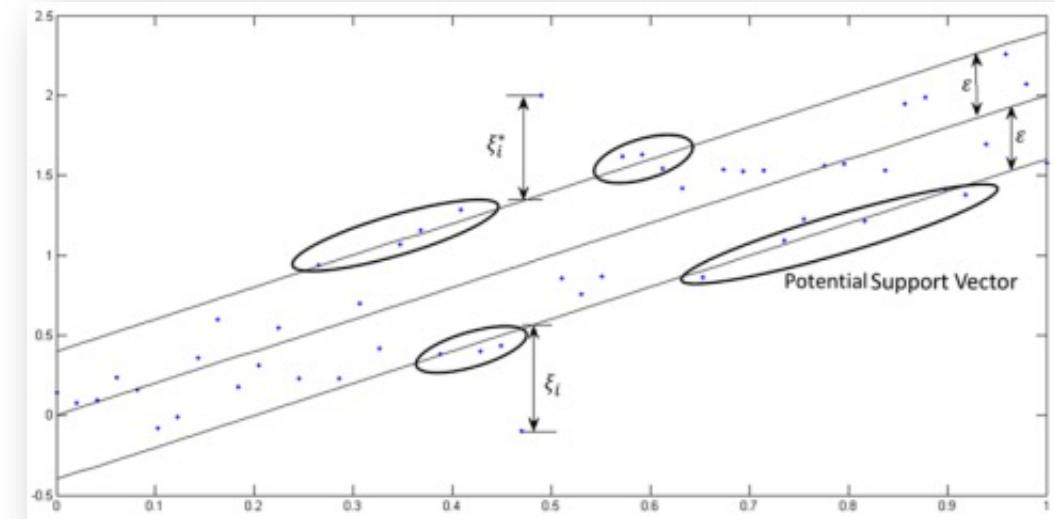


SVR Intuition

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)

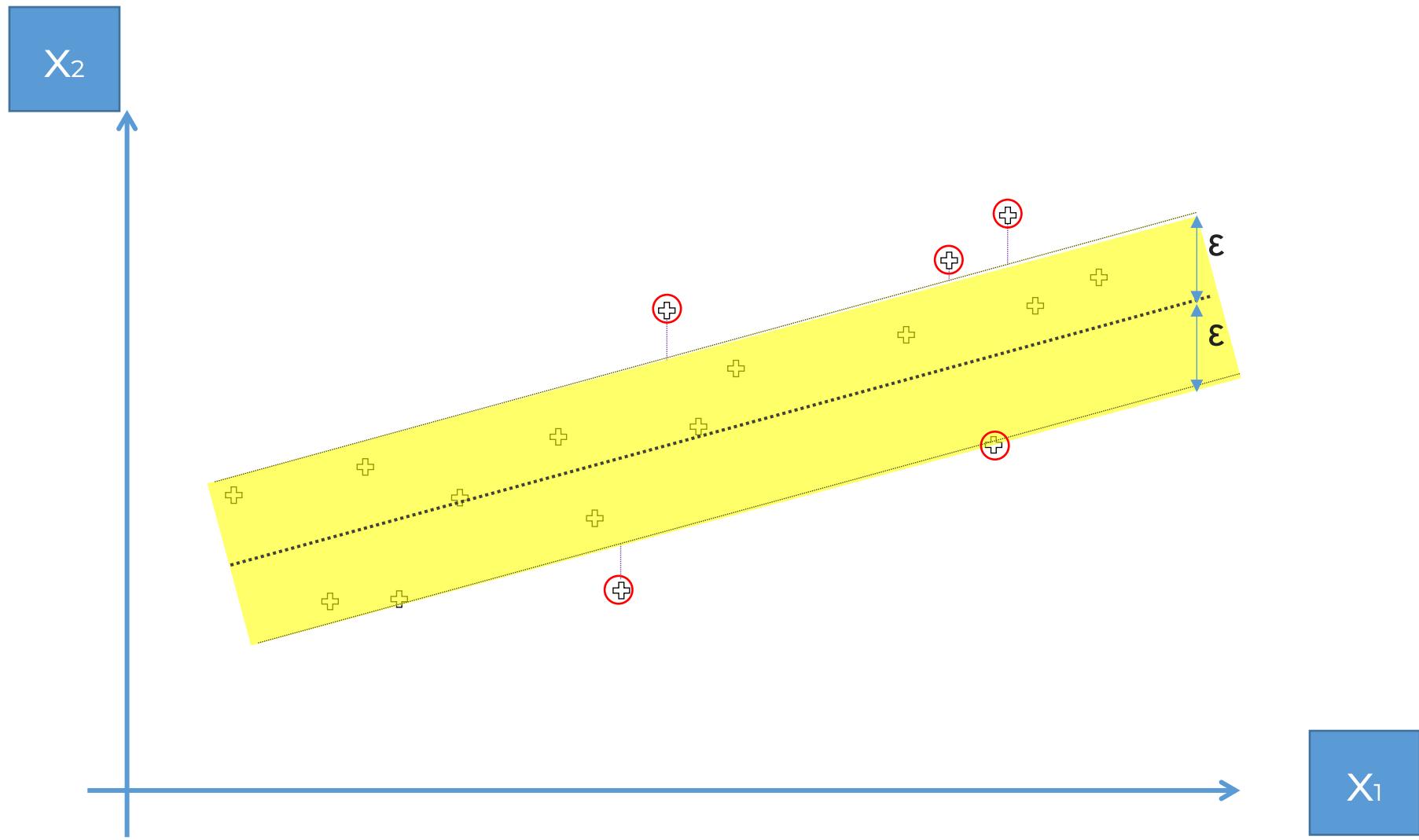


Link:

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

Heads-up about Non-Linear SVR

SVR Intuition



Copy of support_vector_regression.ipynb

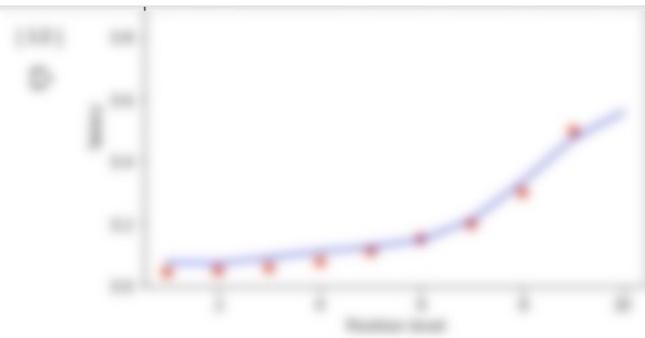
Comment Share



Files

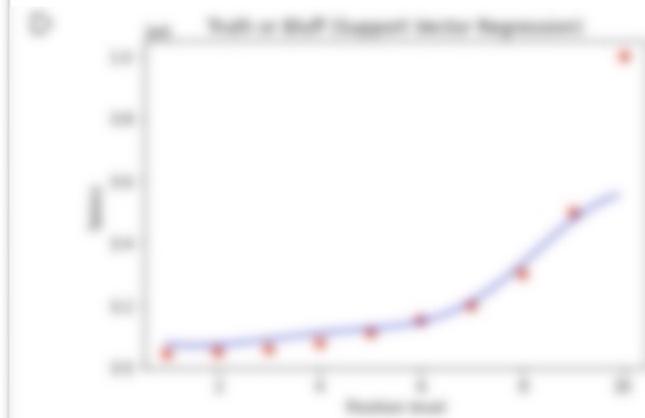
- Upload
- Refresh
- Mount Drive
- ..
- sample_data
- Position_Salaries.csv

+ Code + Text



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()
```



Position_Salaries.csv

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

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

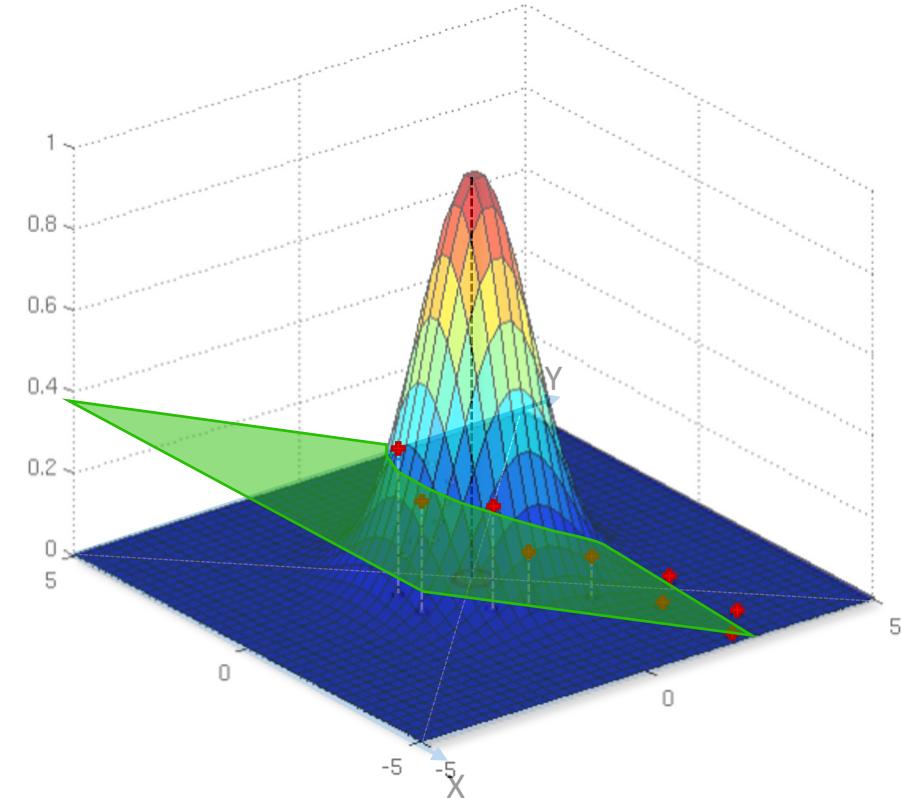
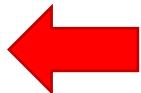
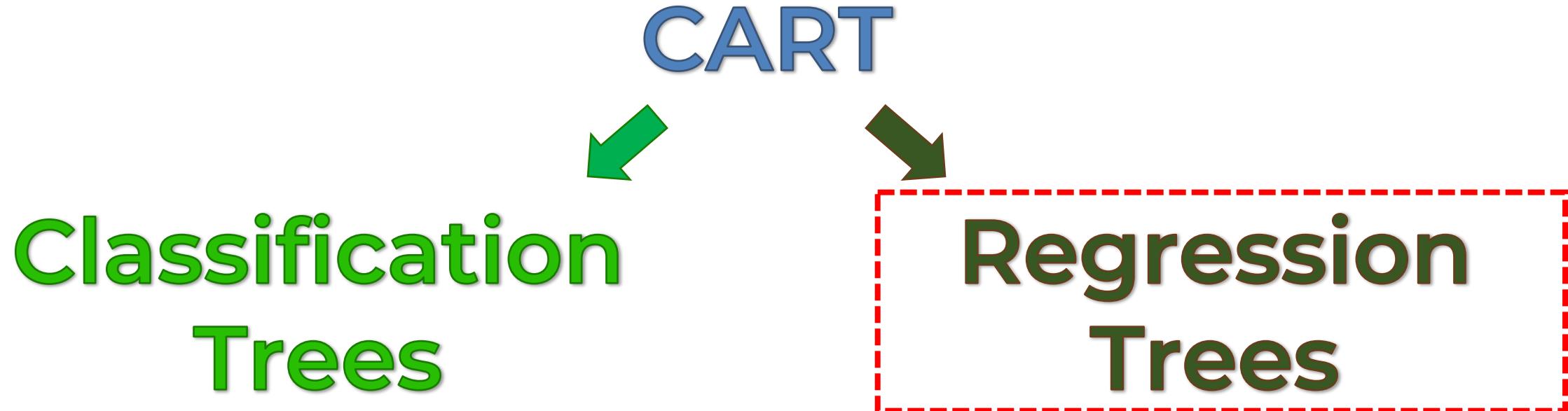


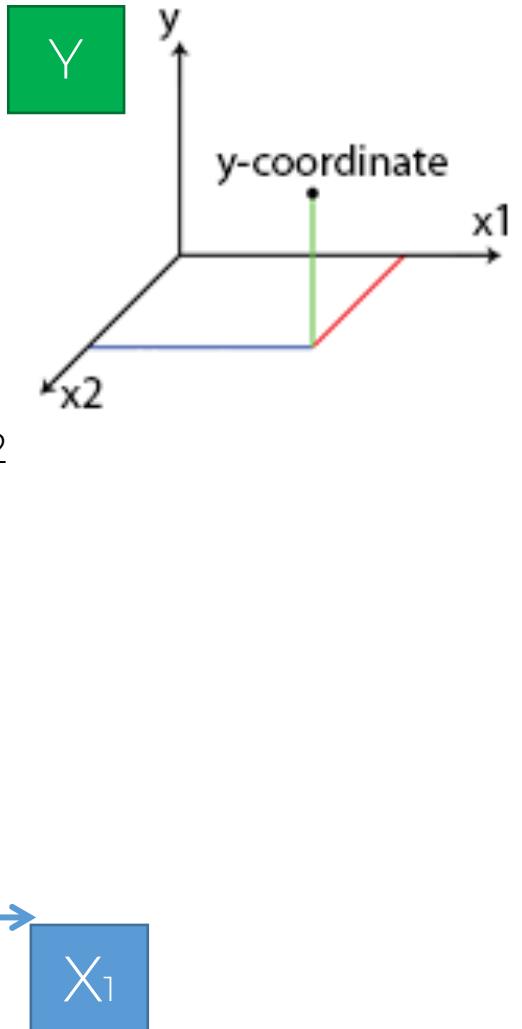
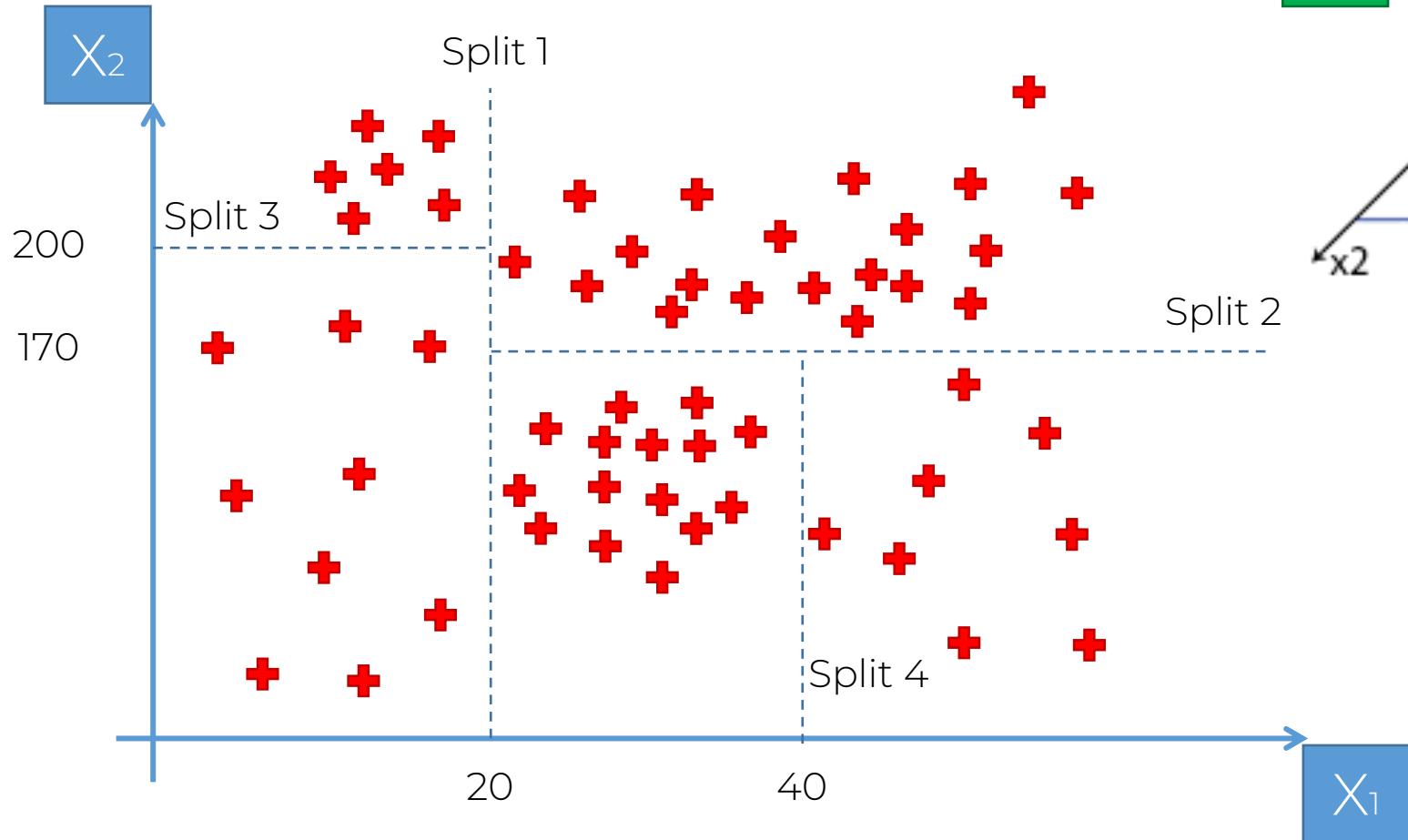
Image source: <http://www.cs.toronto.edu/~duvenaud/cookbook/index.html>

Decision Tree Intuition

Decision Tree Intuition



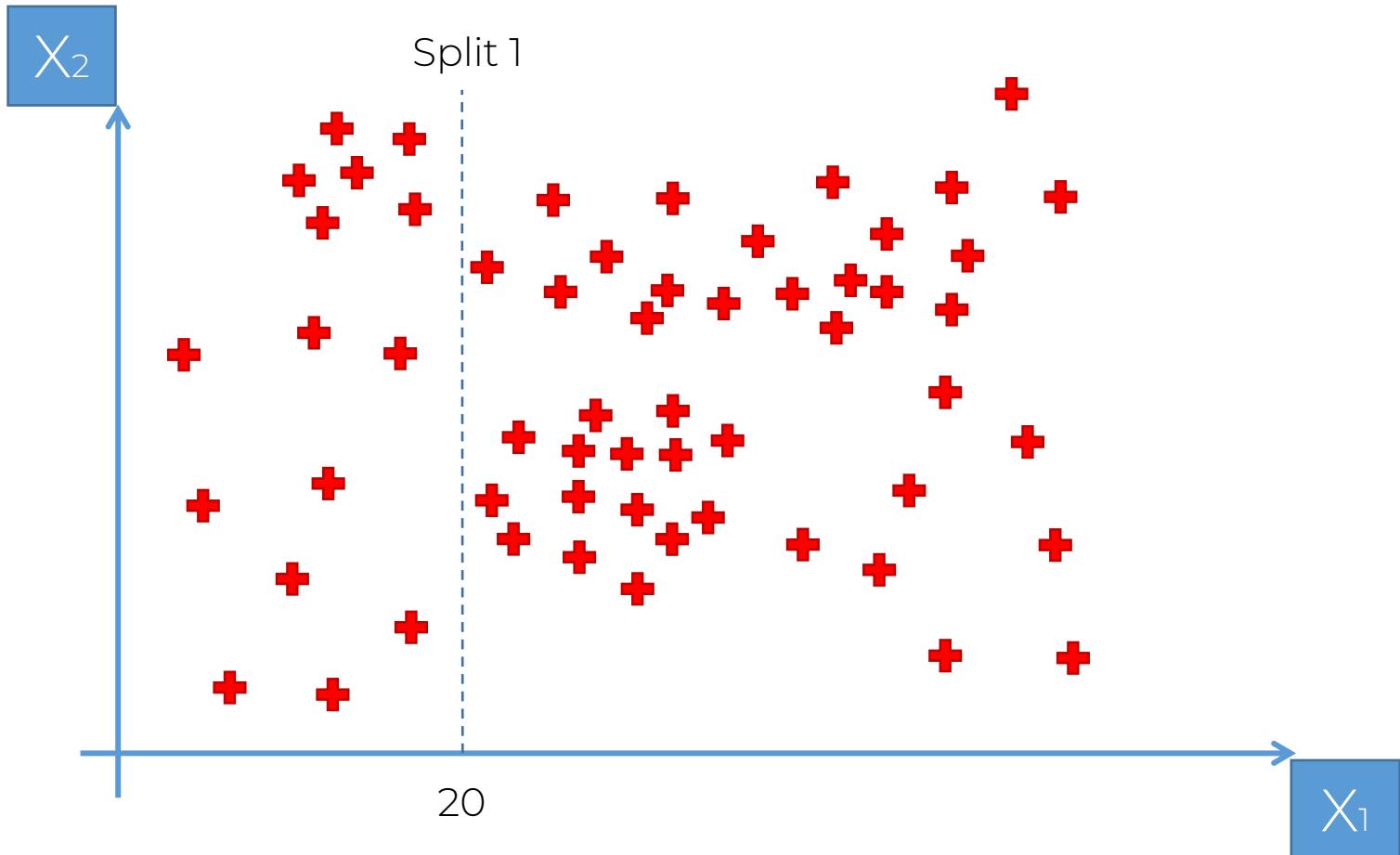
Decision Tree Intuition



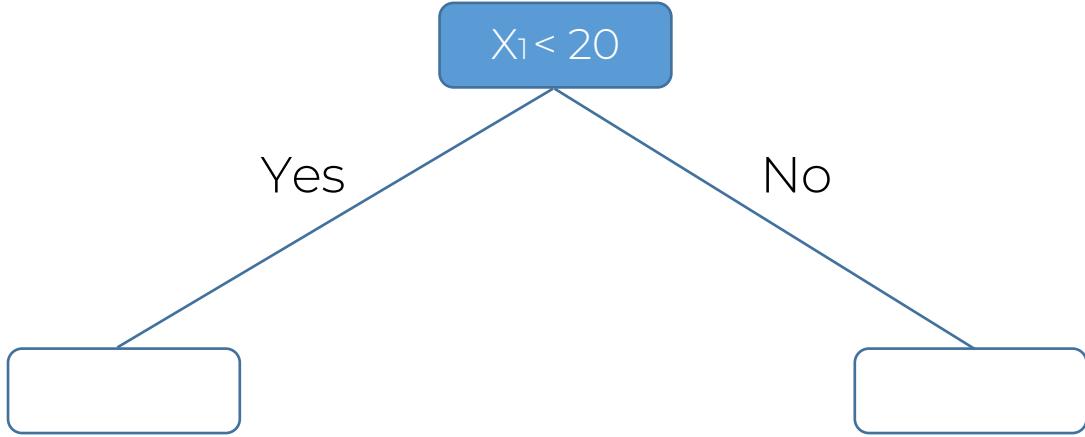
Decision Tree Intuition

Rewind...

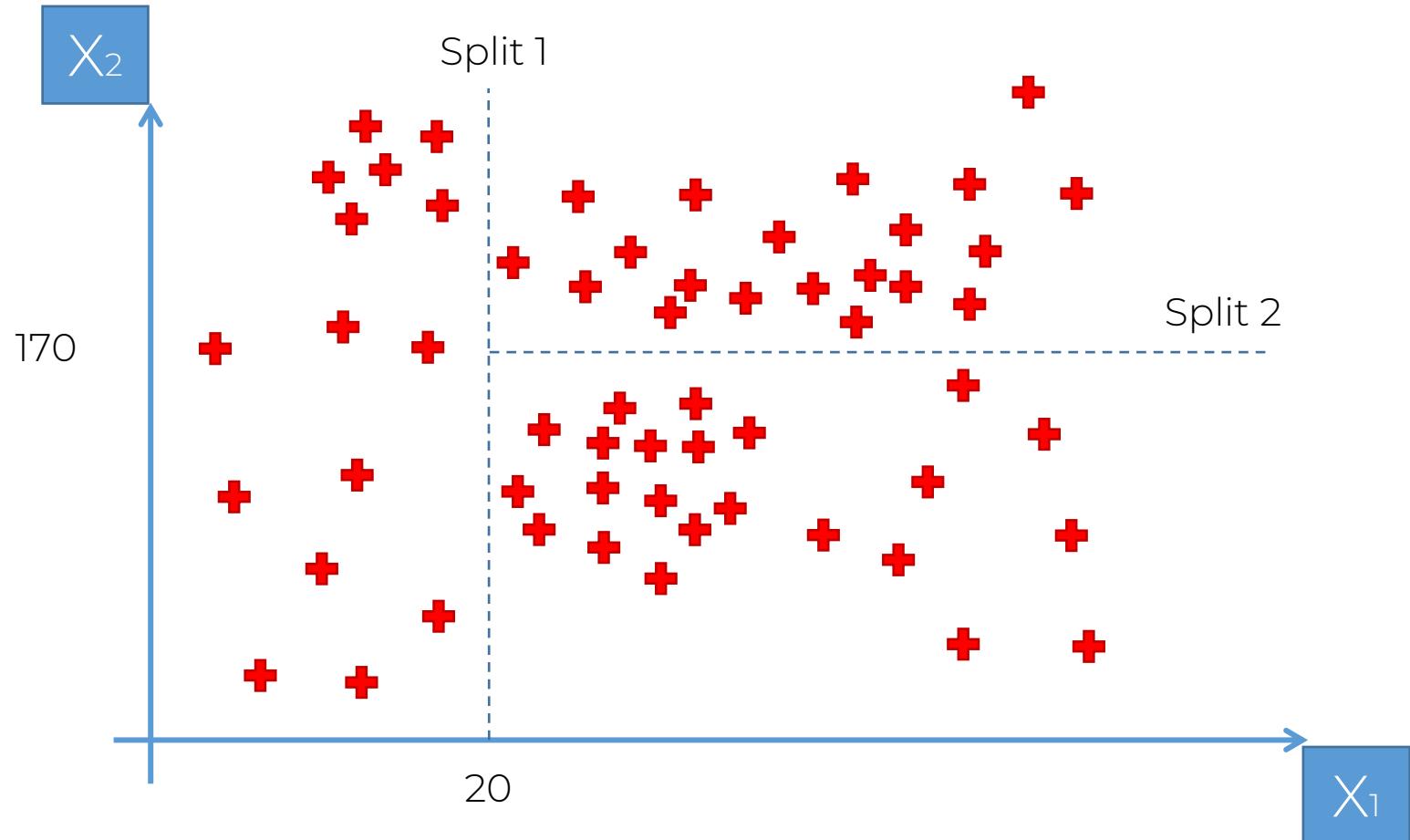
Decision Tree Intuition



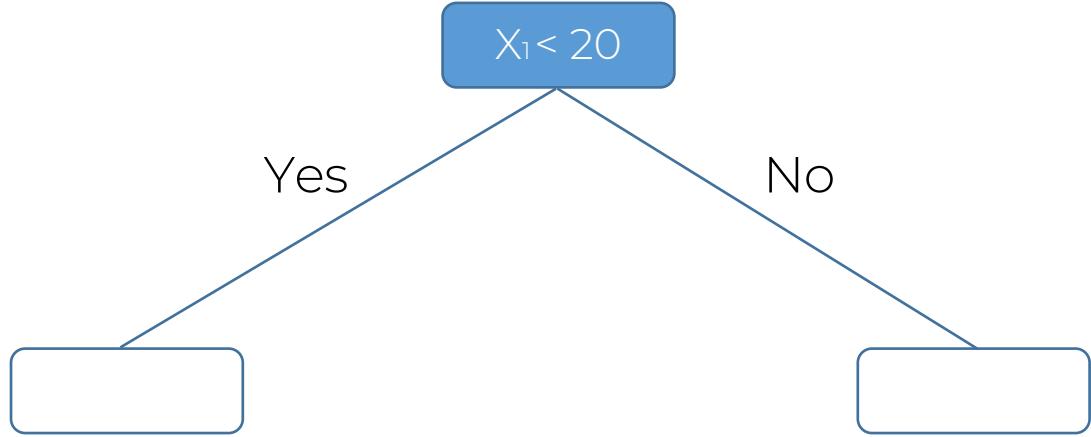
Decision Tree Intuition



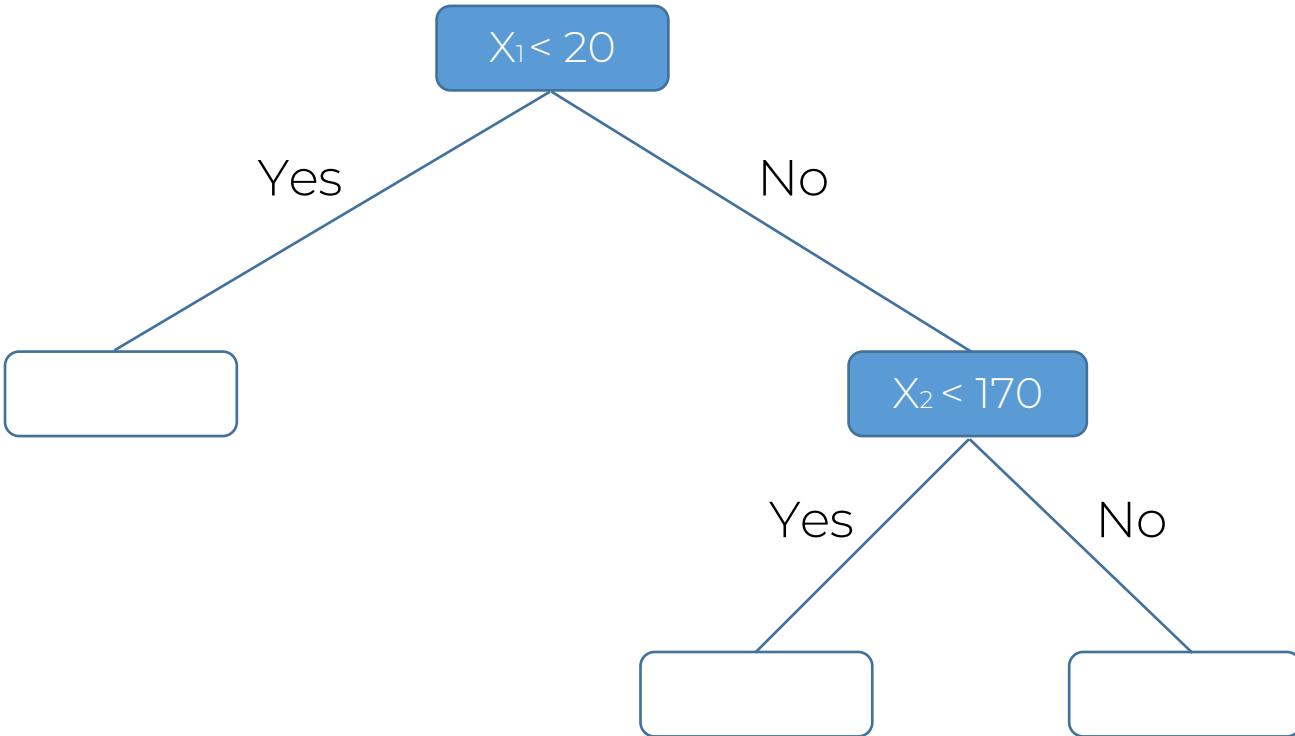
Decision Tree Intuition



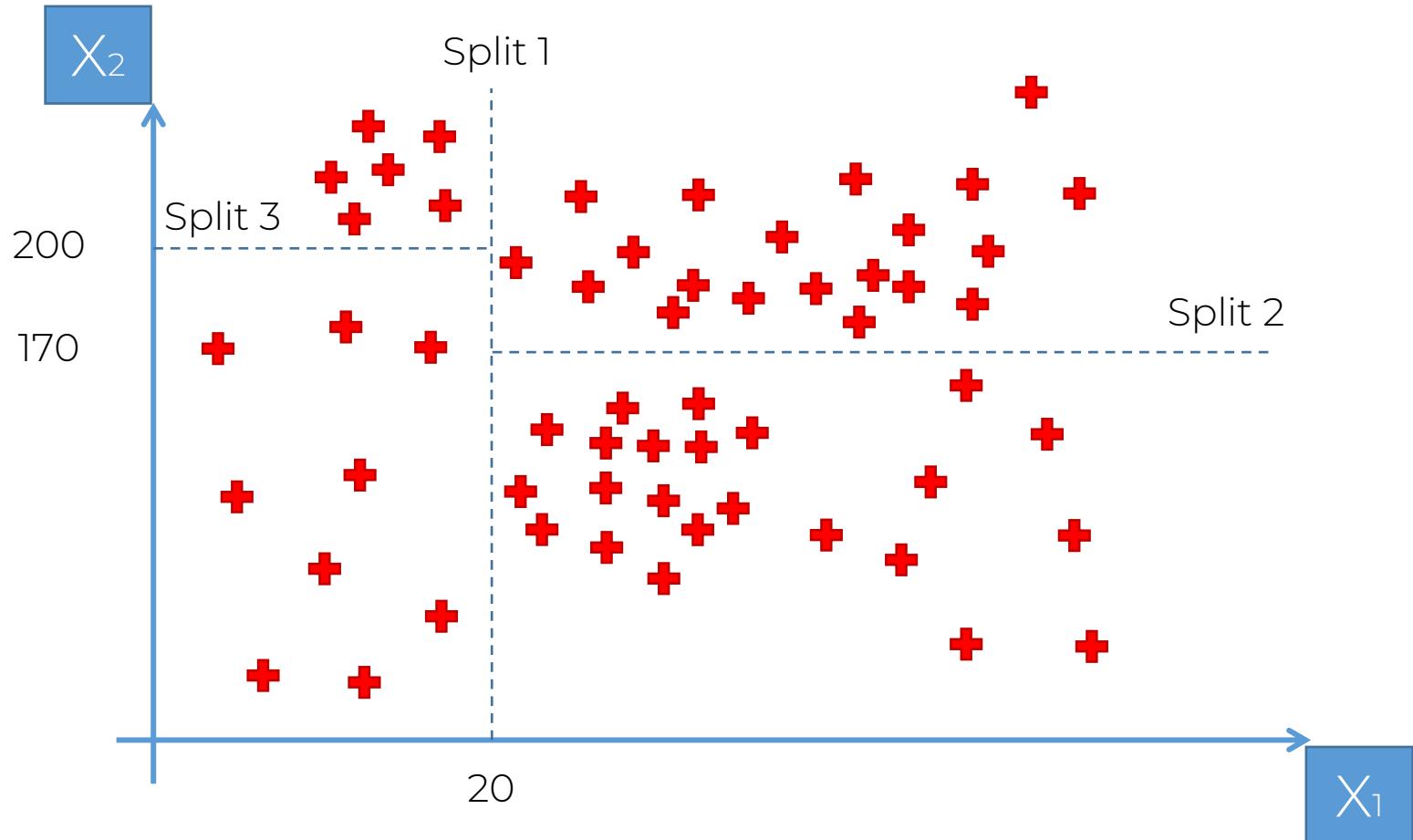
Decision Tree Intuition



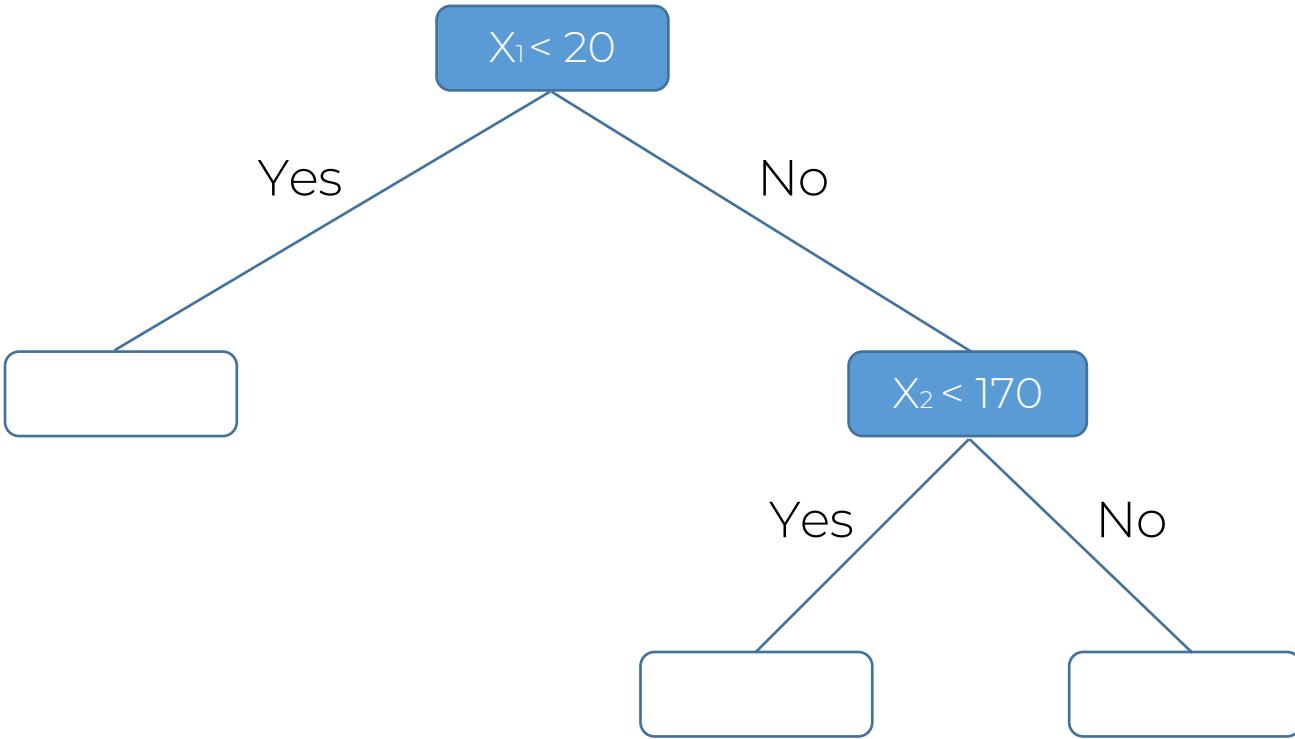
Decision Tree Intuition



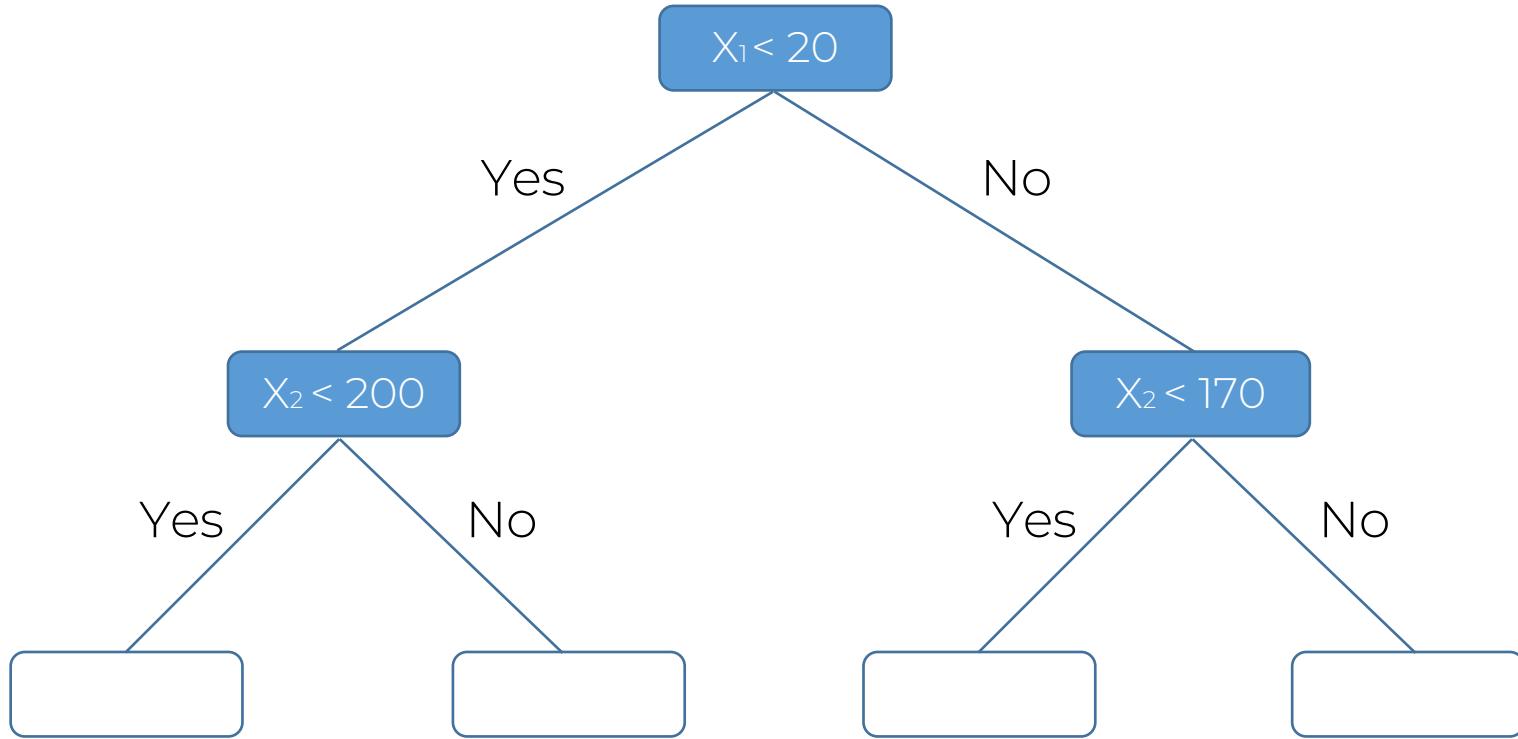
Decision Tree Intuition



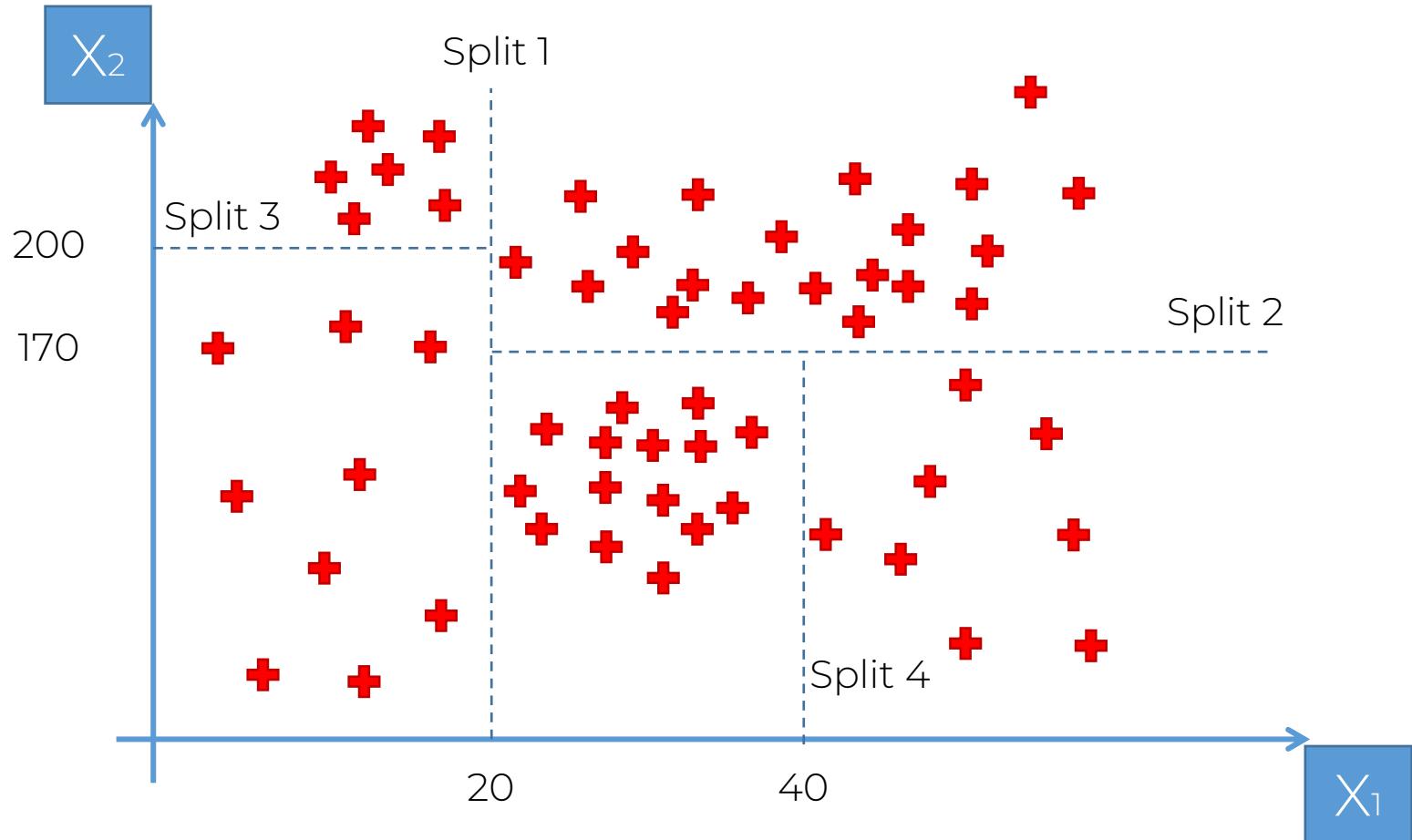
Decision Tree Intuition



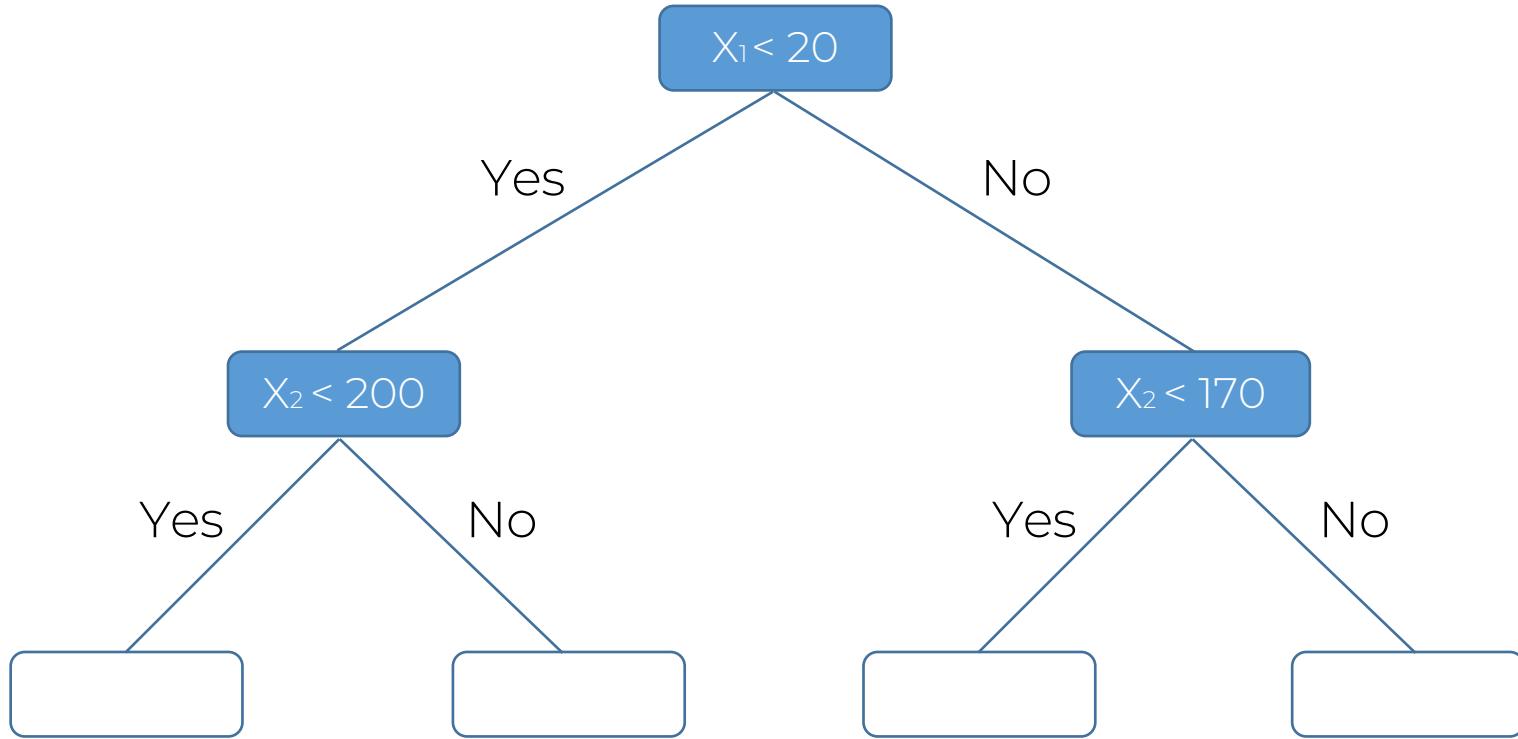
Decision Tree Intuition



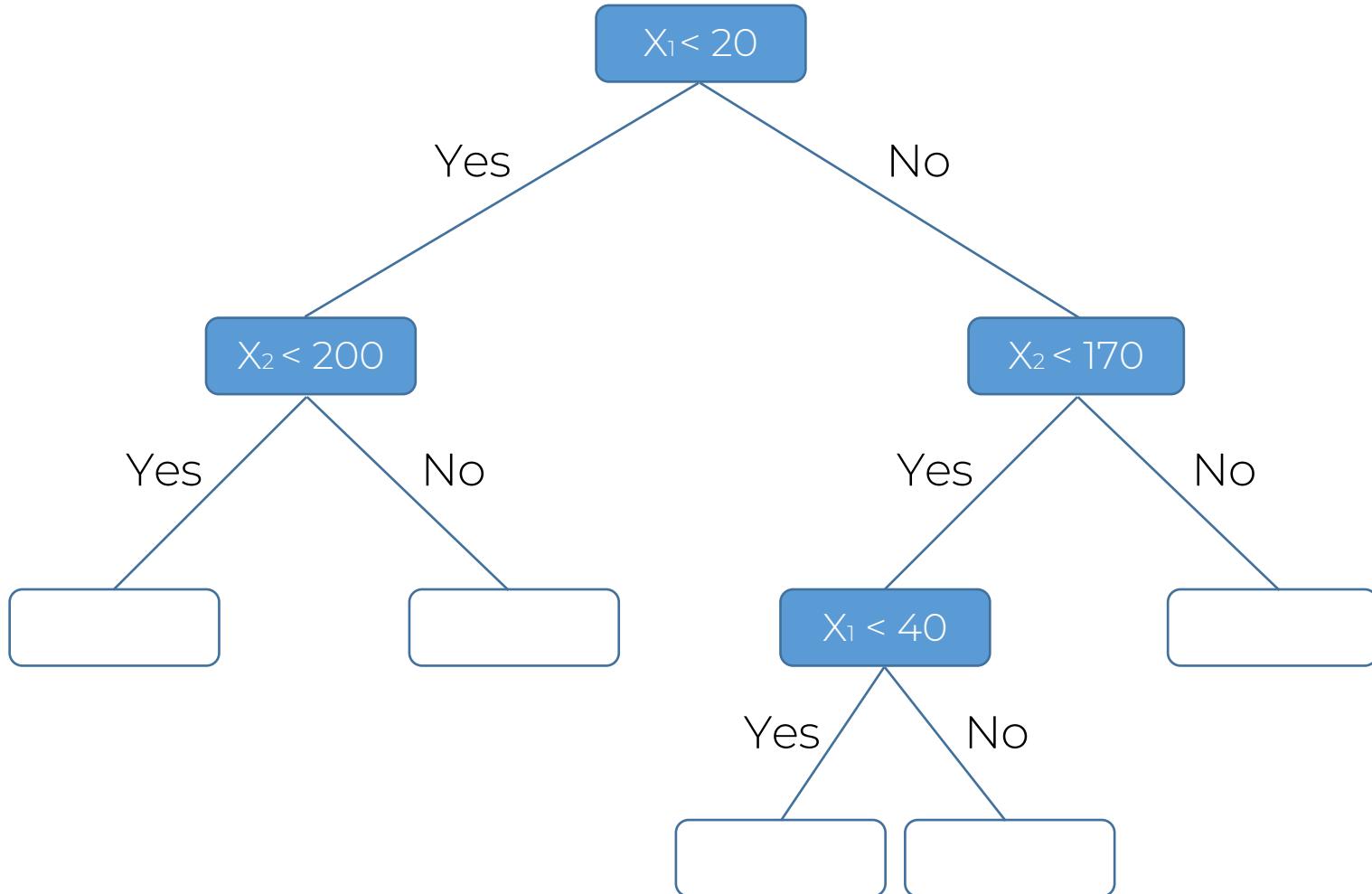
Decision Tree Intuition



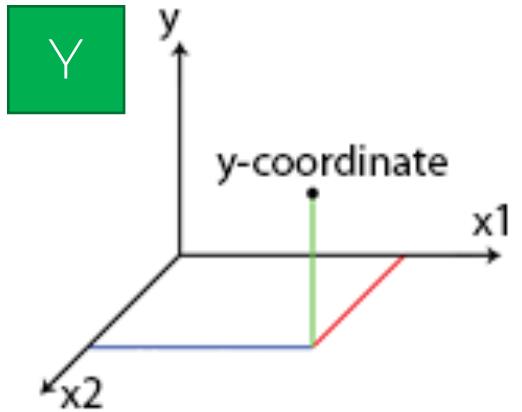
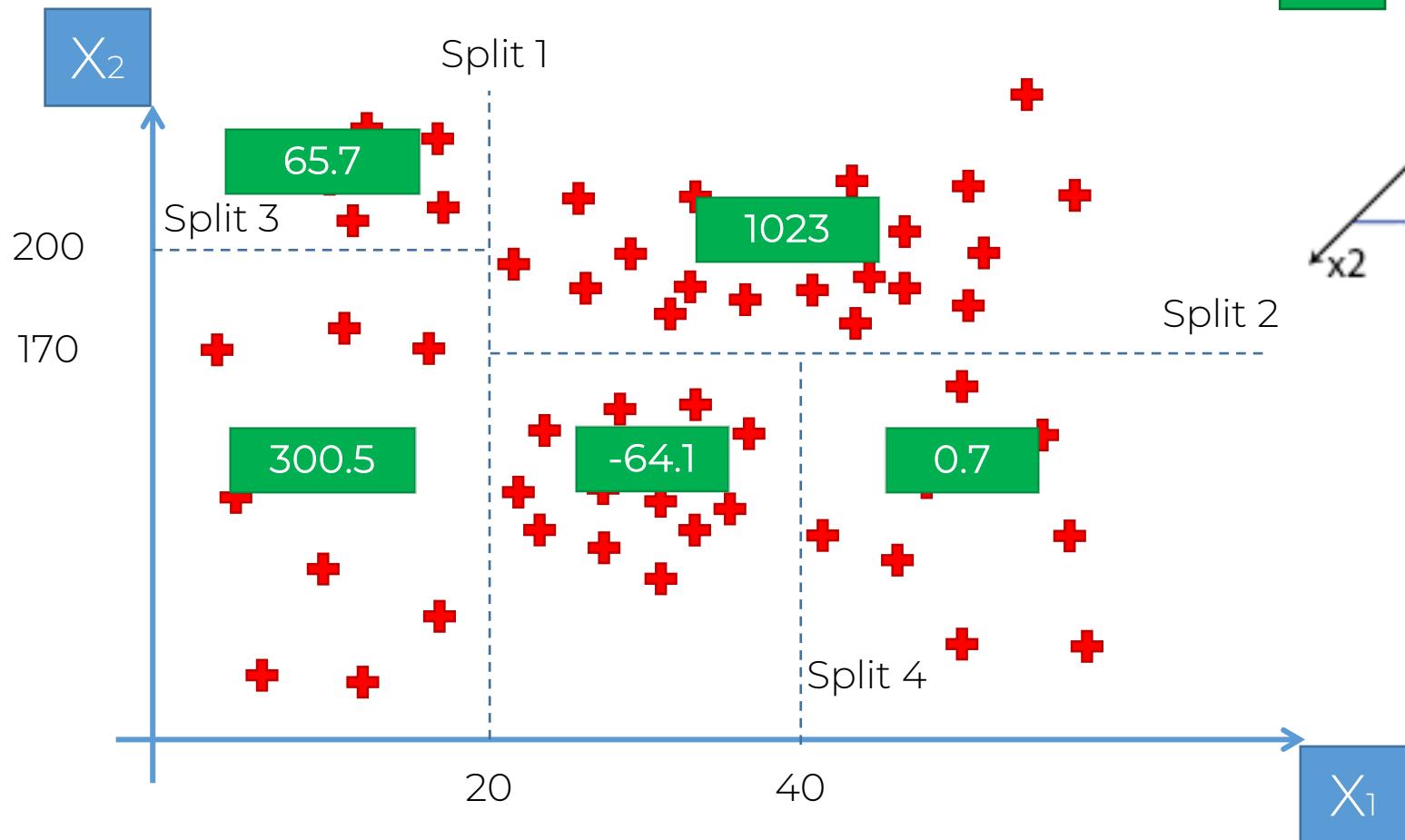
Decision Tree Intuition



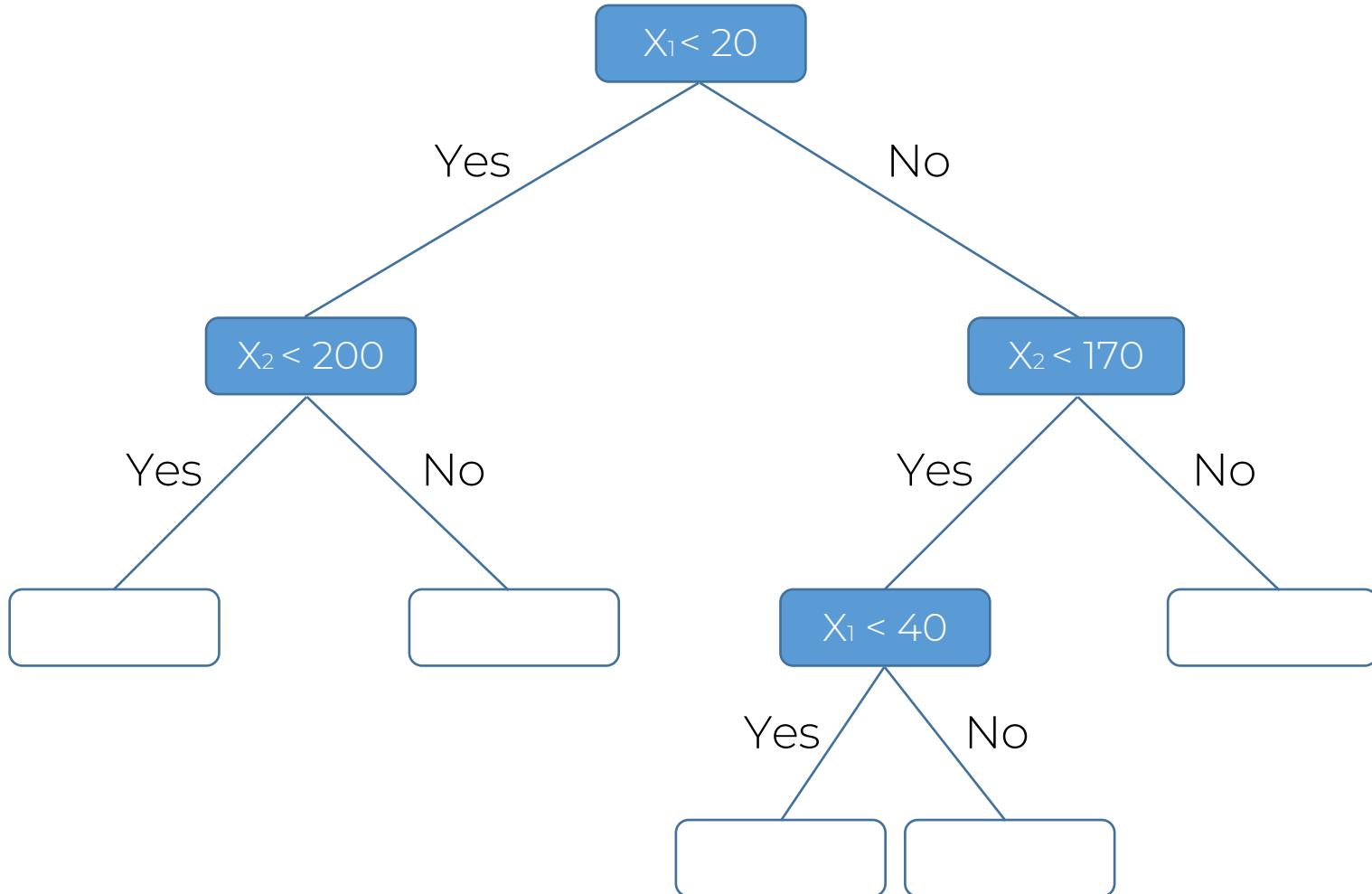
Decision Tree Intuition



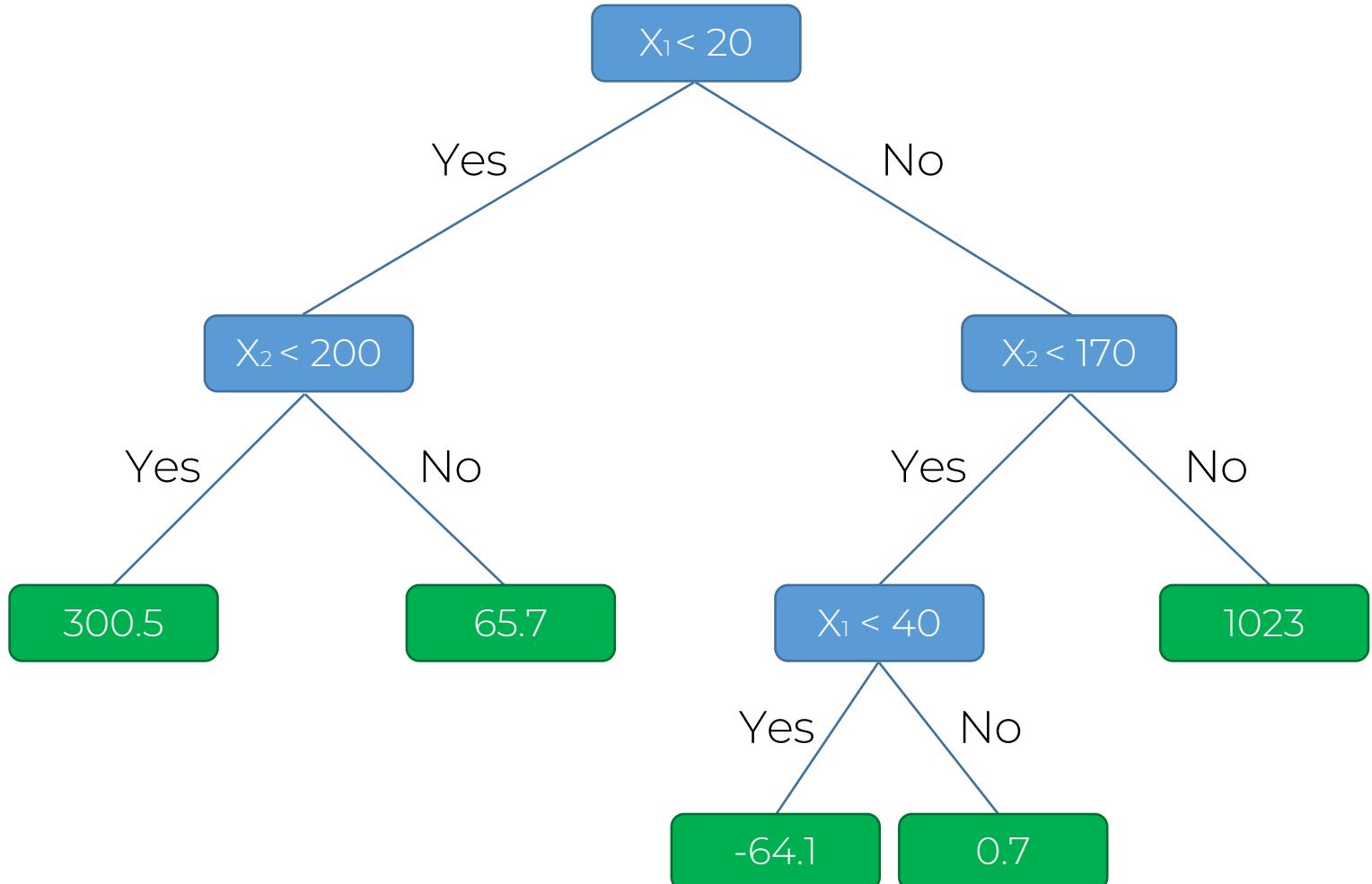
Decision Tree Intuition



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



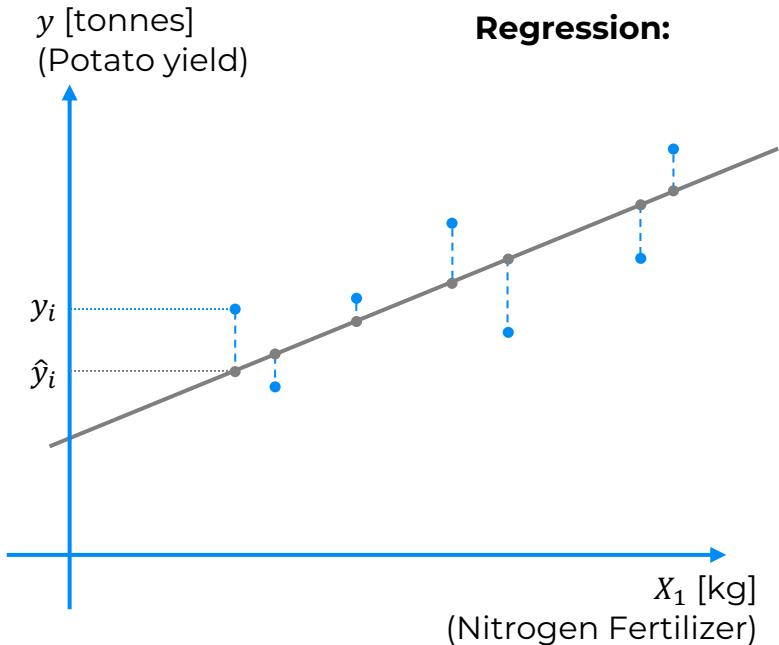
STEP 4: For a new data point, make each one of your Ntree trees predict the value of Y to for the data point in question, and assign the new data point the average across all of the predicted Y values.

R Squared

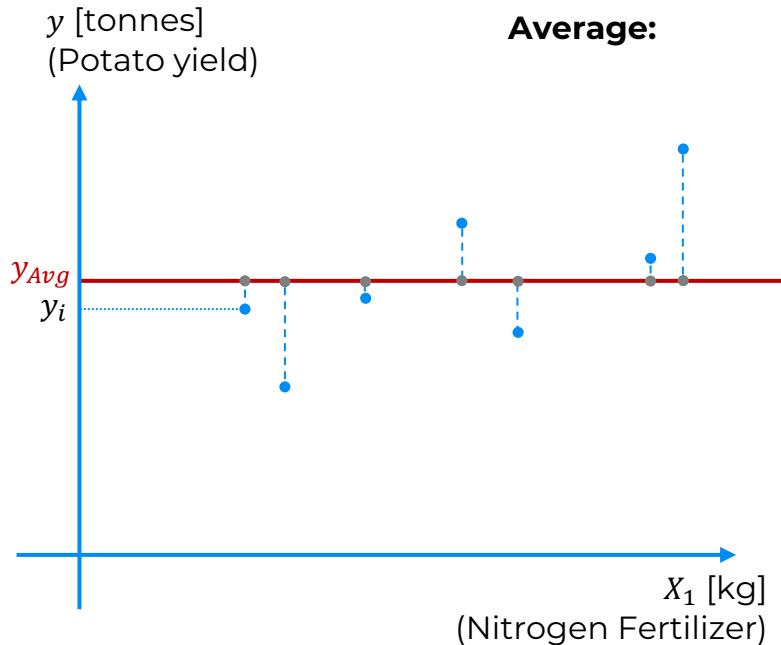




R Squared



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



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

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

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

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_{tot} doesn't change

SS_{res} will decrease or stay the same

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

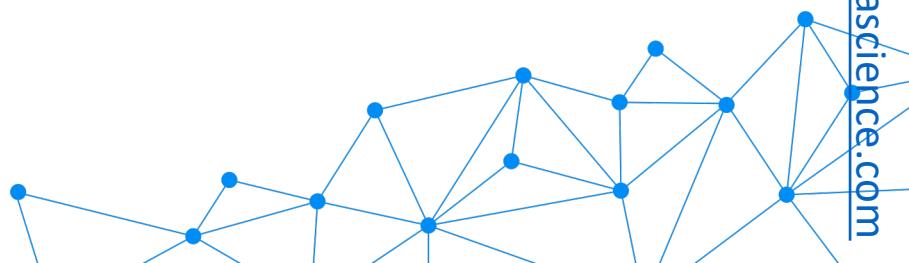
(This is because of Ordinary Least Squares: $SS_{res} \rightarrow \text{Min}$)

Solution:

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

k – number of independent variables

n – sample size



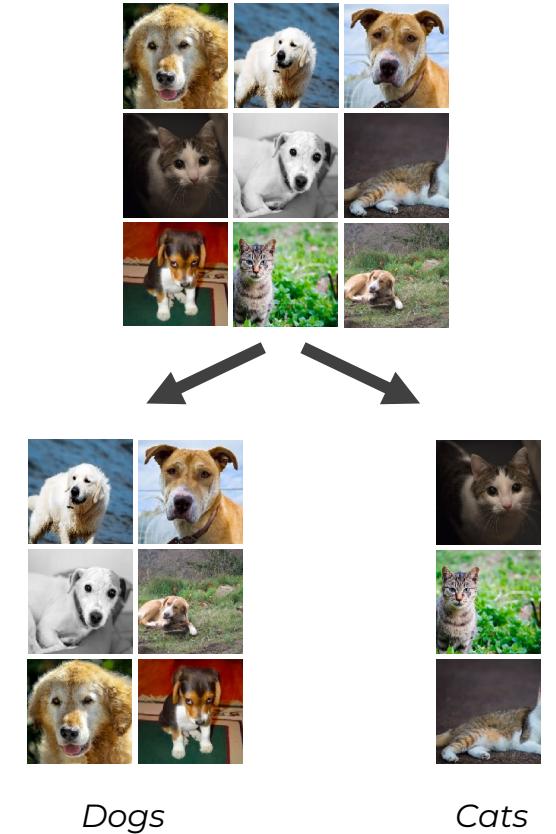
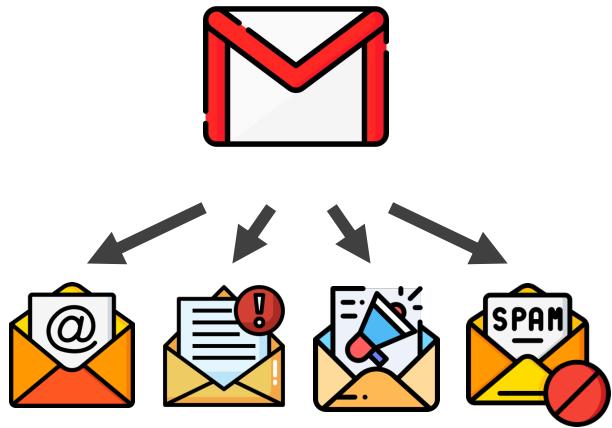
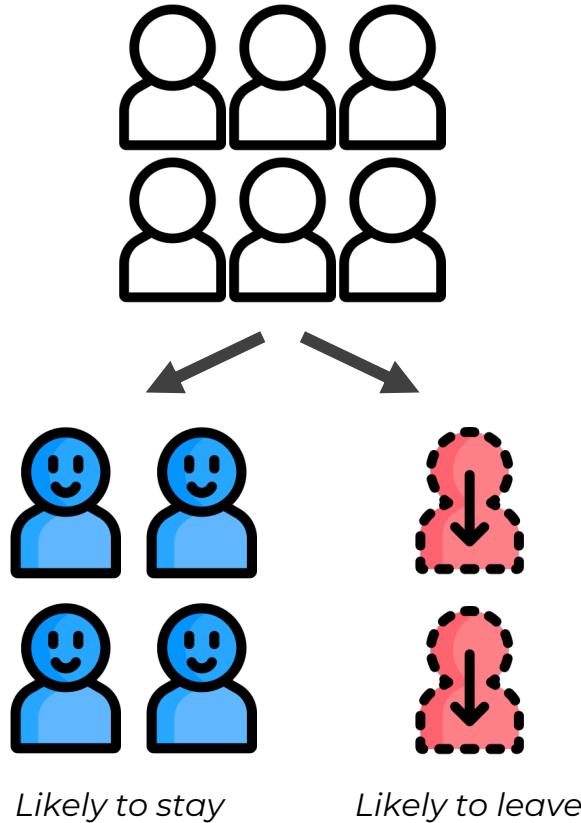


Classification

What is Classification?



Classification: a Machine Learning technique to identify the category of new observations based on training data.





Logistic Regression

Logistic Regression

Logistic regression: predict a categorical dependent variable from a number of independent variables.

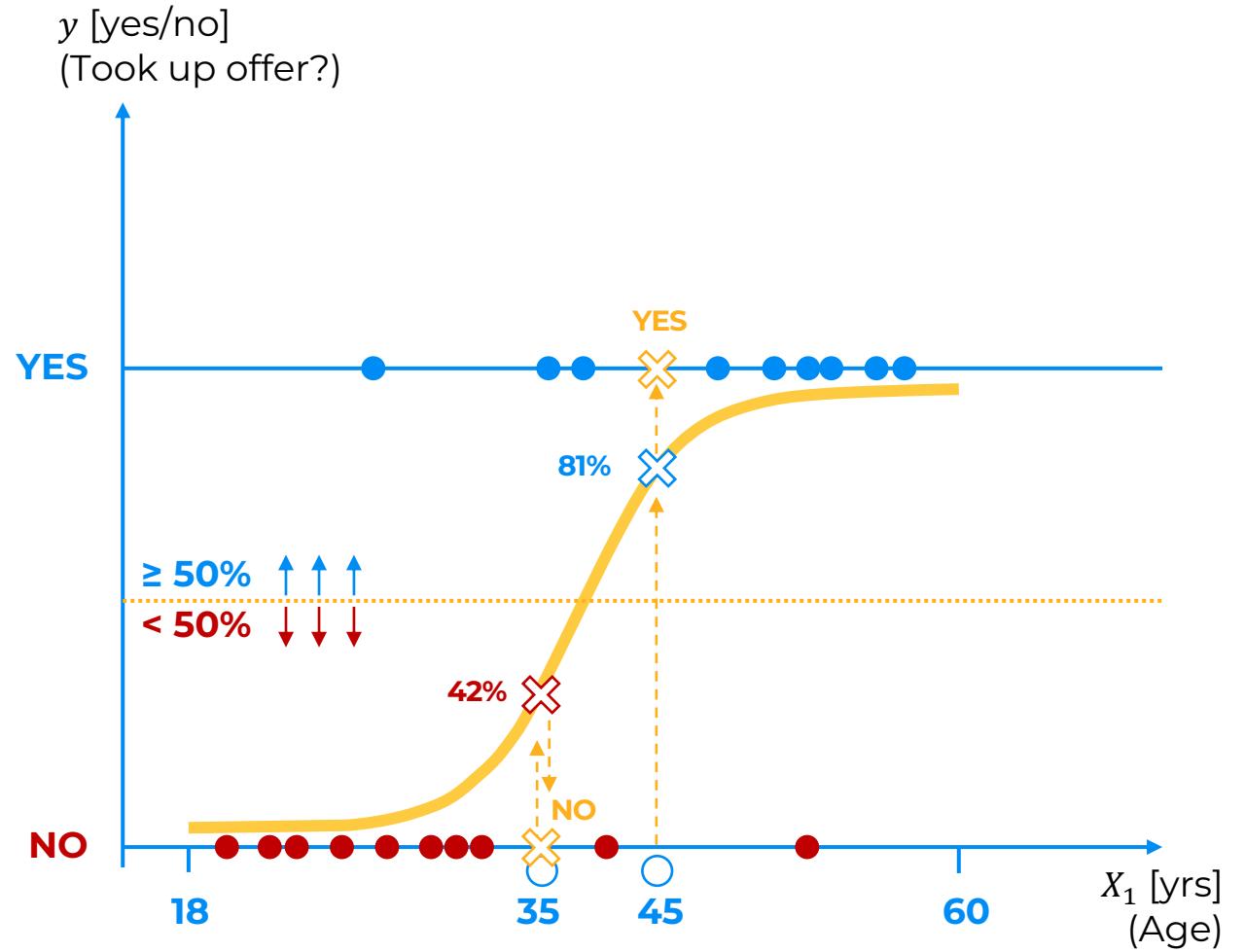


Will purchase
health insurance:
Yes / No



Age

$$\ln \frac{p}{1-p} = b_0 + b_1 X_1$$



Logistic Regression



NOT FOR DISTRIBUTION

© SUPERDATASCIENCE

www.superdatascience.com

Will purchase
health insurance:
Yes / No



Age



Income

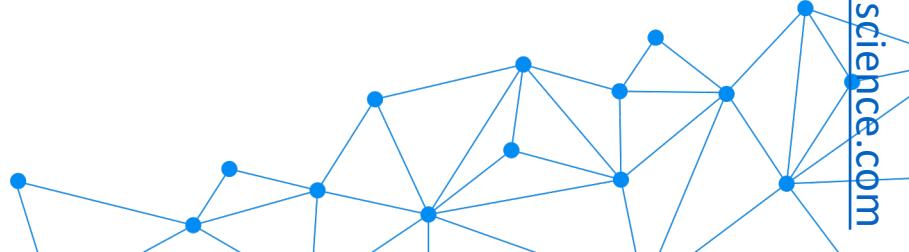


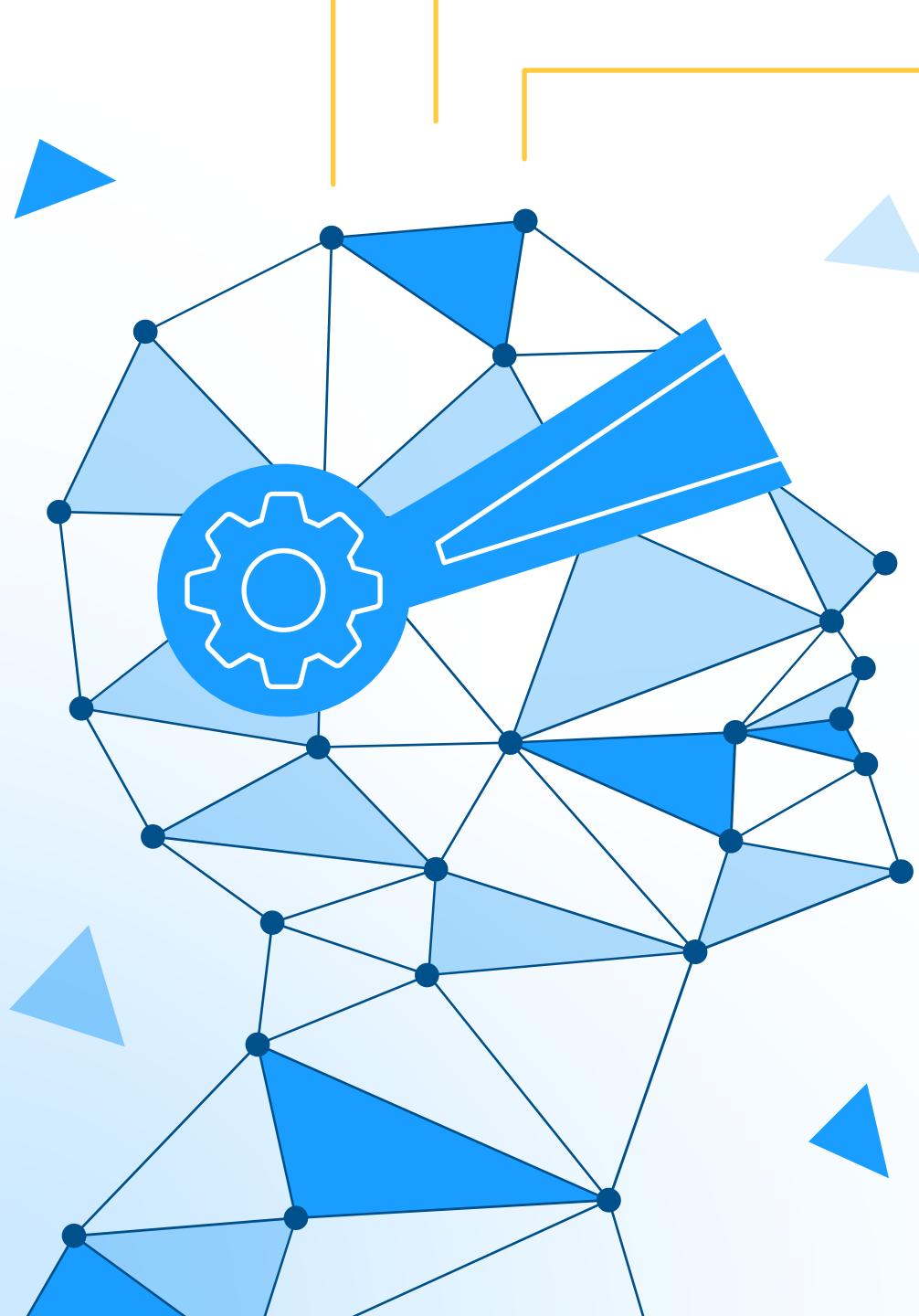
Level of
Education



Family or
Single

$$\ln \frac{p}{1-p} = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 X_4$$





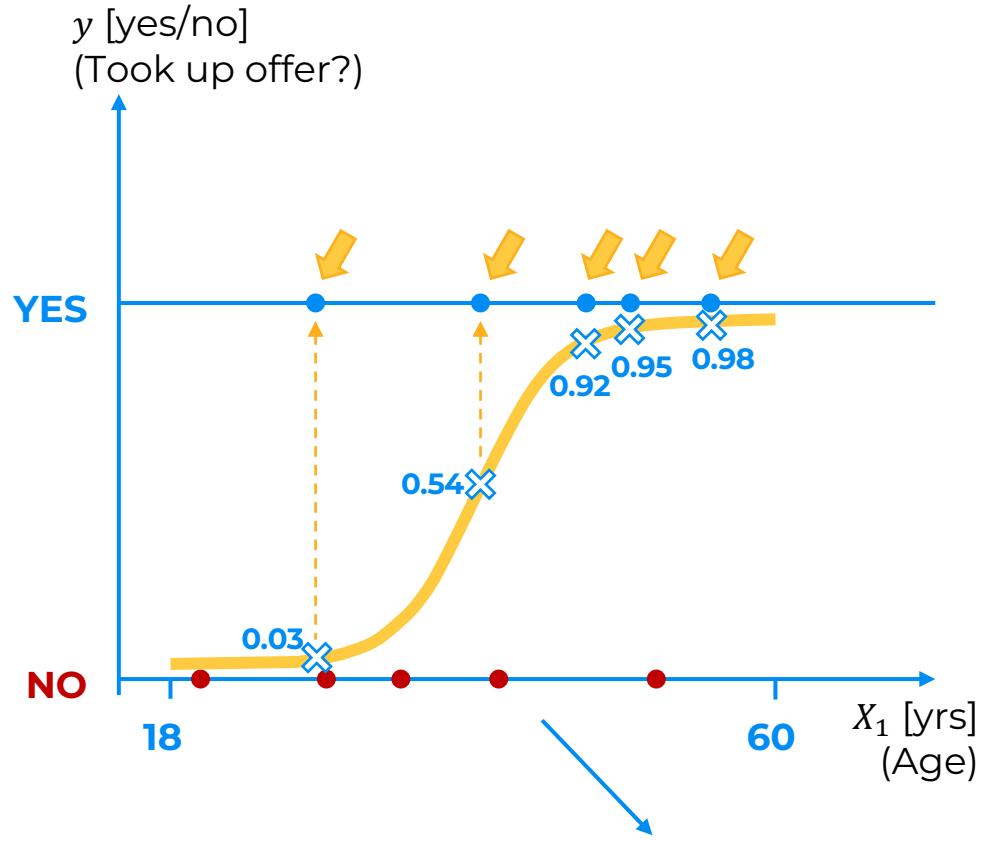
Maximum Likelihood

Maximum Likelihood



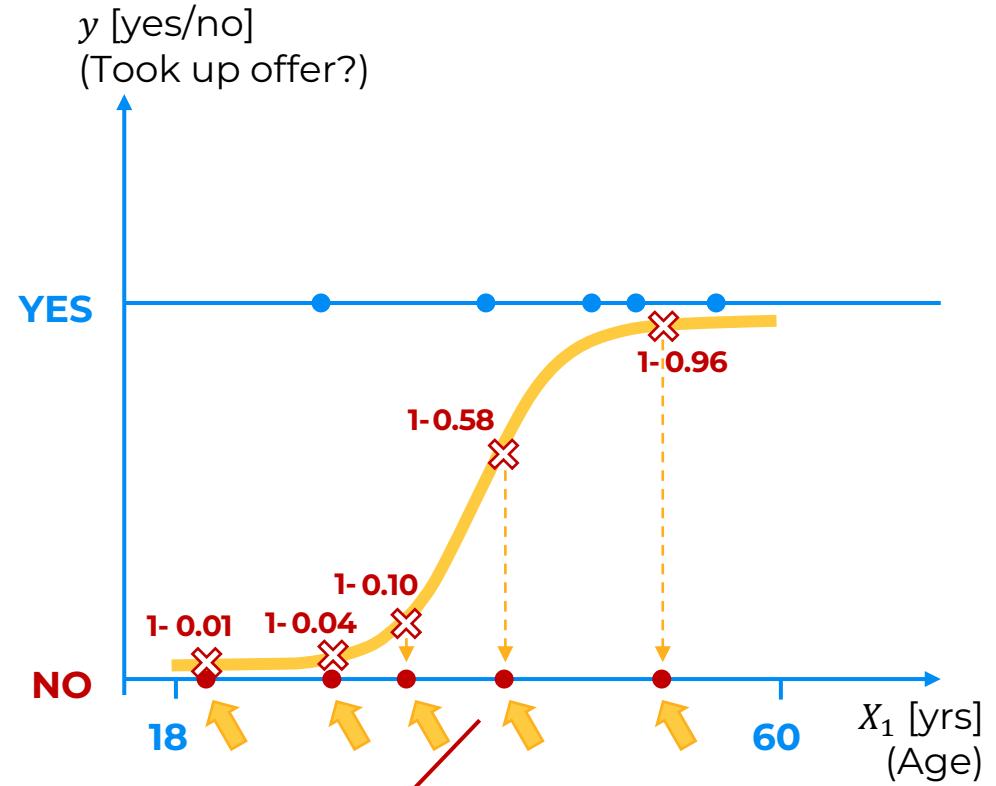
NOT FOR DISTRIBUTION

© SUPERDATASCIENCE

www.superdatascience.com

$$\text{Likelihood} = 0.03 \times 0.54 \times 0.92 \times 0.95 \times 0.98 \times (1 - 0.01) \times (1 - 0.04) \times (1 - 0.10) \times (1 - 0.58) \times (1 - 0.96)$$

$$\text{Likelihood} = \mathbf{0.00019939}$$



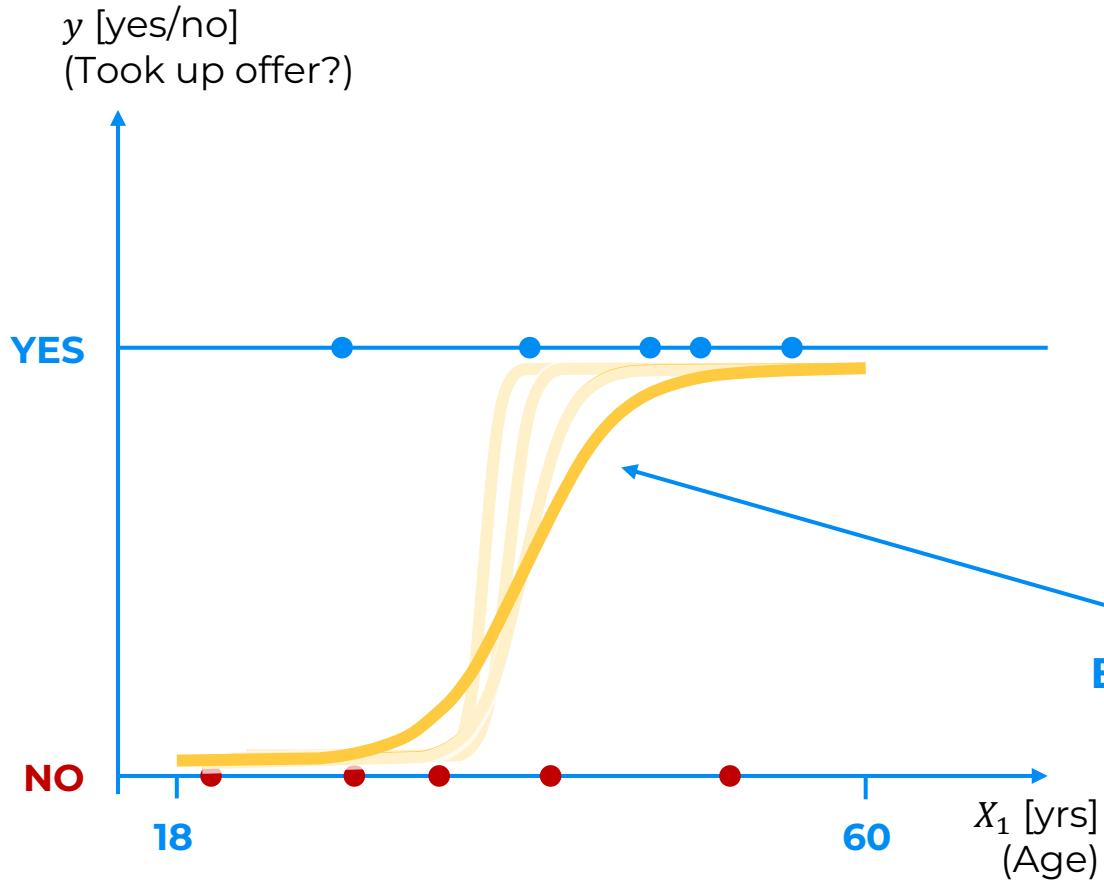
$$\text{Likelihood} = 0.00019939$$

Maximum Likelihood



NOT FOR DISTRIBUTION

© SUPERDATASCIENCE

www.superdatascience.com

Likelihood = 0.00007418

Likelihood = 0.00012845

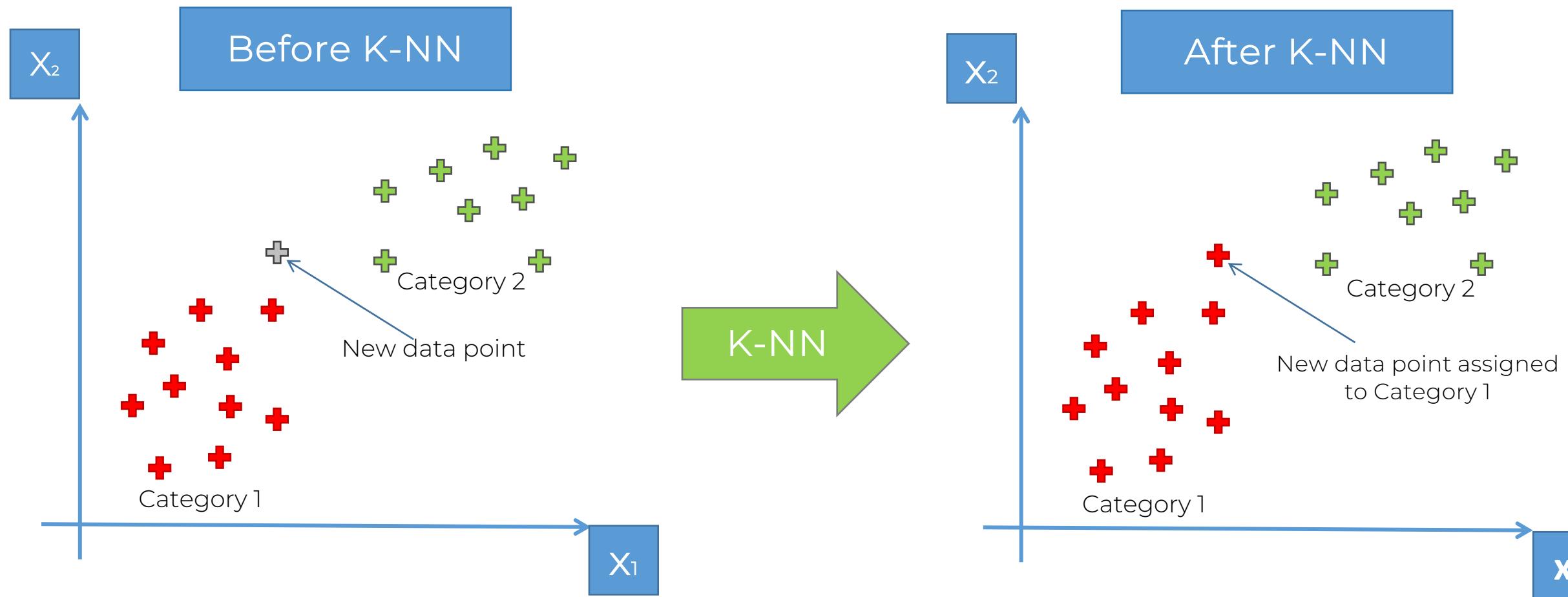
Likelihood = 0.00016553

Likelihood = 0.00019939

Best Curve <= Maximum Likelihood

K-NN Intuition

What K-NN does for you



How did it do that ?

STEP 1: Choose the number K of neighbors



STEP 2: Take the K nearest neighbors of the new data point, according to the Euclidean distance



STEP 3: Among these K neighbors, count the number of data points in each category



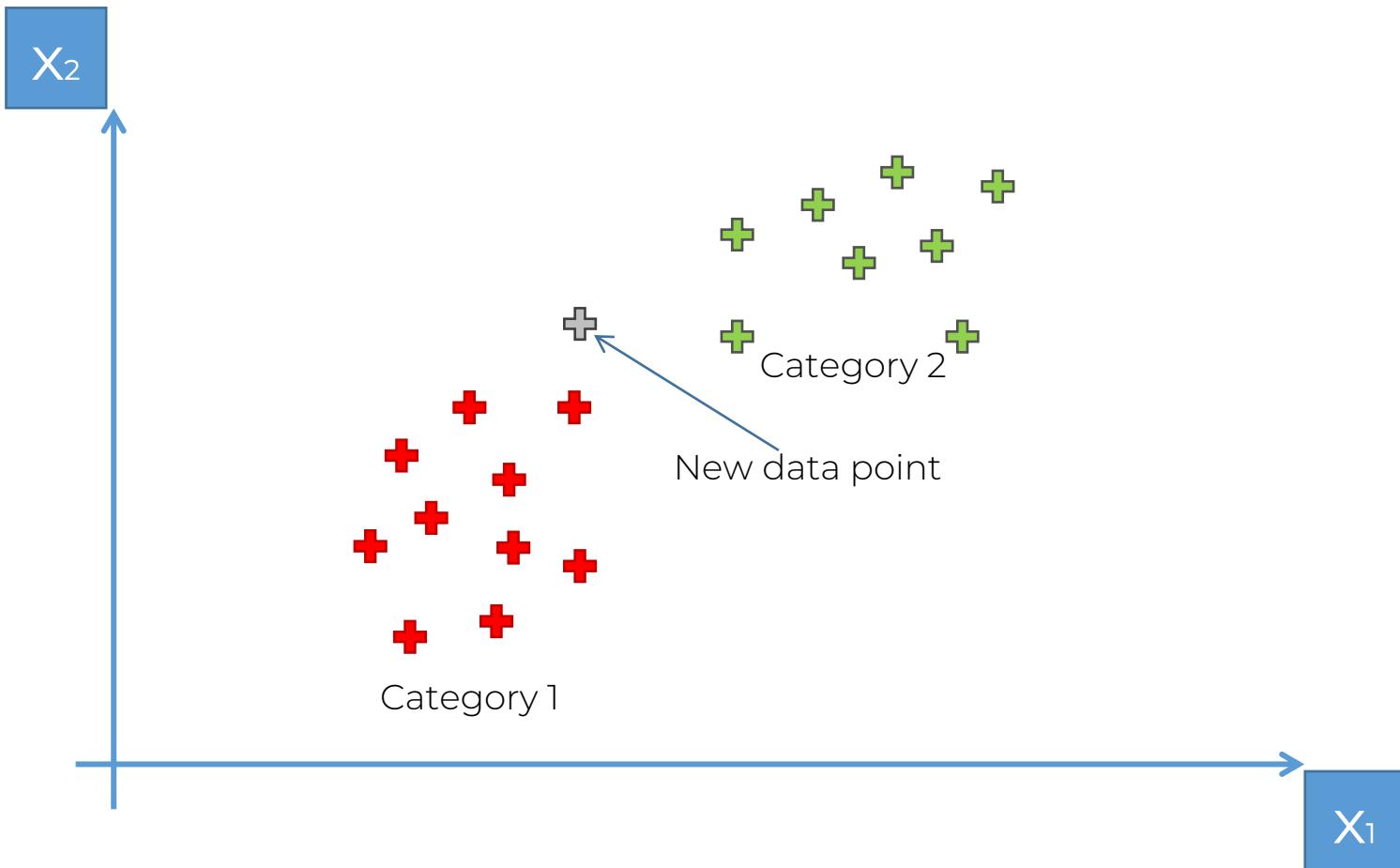
STEP 4: Assign the new data point to the category where you counted the most neighbors



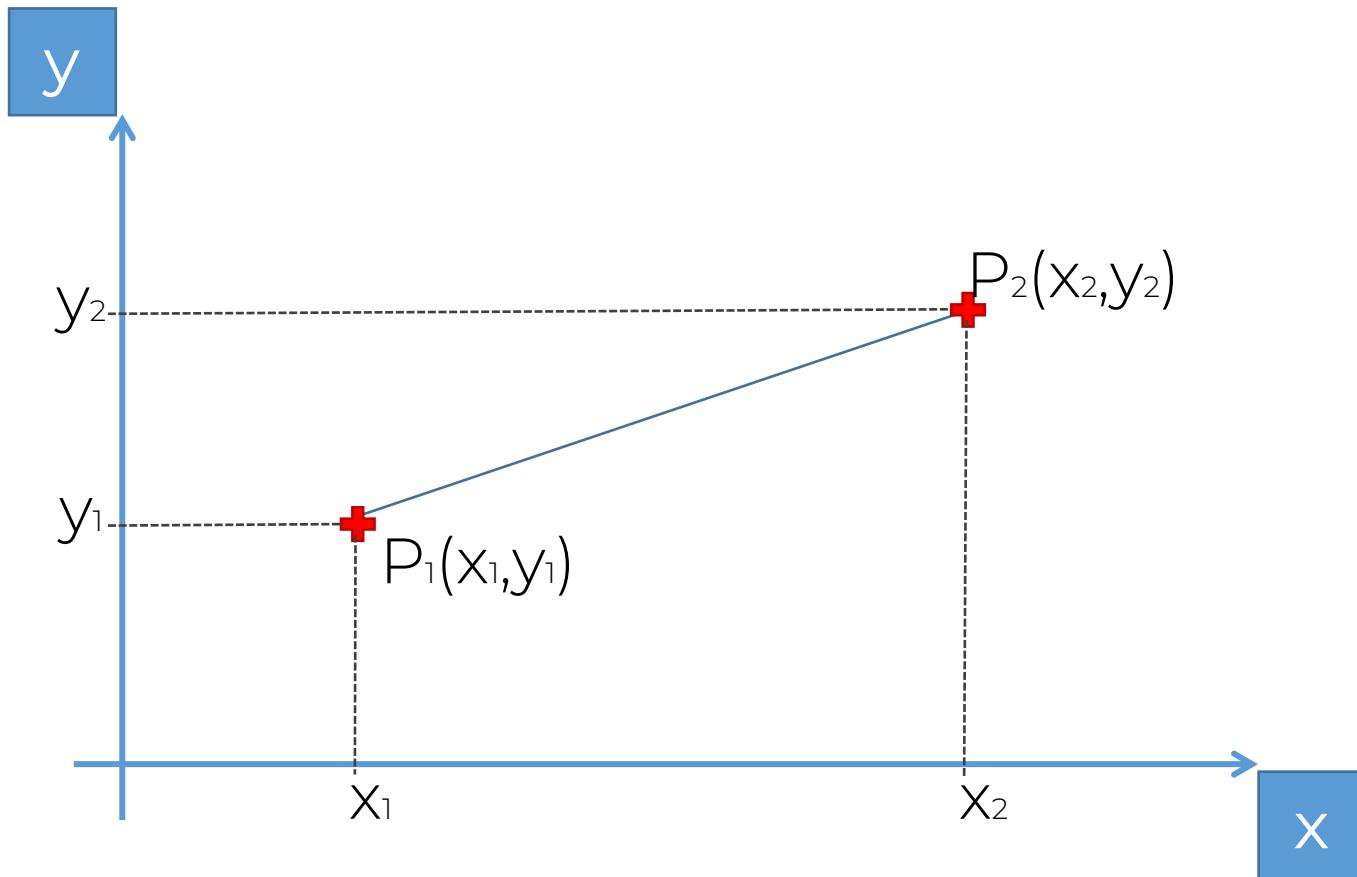
Your Model is Ready

K-NN algorithm

STEP 1: Choose the number K of neighbors: K = 5

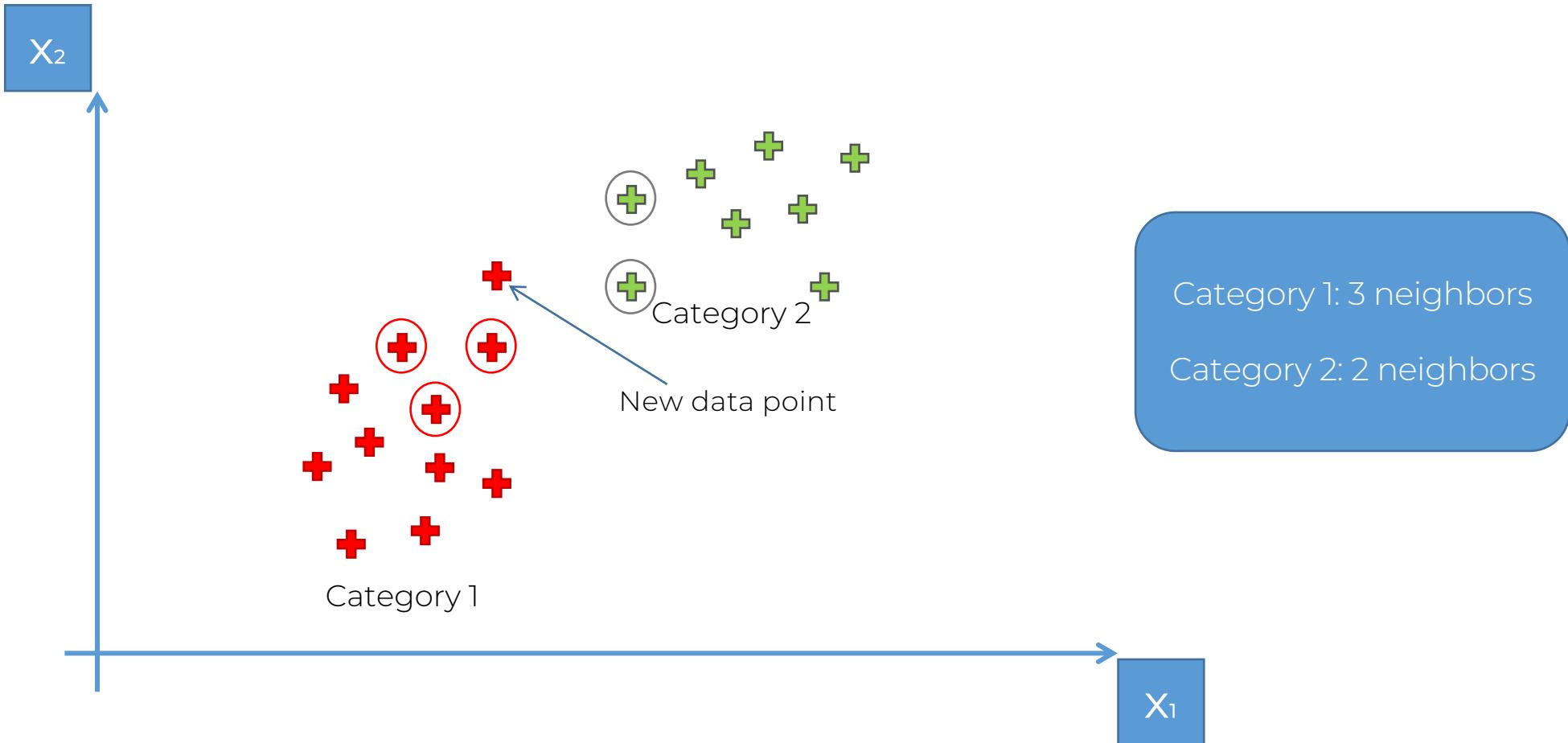


Euclidean Distance



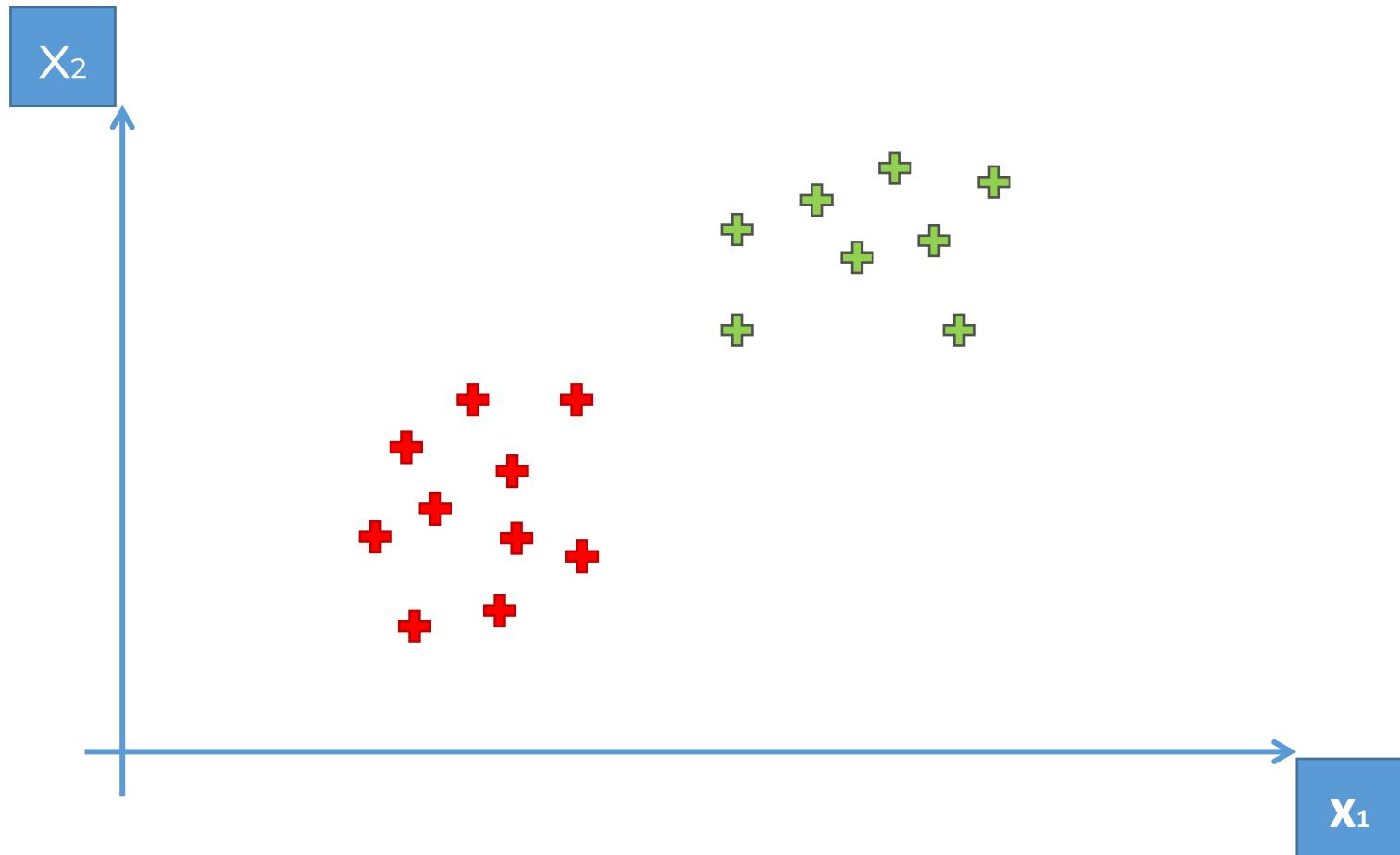
$$\text{Euclidean Distance between } P_1 \text{ and } P_2 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

K-NN algorithm

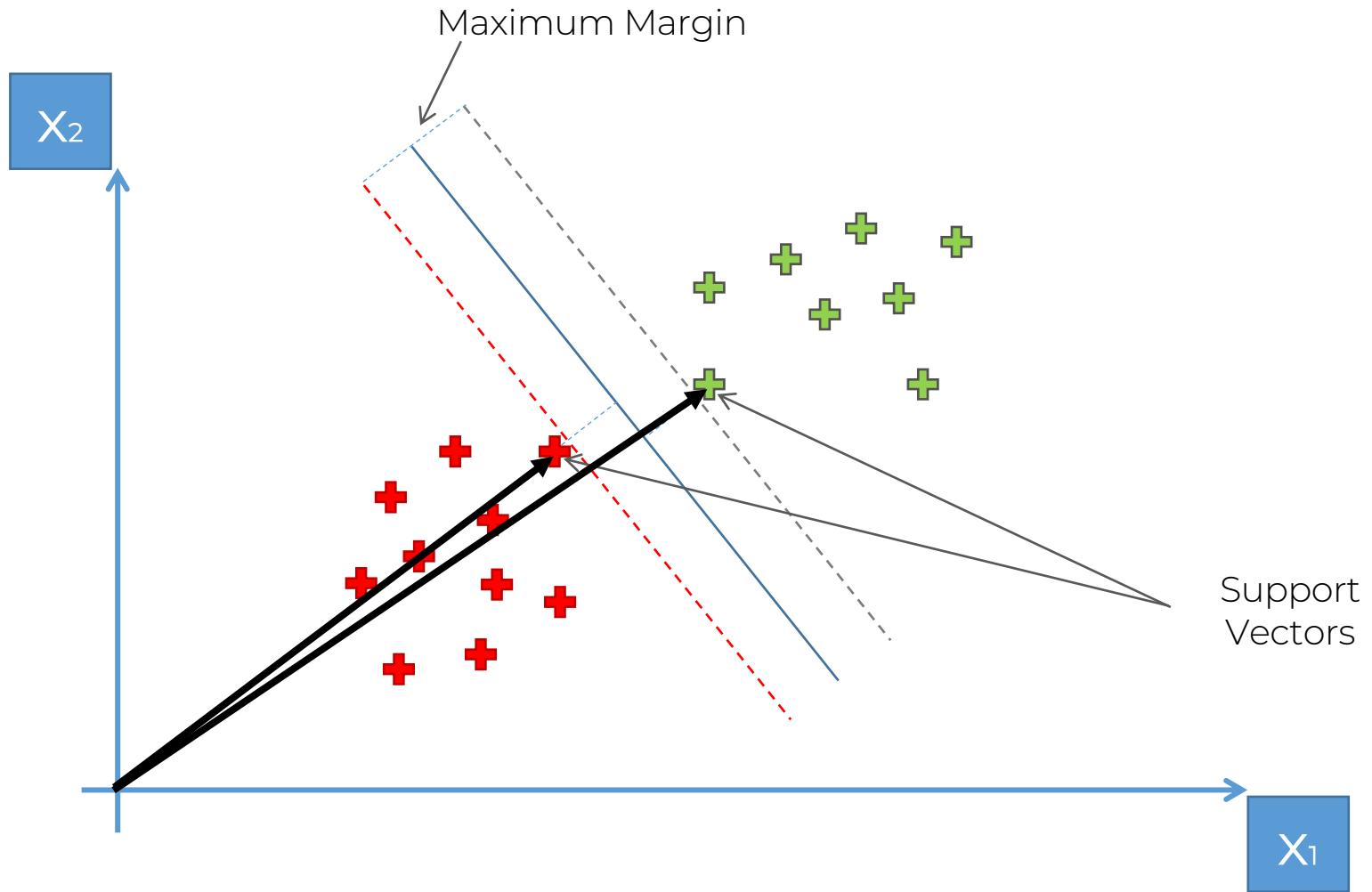


SVM Intuition

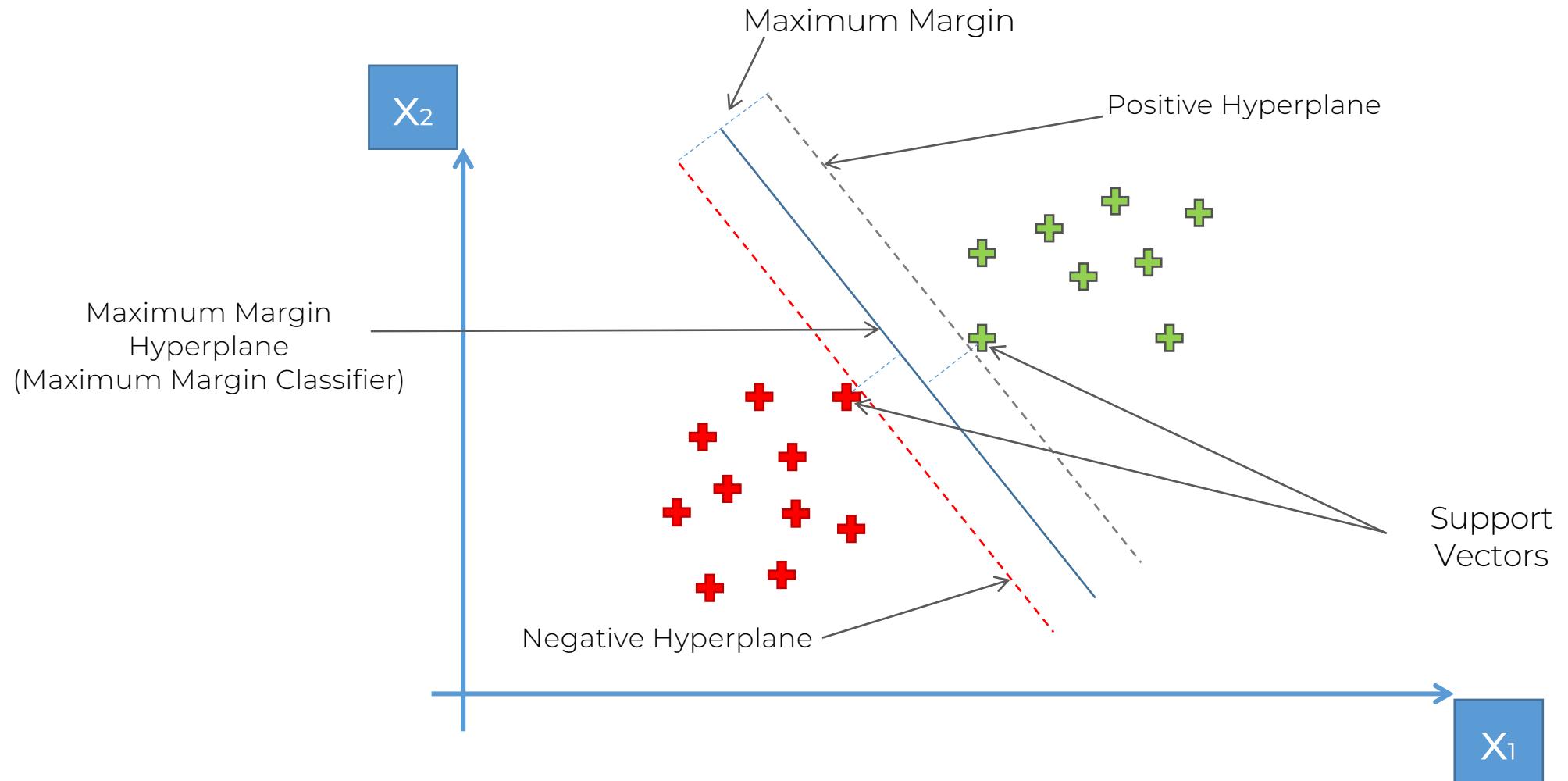
How to separate these points ?



Support Vectors

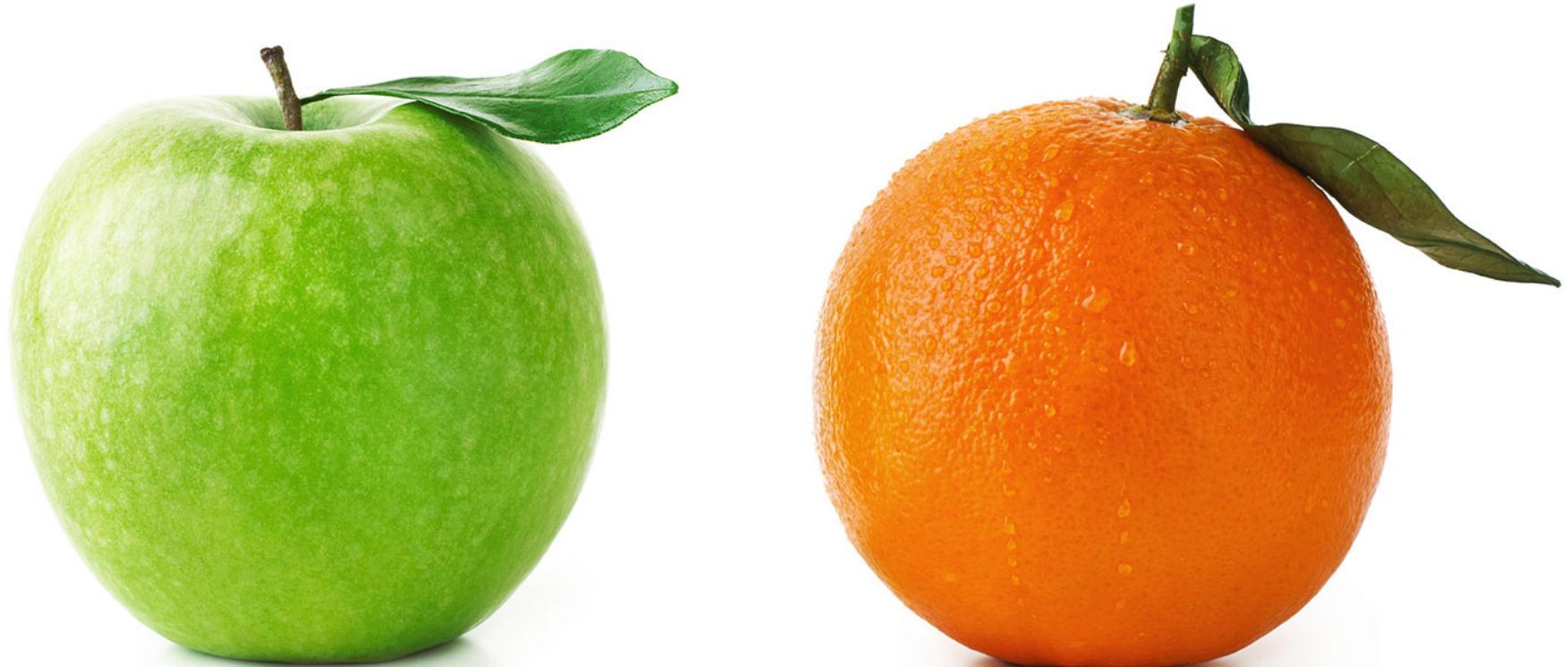


Hyperplanes

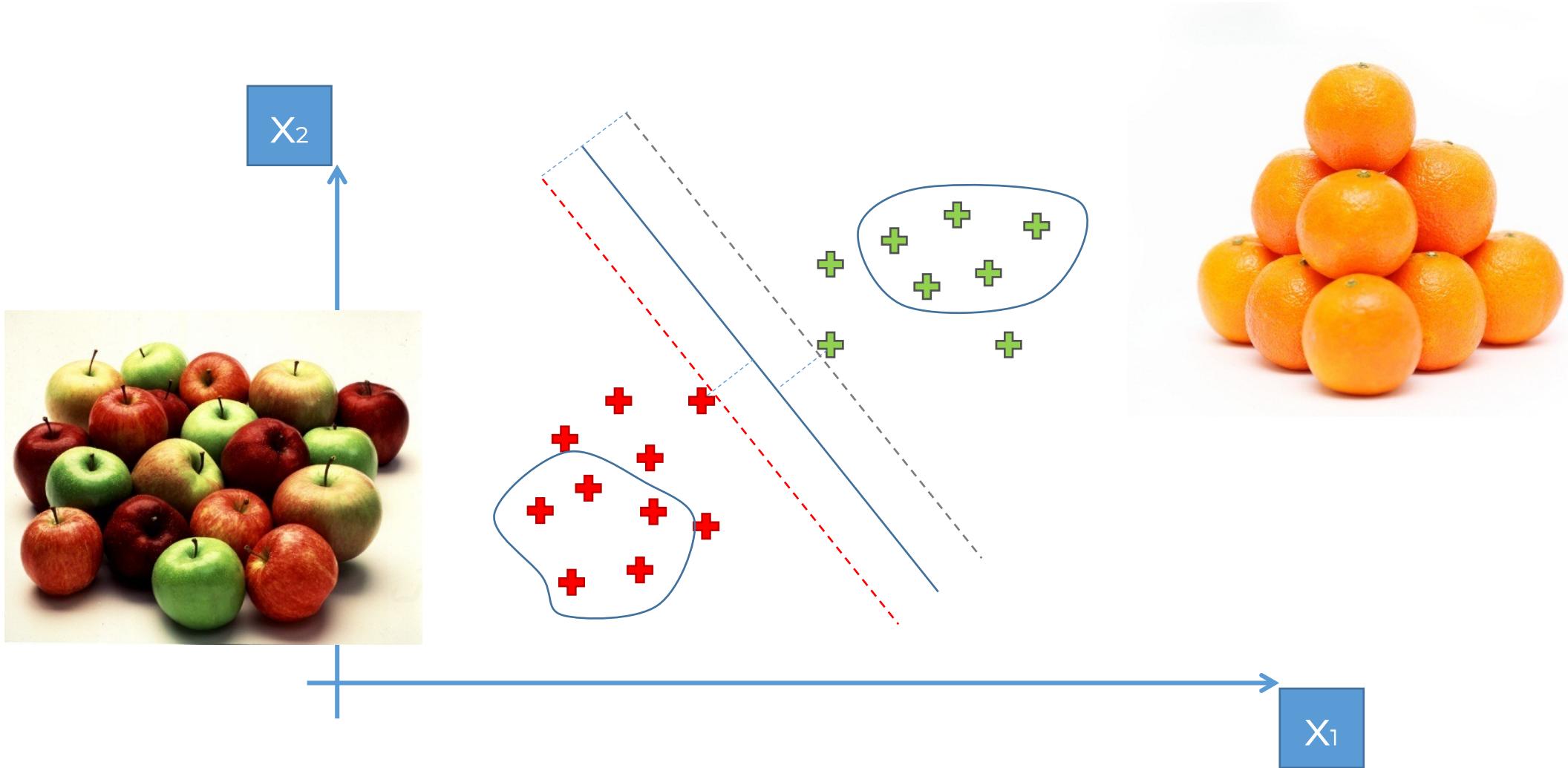


What's So Special About SVMs?

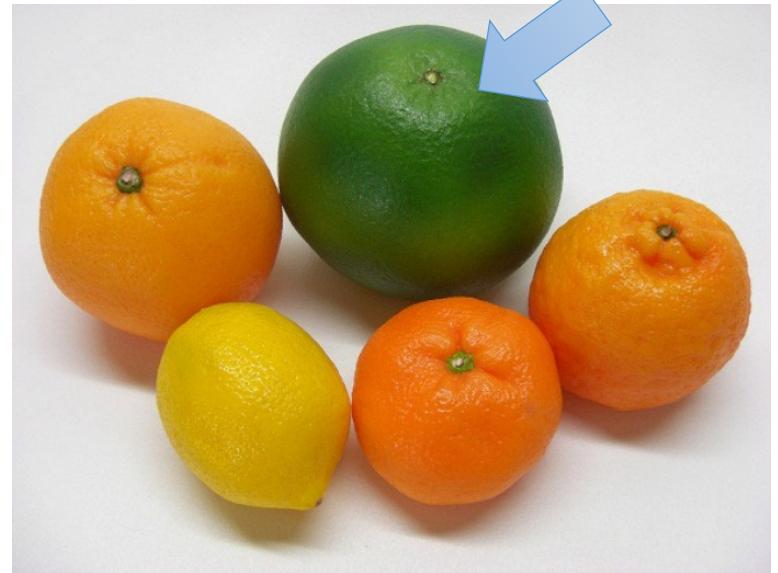
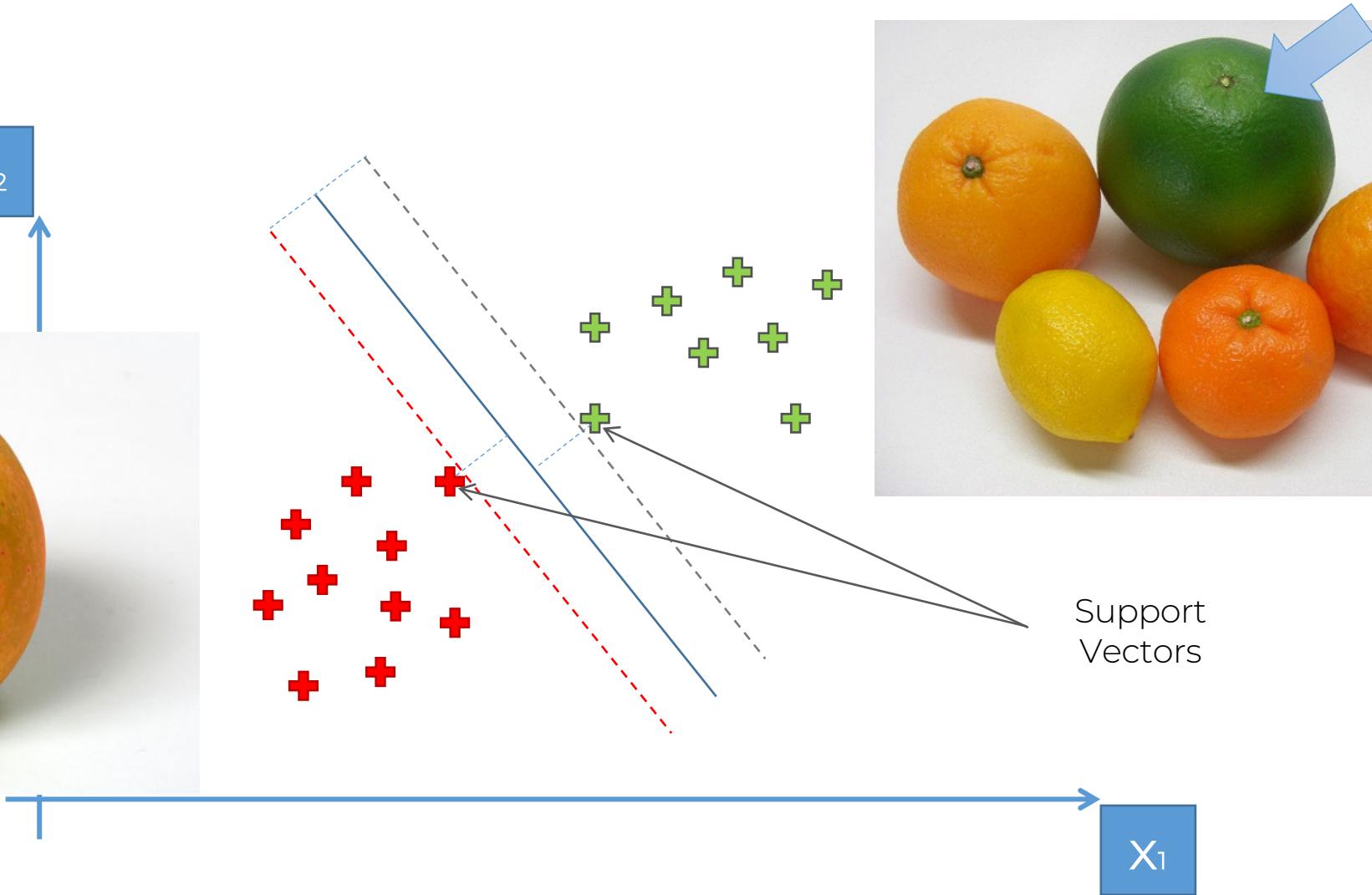
What's So Special About SVMs?



What's So Special About SVMs?



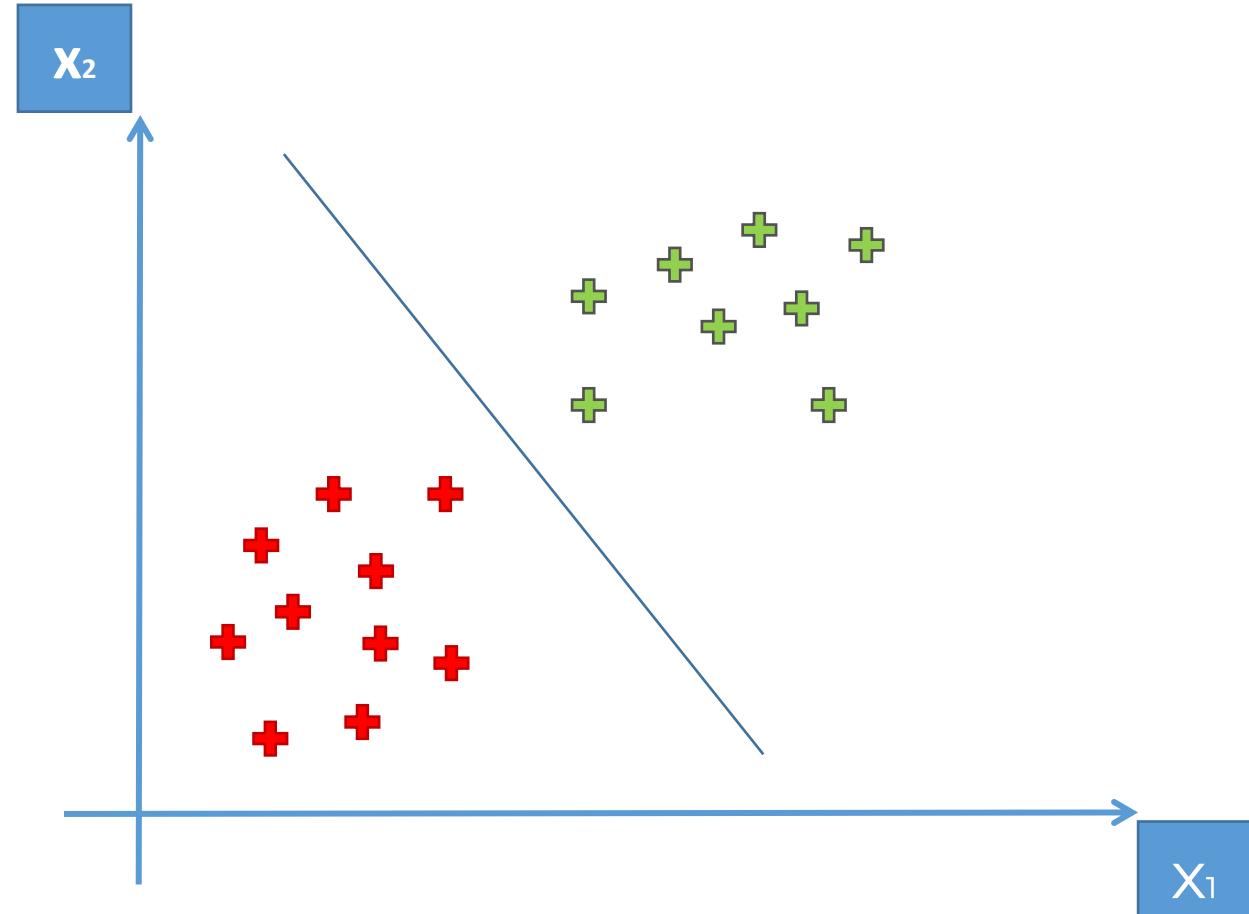
What's So Special About SVMs?



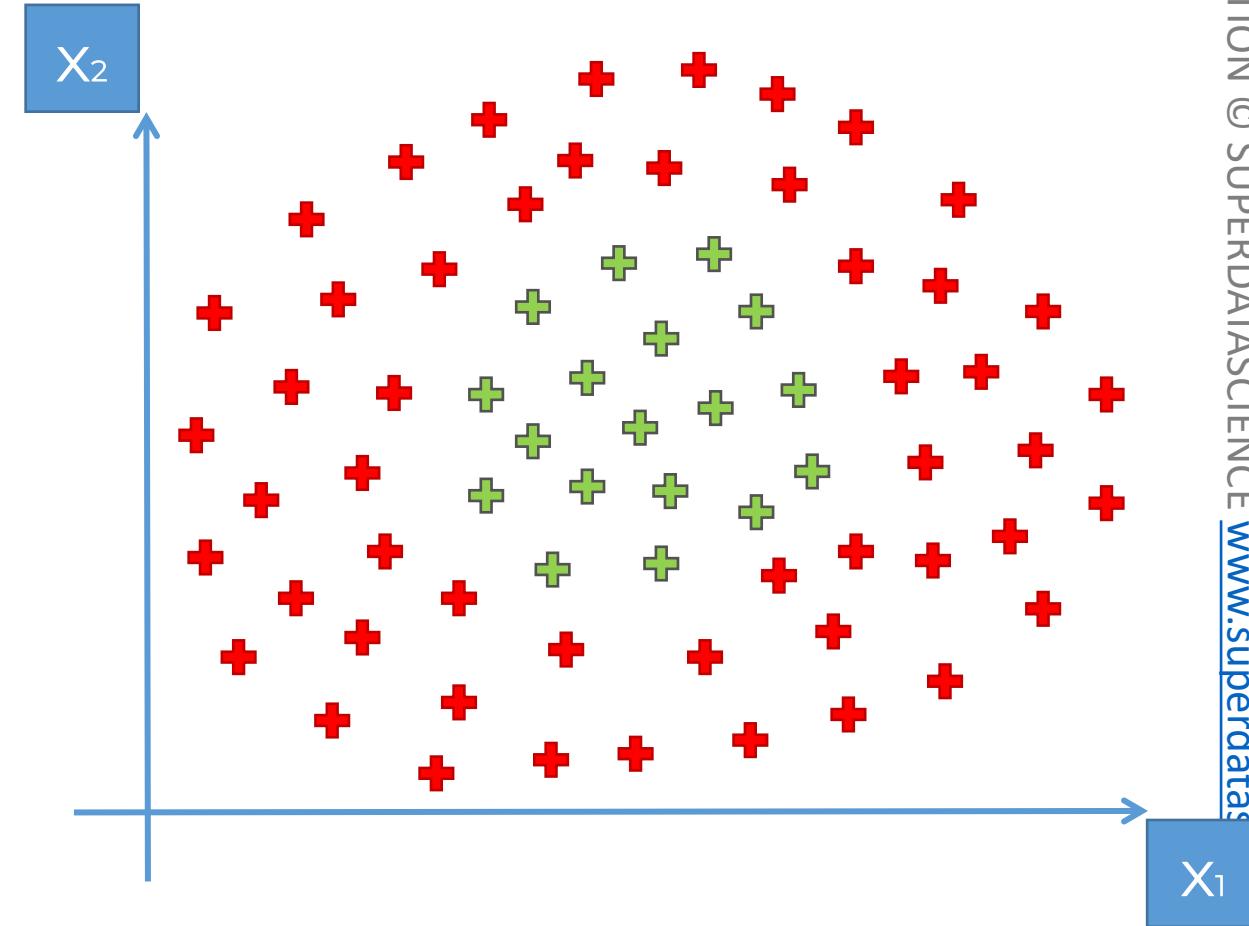
Kernel SVM Intuition

Linear Separability

Linearly Separable



Not Linearly Separable



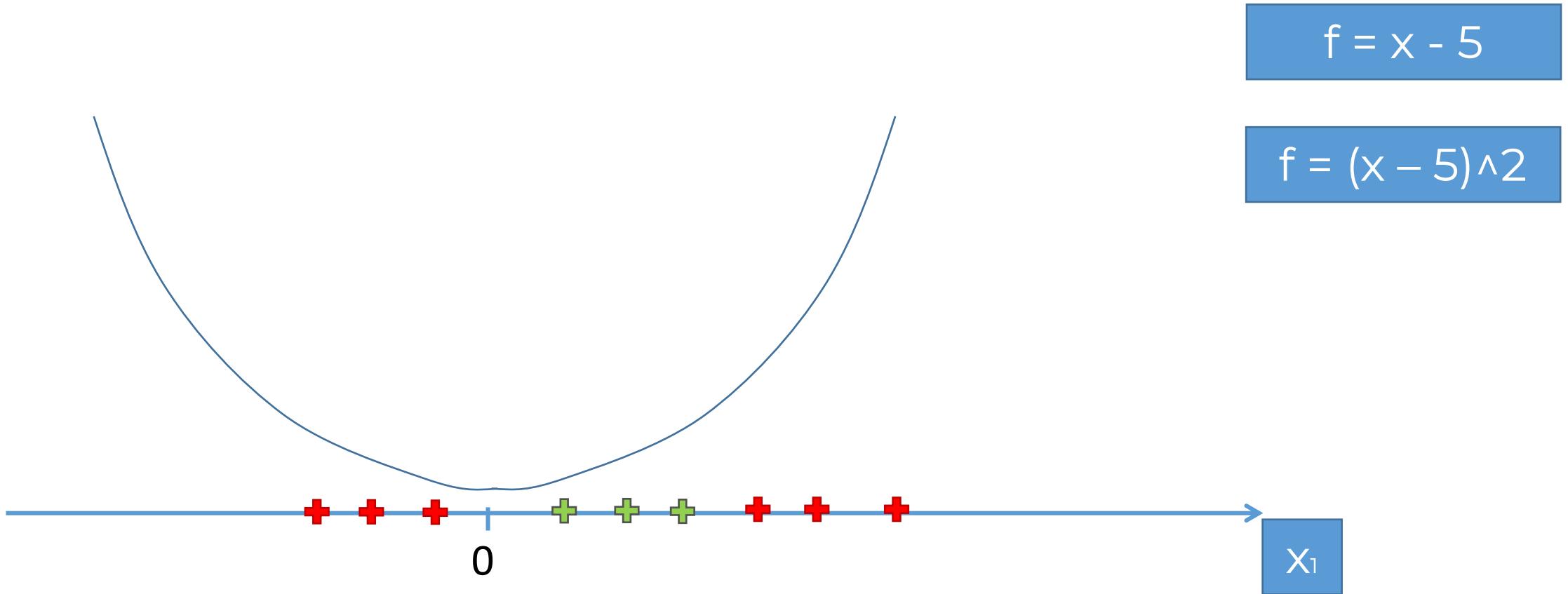
A Higher-Dimensional Space

Mapping to a Higher Dimension

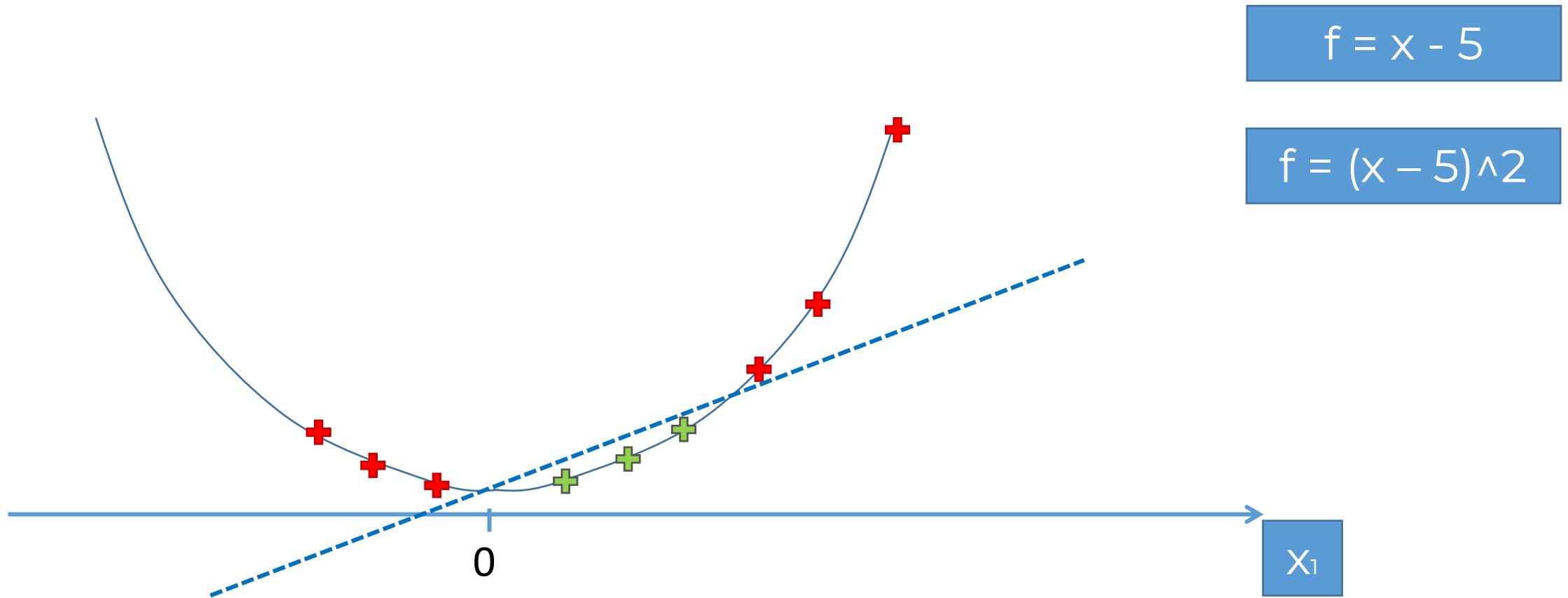
$$f = x - 5$$



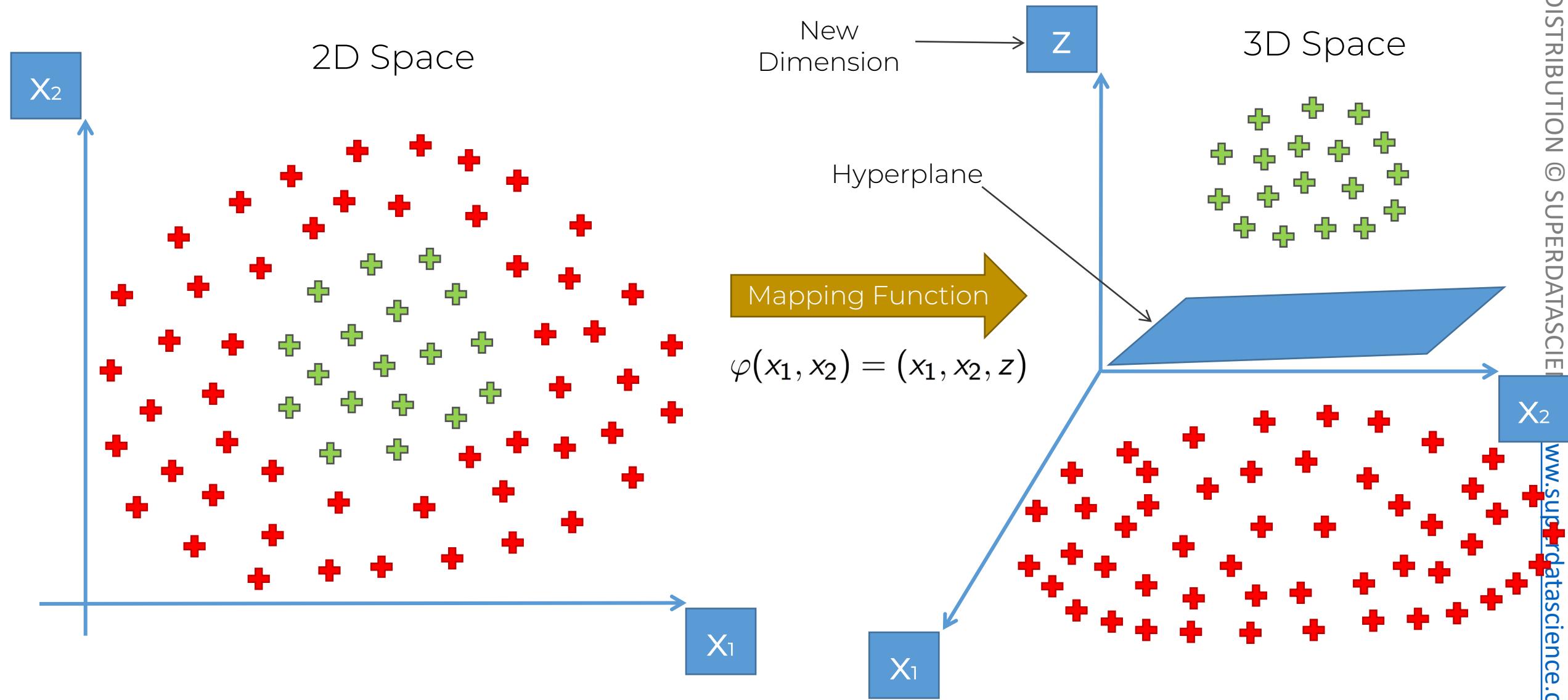
Mapping to a Higher Dimension



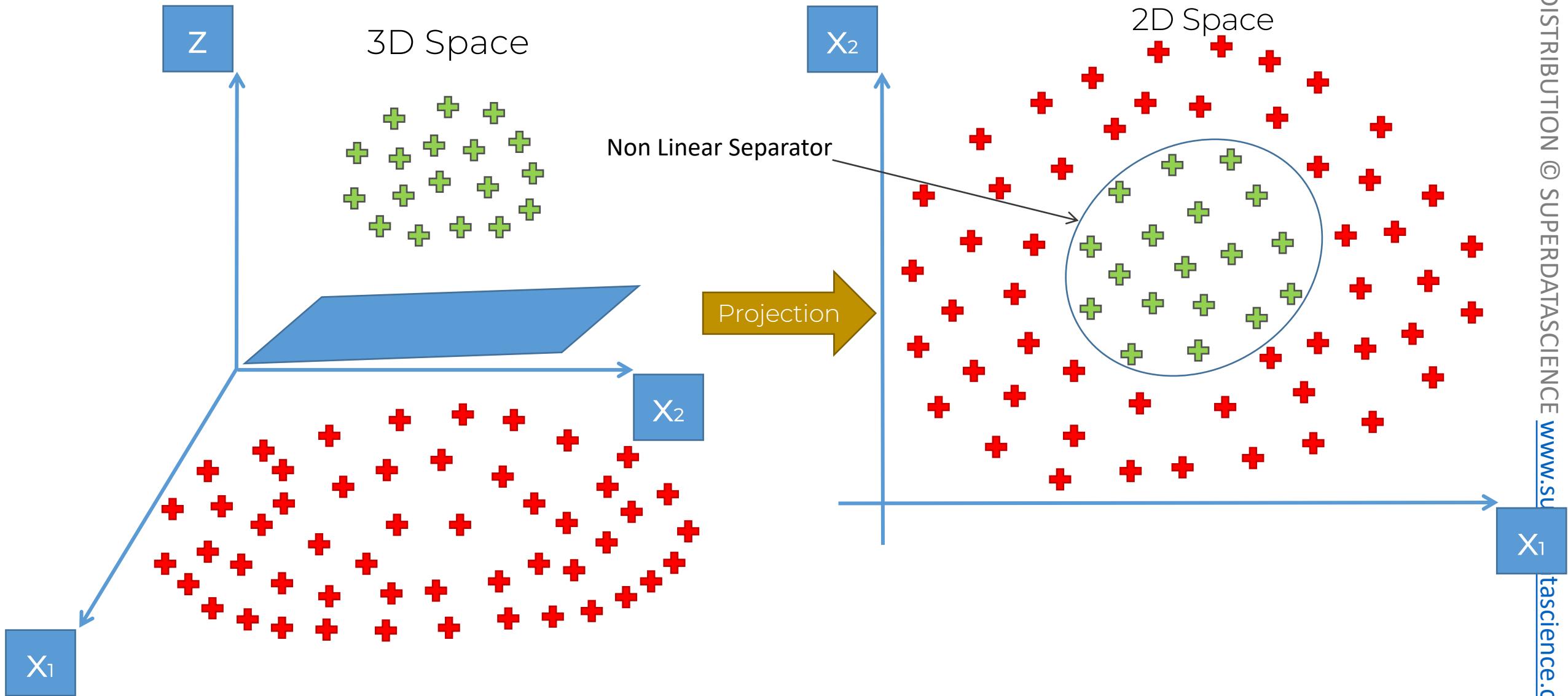
Mapping to a Higher Dimension



Mapping to a Higher Dimension



Projecting back to 2D Space



But there is a catch...

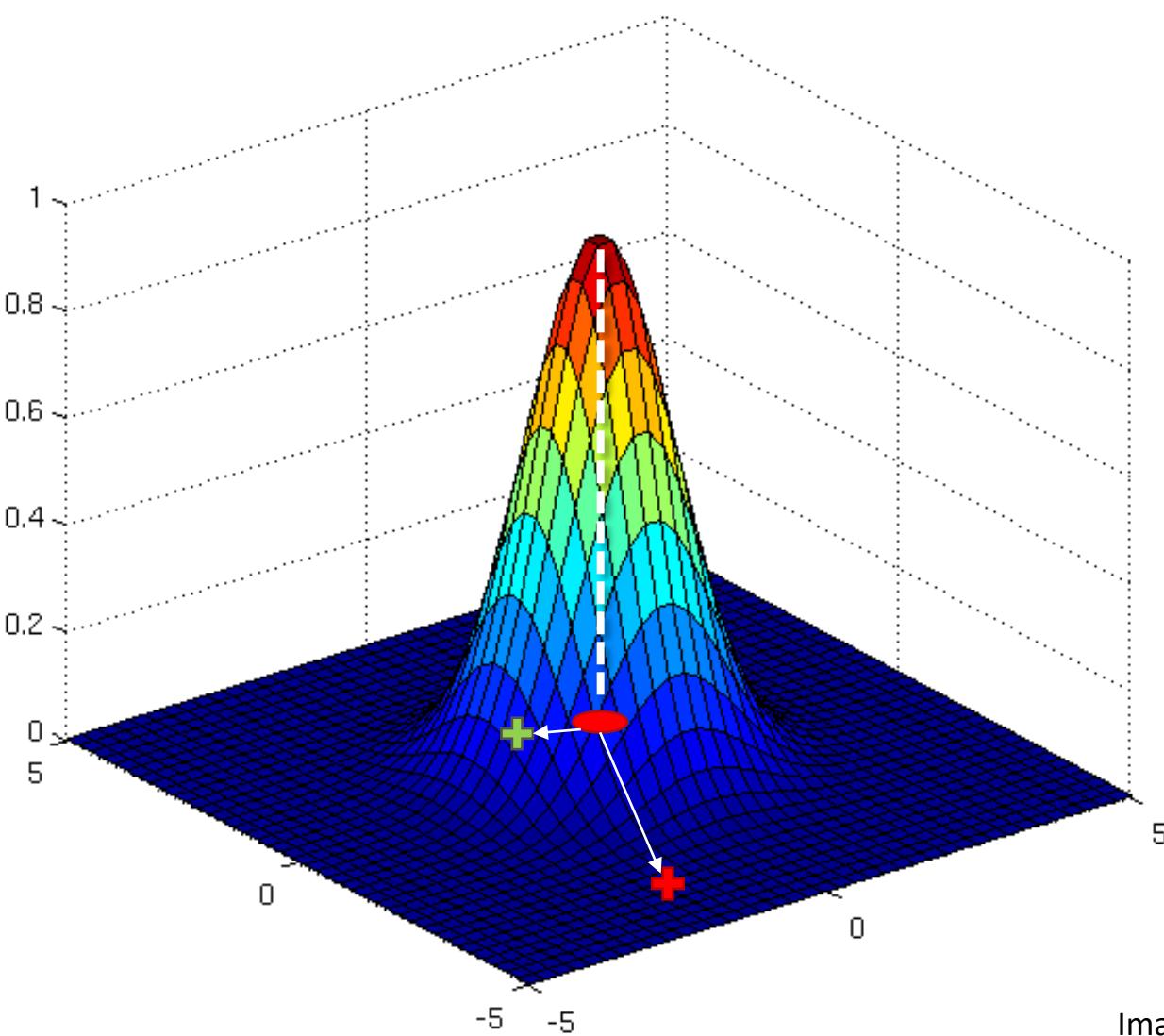
Mapping to a Higher Dimensional Space
can be highly compute-intensive

The Kernel Trick

The Gaussian RBF Kernel

$$K(\vec{x}, \vec{l}^i) = e^{-\frac{\|\vec{x} - \vec{l}^i\|^2}{2\sigma^2}}$$

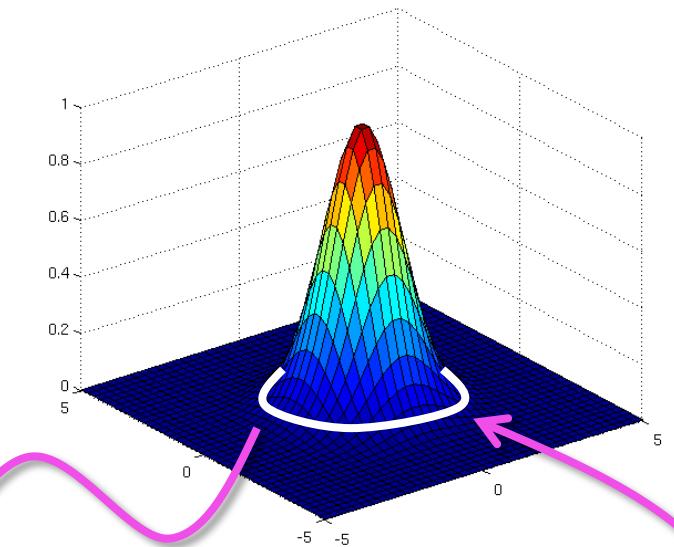
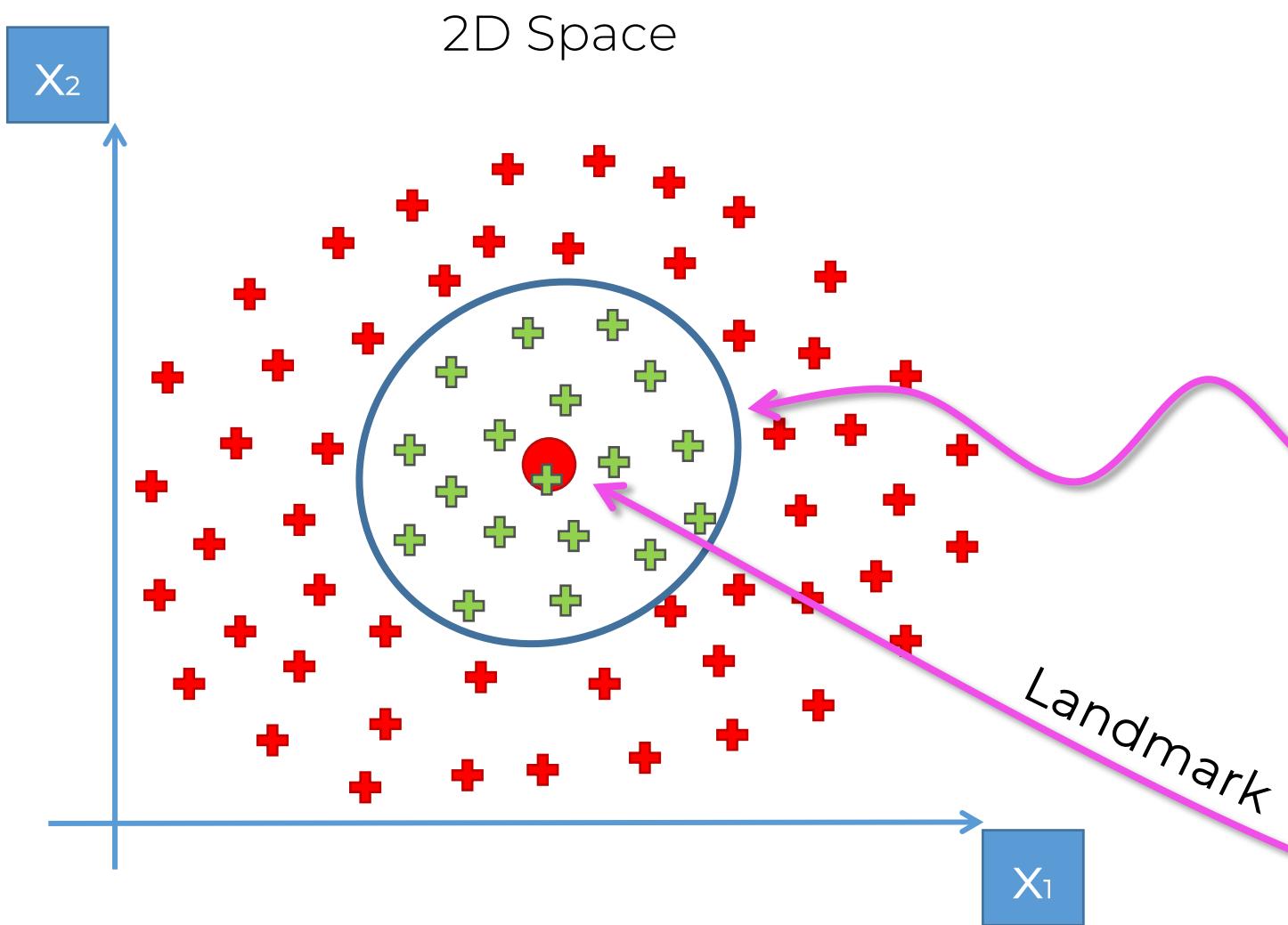
The Gaussian RBF Kernel



$$K(\vec{x}, \vec{l}^i) = e^{-\frac{\|\vec{x} - \vec{l}^i\|^2}{2\sigma^2}}$$

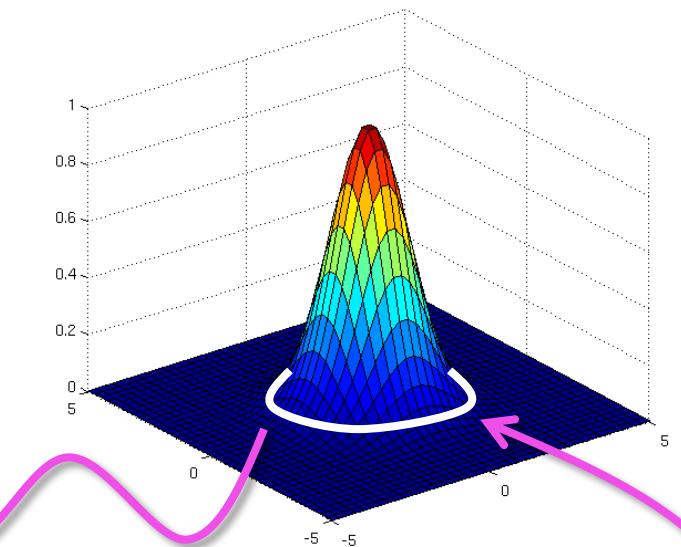
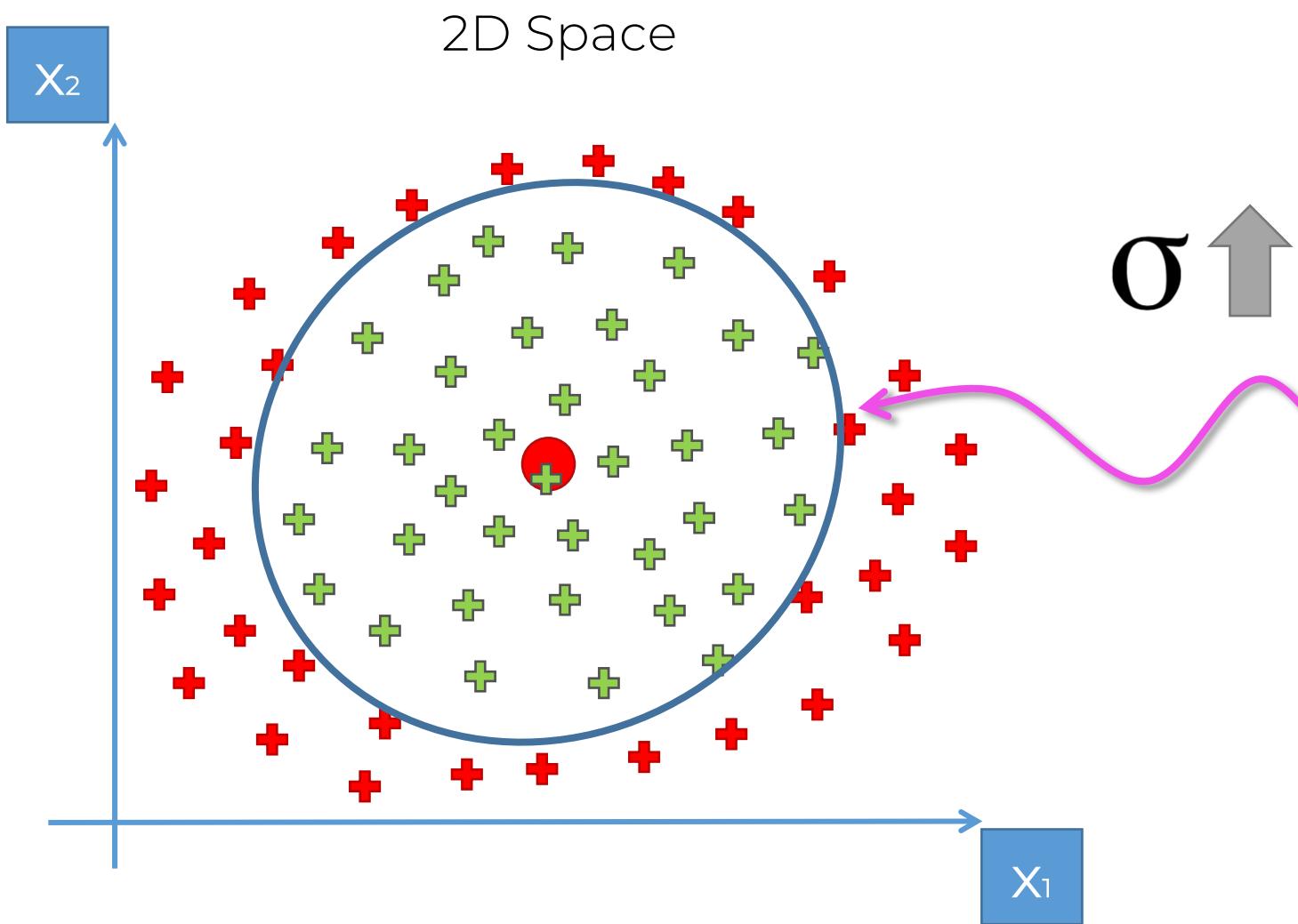
Image source: <http://www.cs.toronto.edu/~duvenaud/cookbook/index.html>

The Gaussian RBF Kernel



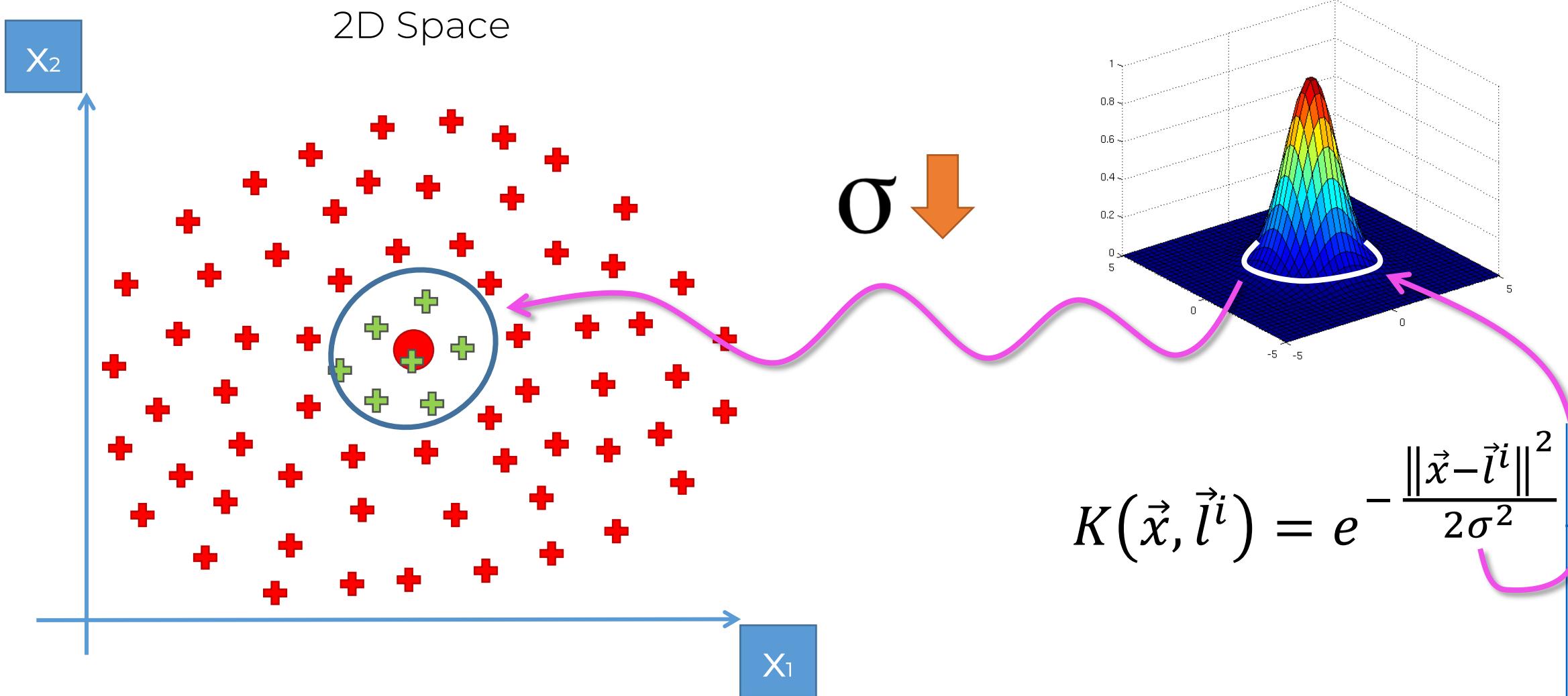
$$K(\vec{x}, \vec{l}^i) = e^{-\frac{\|\vec{x} - \vec{l}^i\|^2}{2\sigma^2}}$$

The Gaussian RBF Kernel

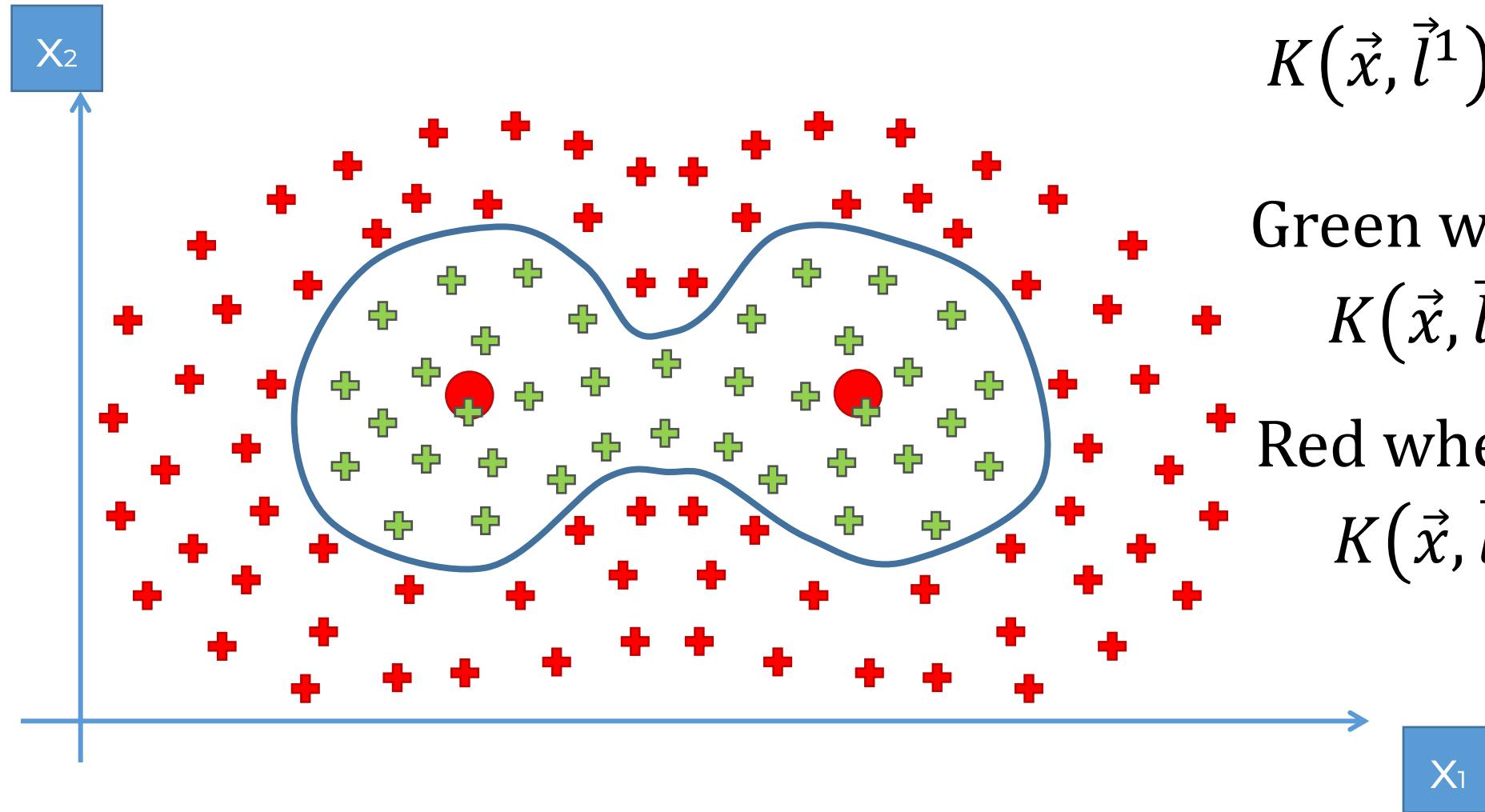


$$K(\vec{x}, \vec{l}^i) = e^{-\frac{\|\vec{x} - \vec{l}^i\|^2}{2\sigma^2}}$$

The Gaussian RBF Kernel

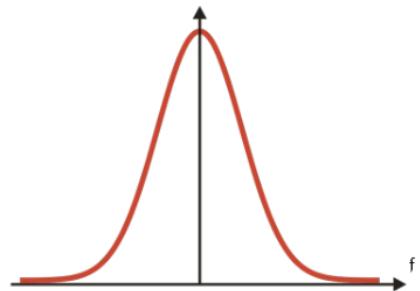


The Gaussian RBF Kernel



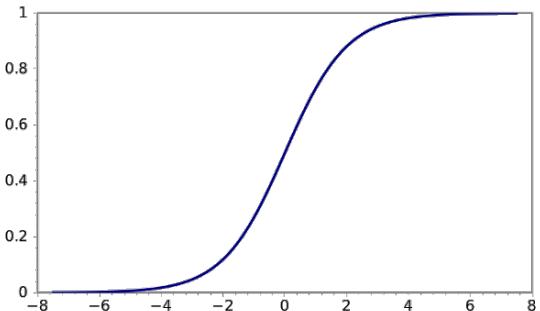
Types of Kernel Functions

Types of Kernel Functions



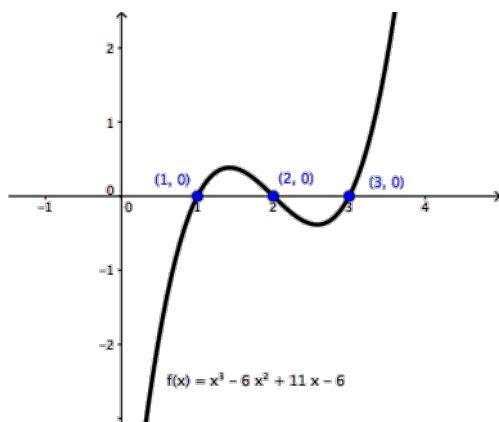
Gaussian RBF Kernel

$$K(\vec{x}, \vec{l}^i) = e^{-\frac{\|\vec{x} - \vec{l}^i\|^2}{2\sigma^2}}$$



Sigmoid Kernel

$$K(X, Y) = \tanh(\gamma \cdot X^T Y + r)$$



Polynomial Kernel

$$K(X, Y) = (\gamma \cdot X^T Y + r)^d, \gamma >$$

Non-Linear SVR (Advanced)

Heads-up about Non-Linear SVR

Section on SVR:

- SVR Intuition



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

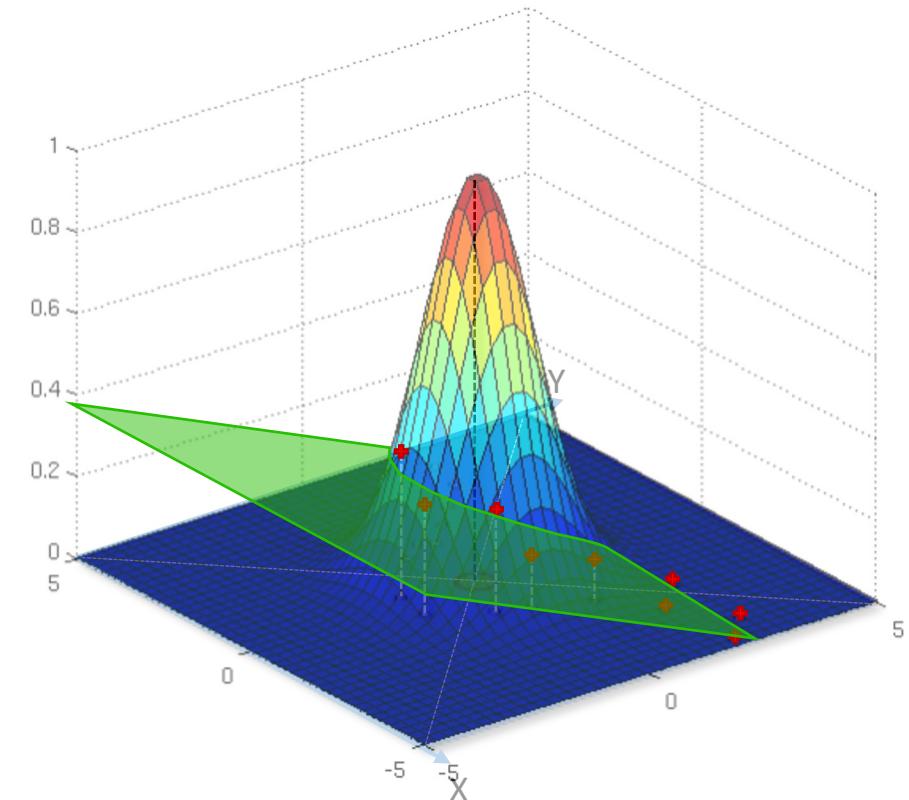
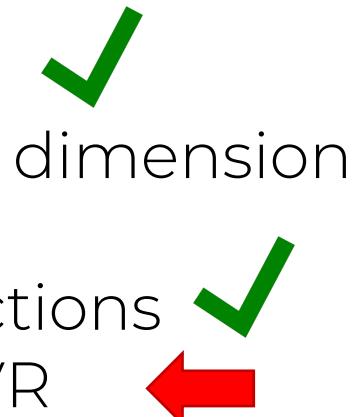
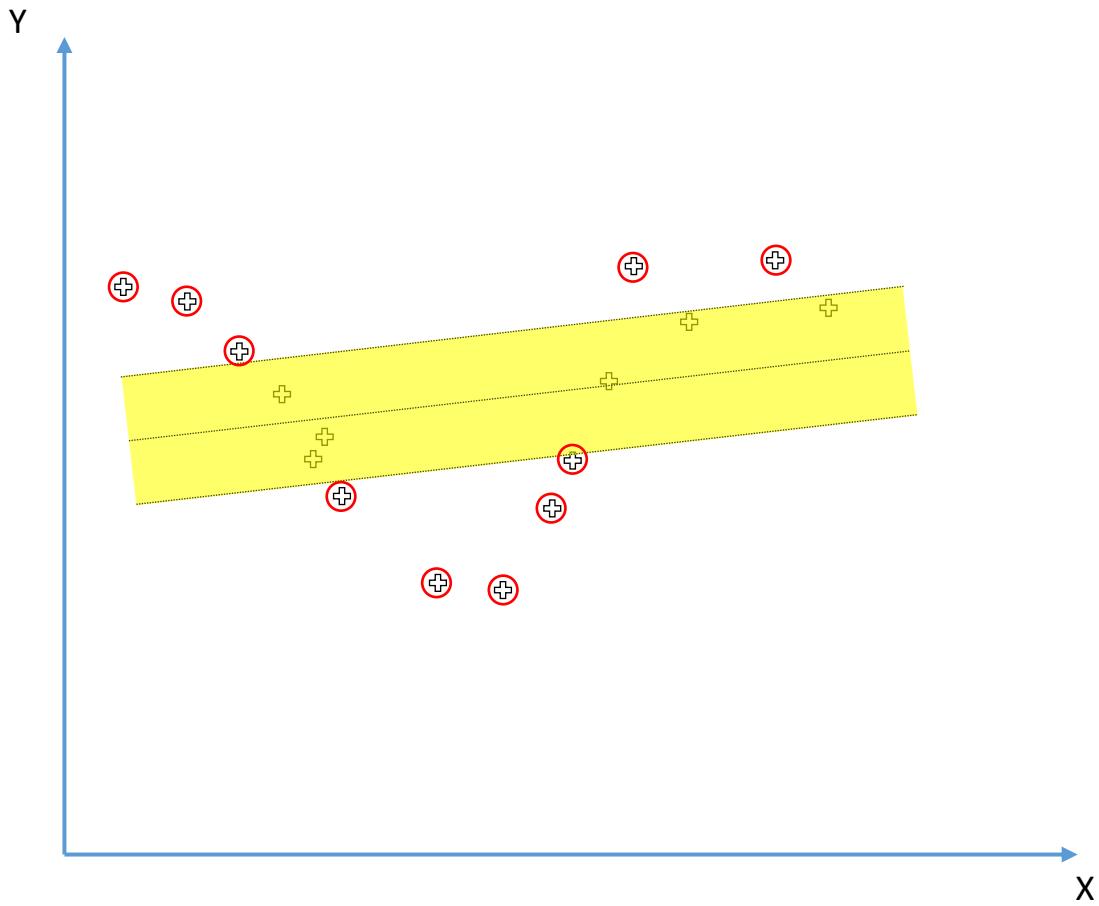


Image source: <http://www.cs.toronto.edu/~duvenaud/cookbook/index.html>

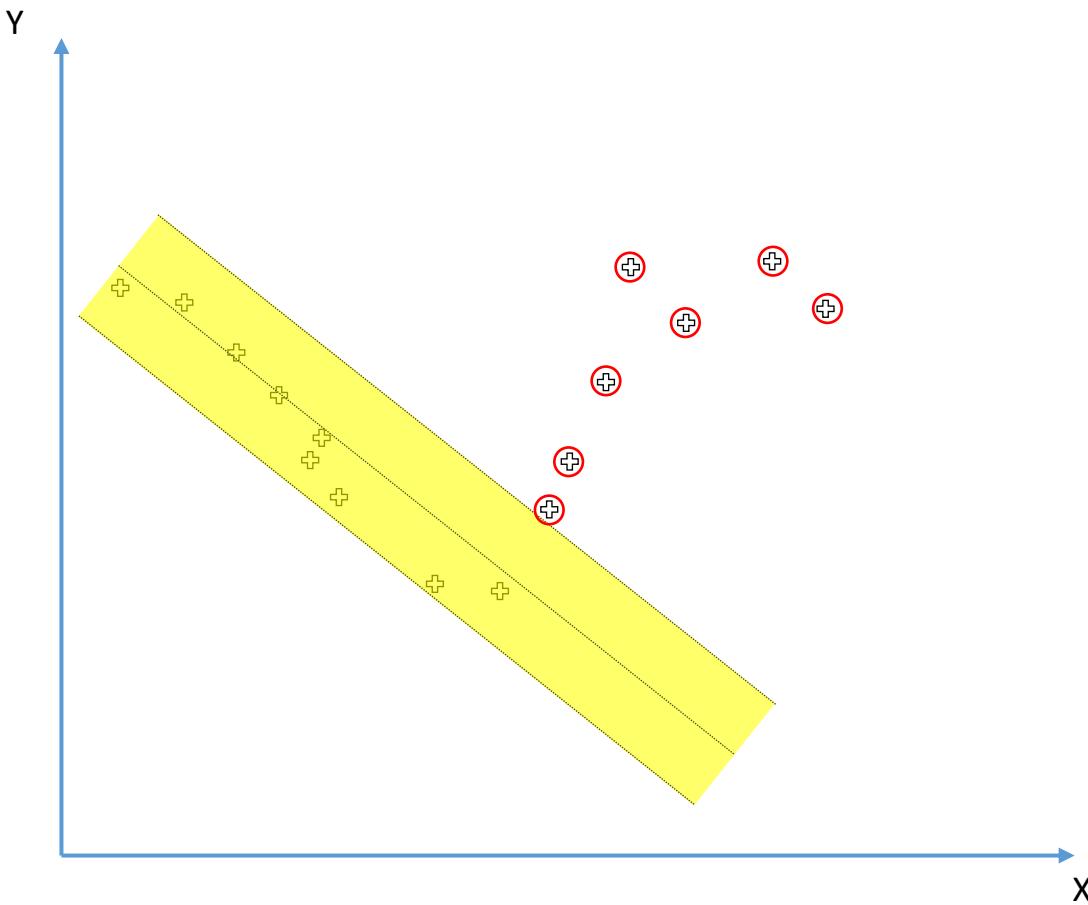
Non-Linear SVR



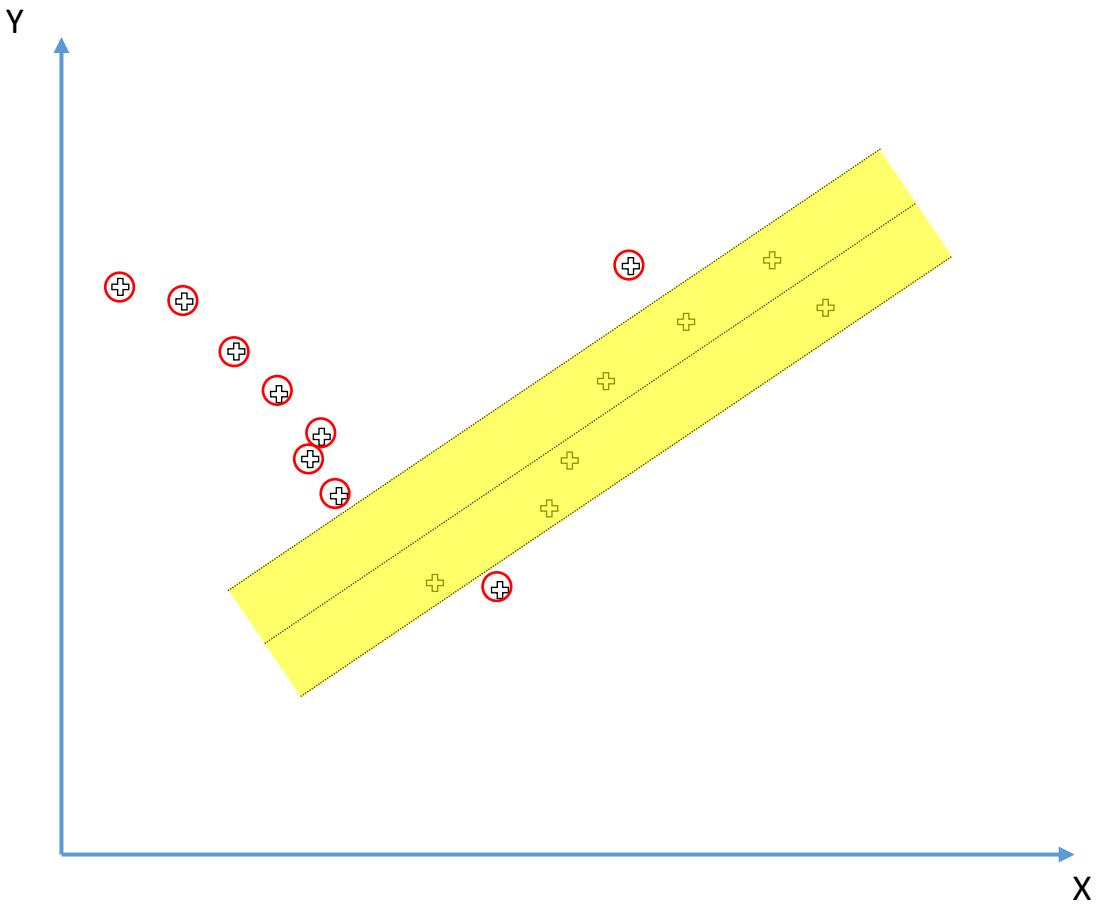
Non-Linear SVR



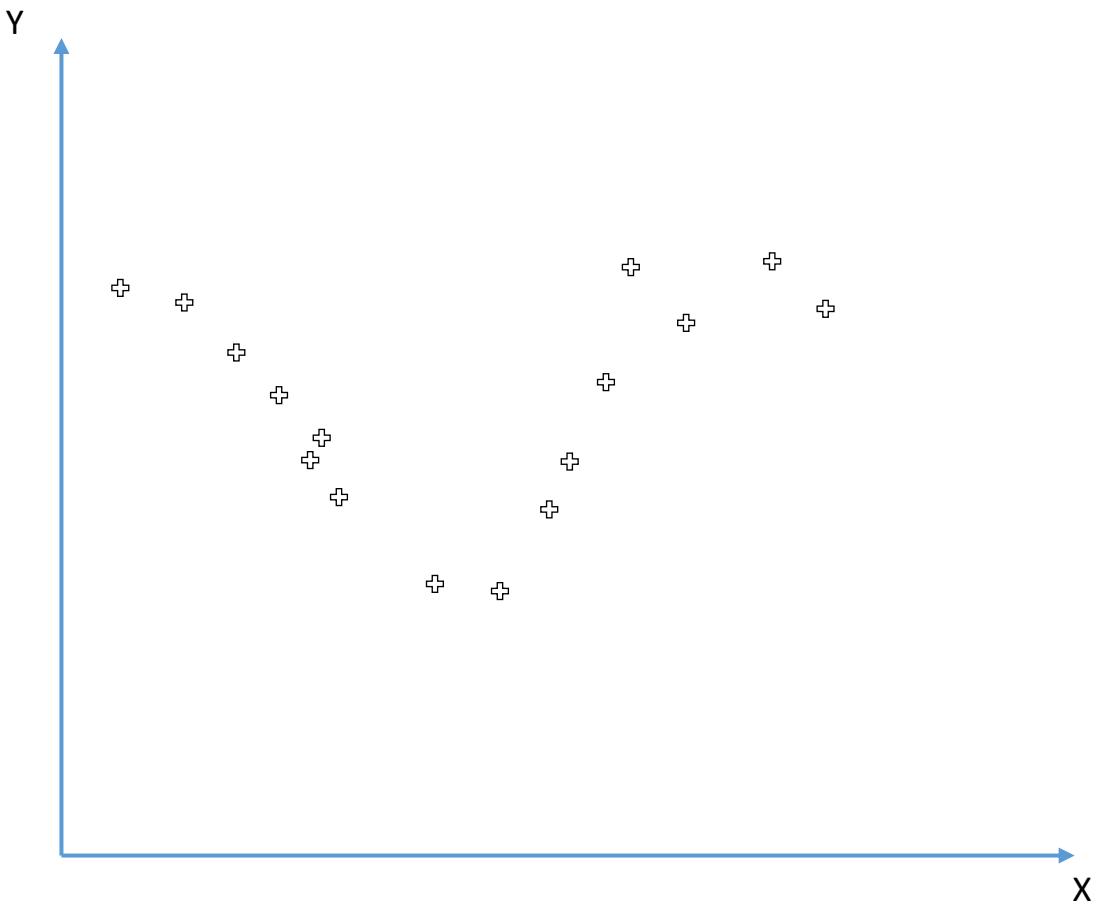
Non-Linear SVR



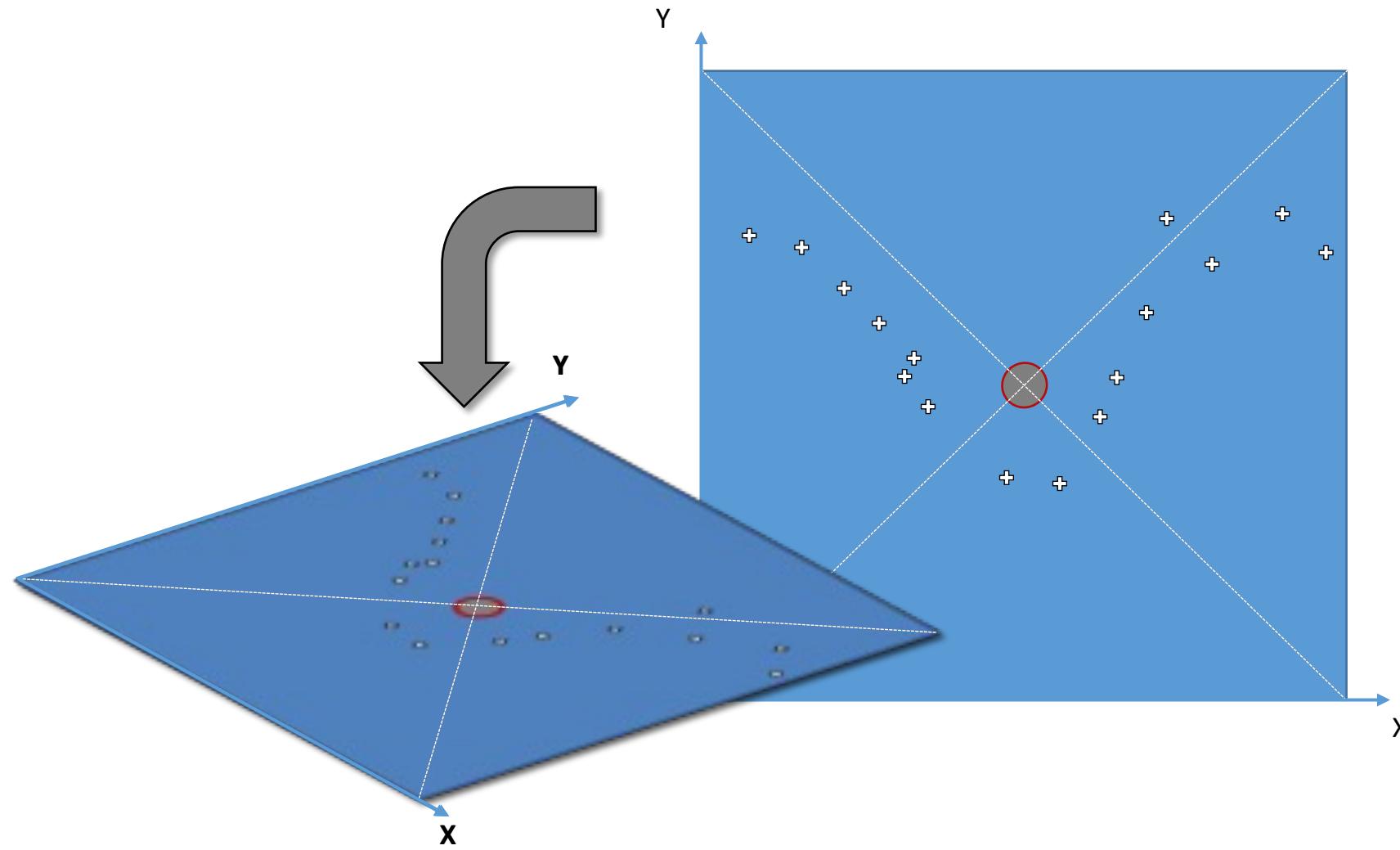
Non-Linear SVR



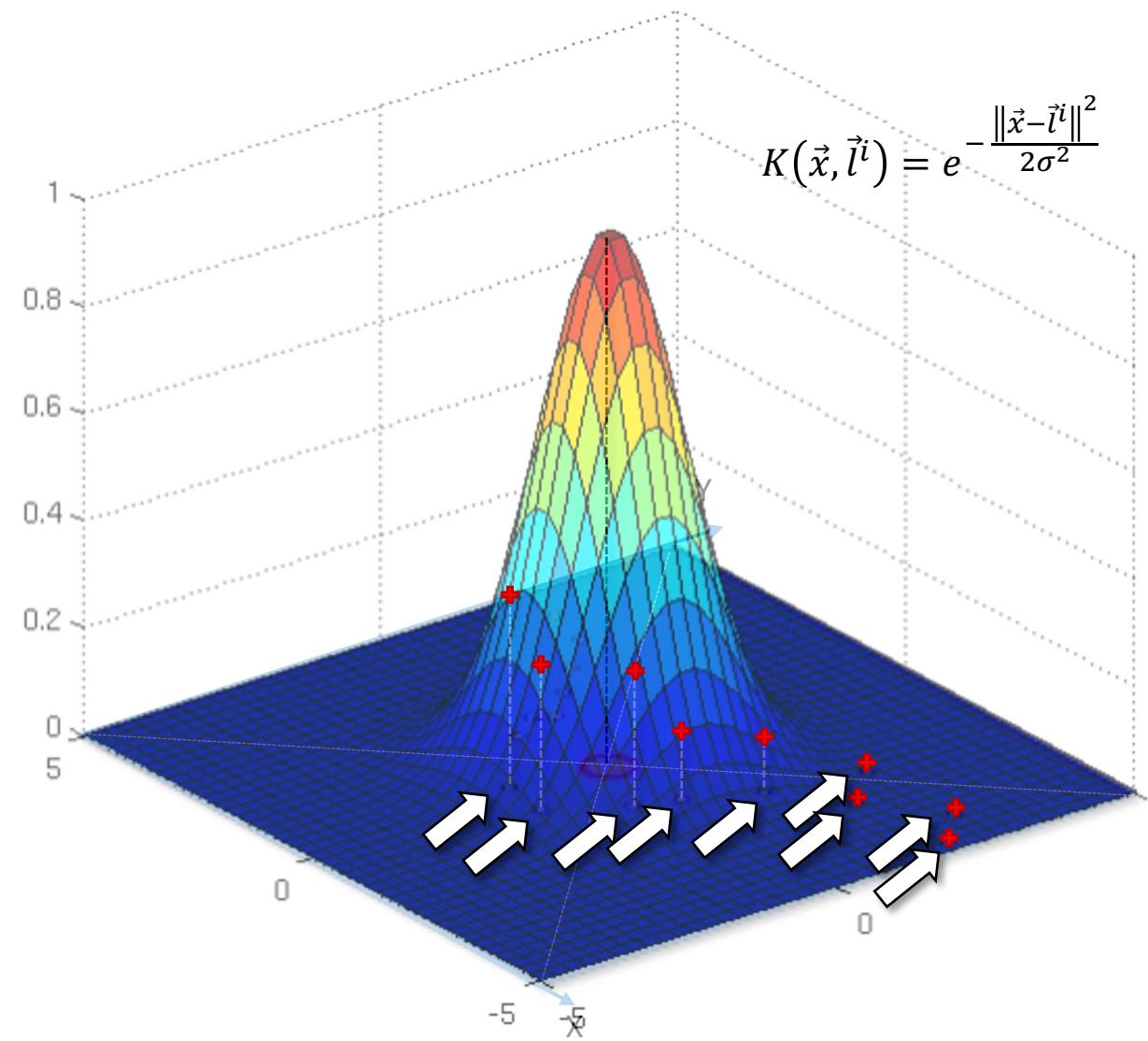
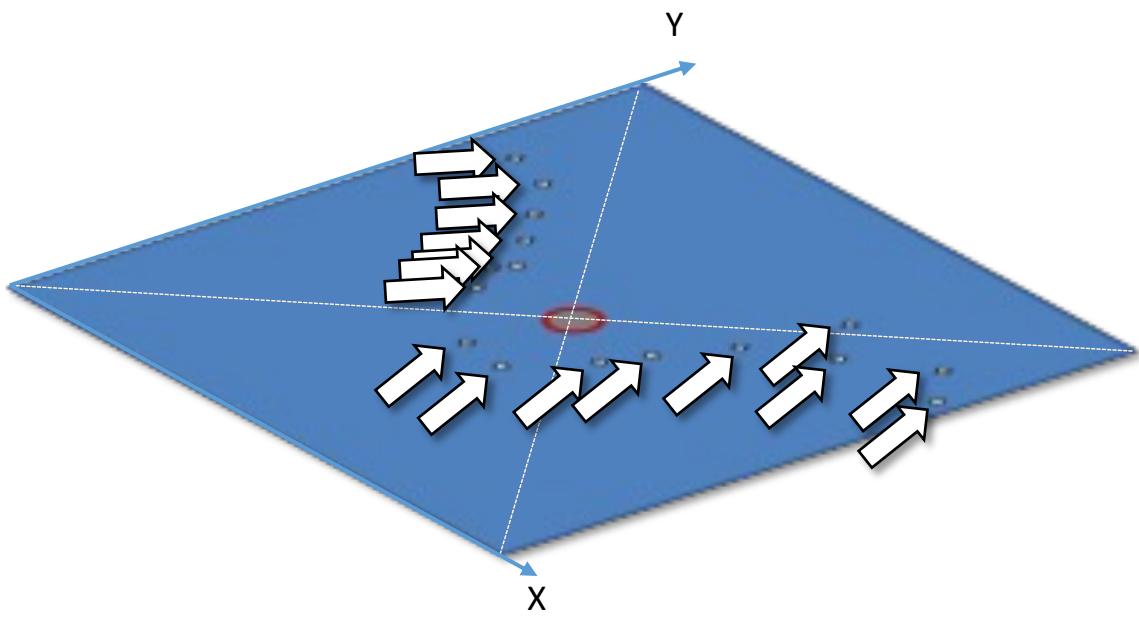
Non-Linear SVR



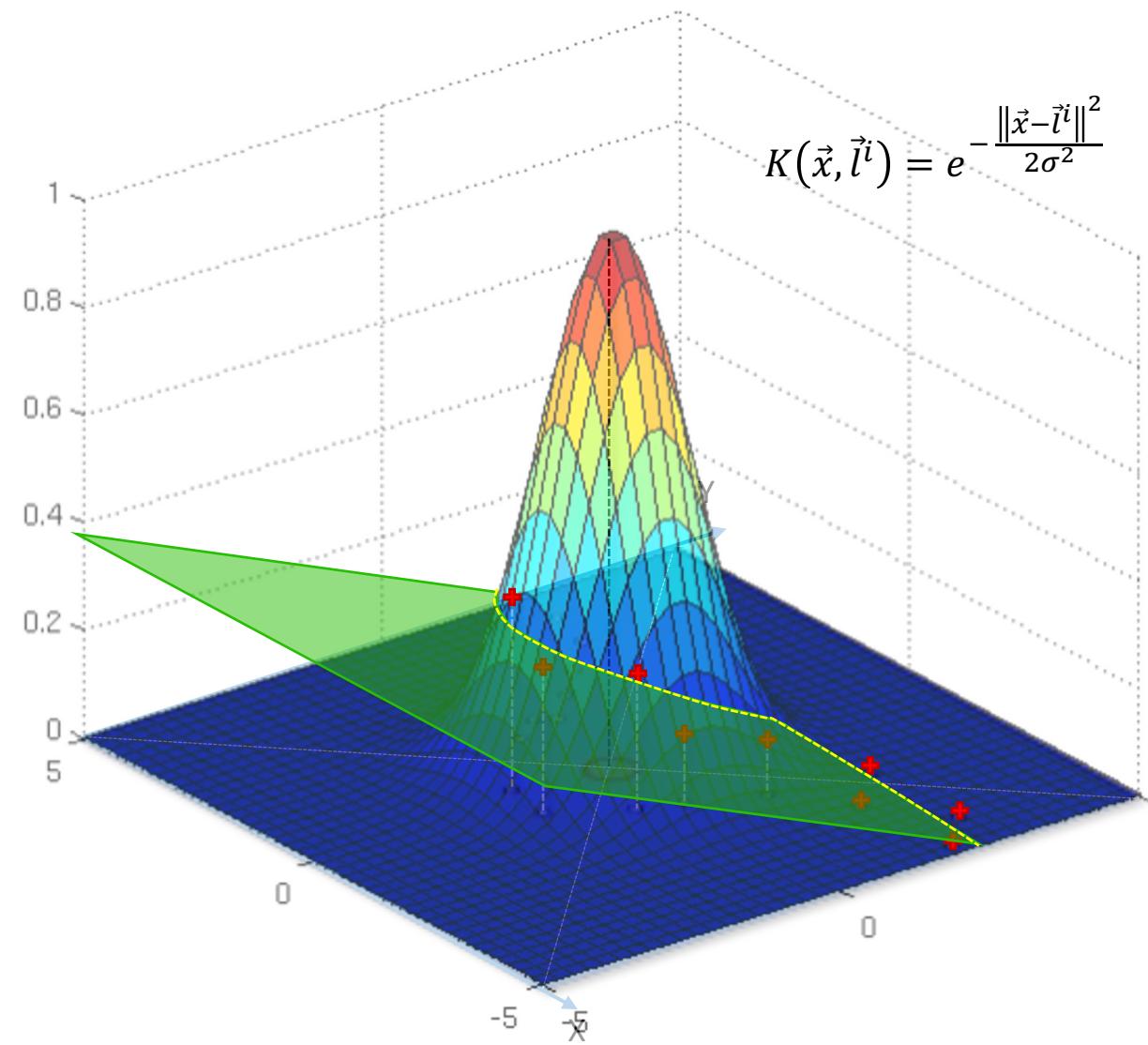
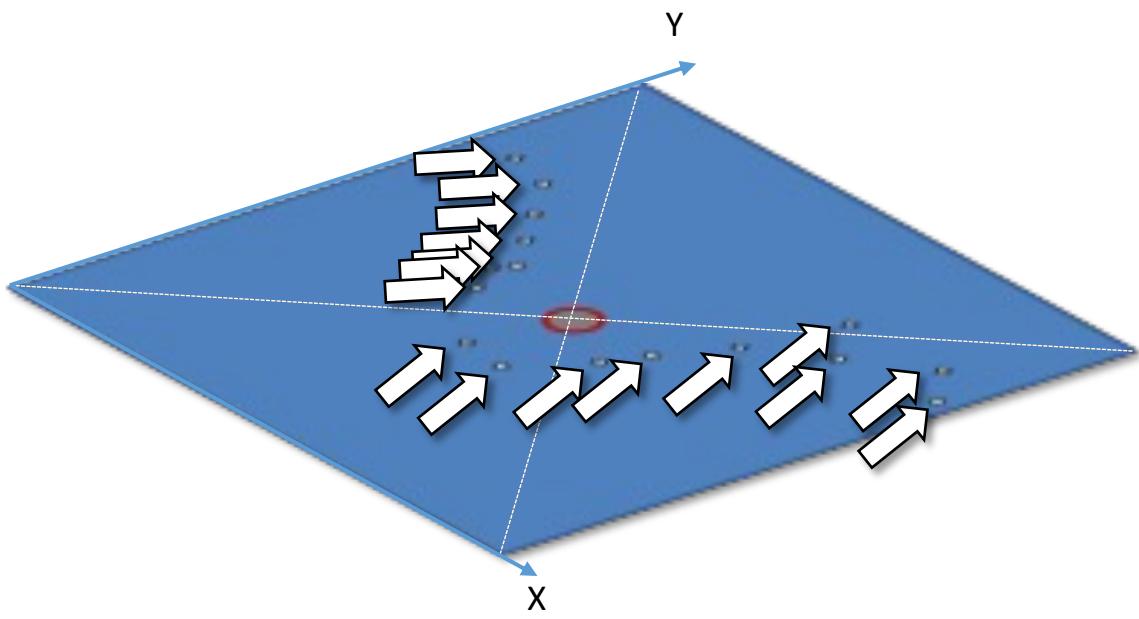
Non-Linear SVR



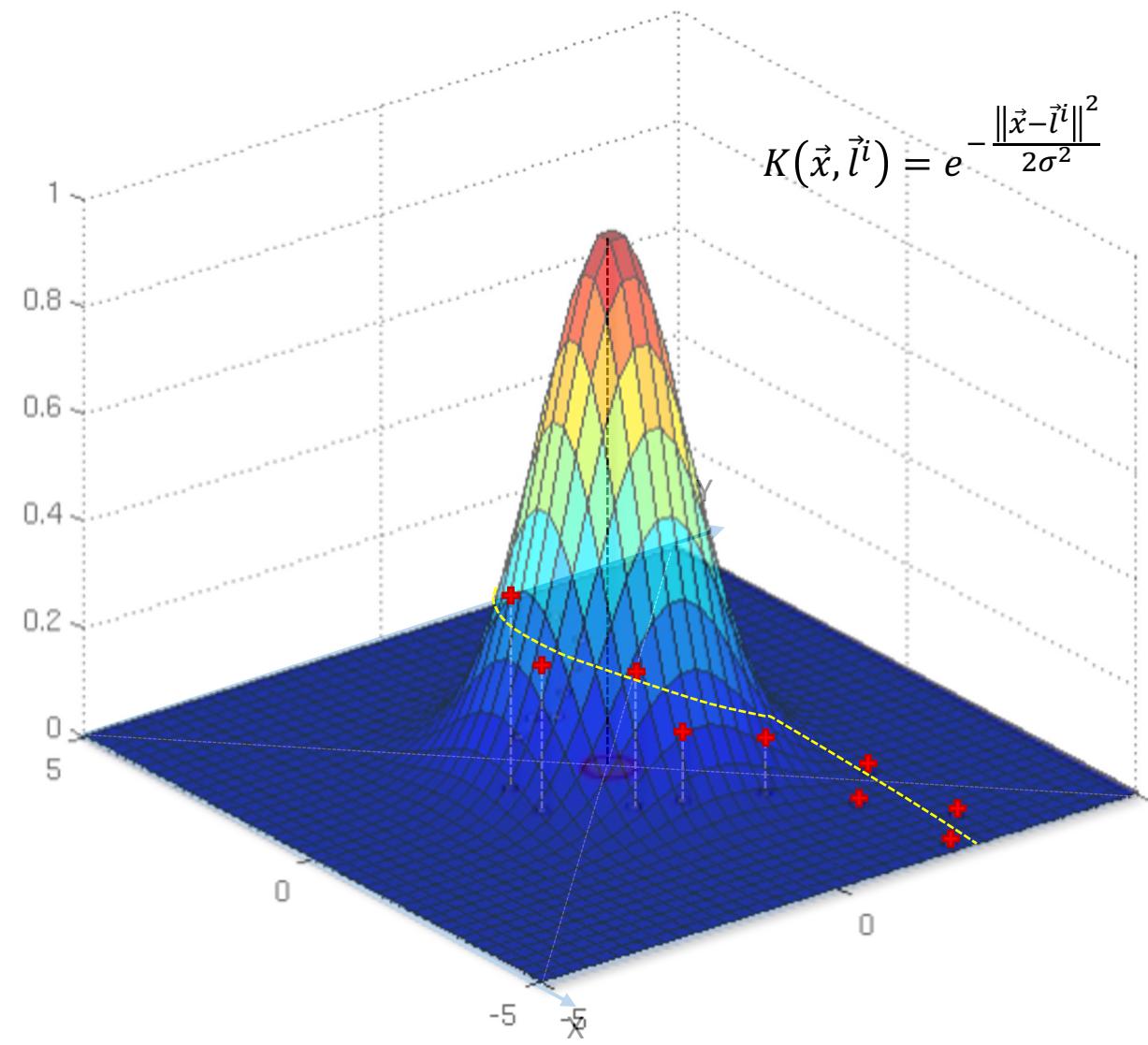
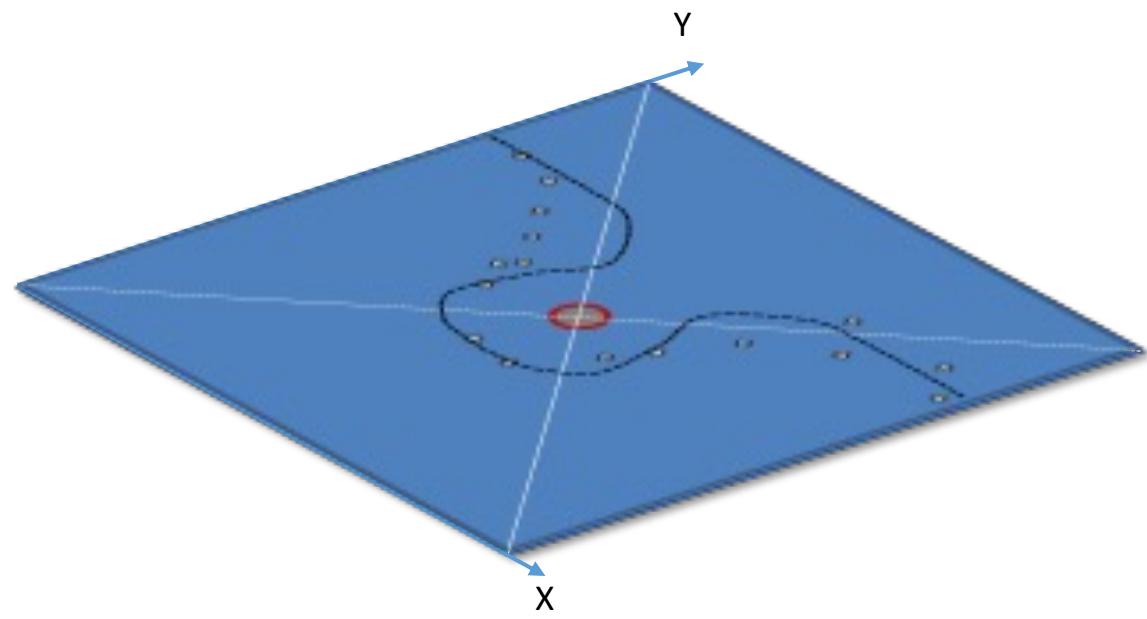
Non-Linear SVR



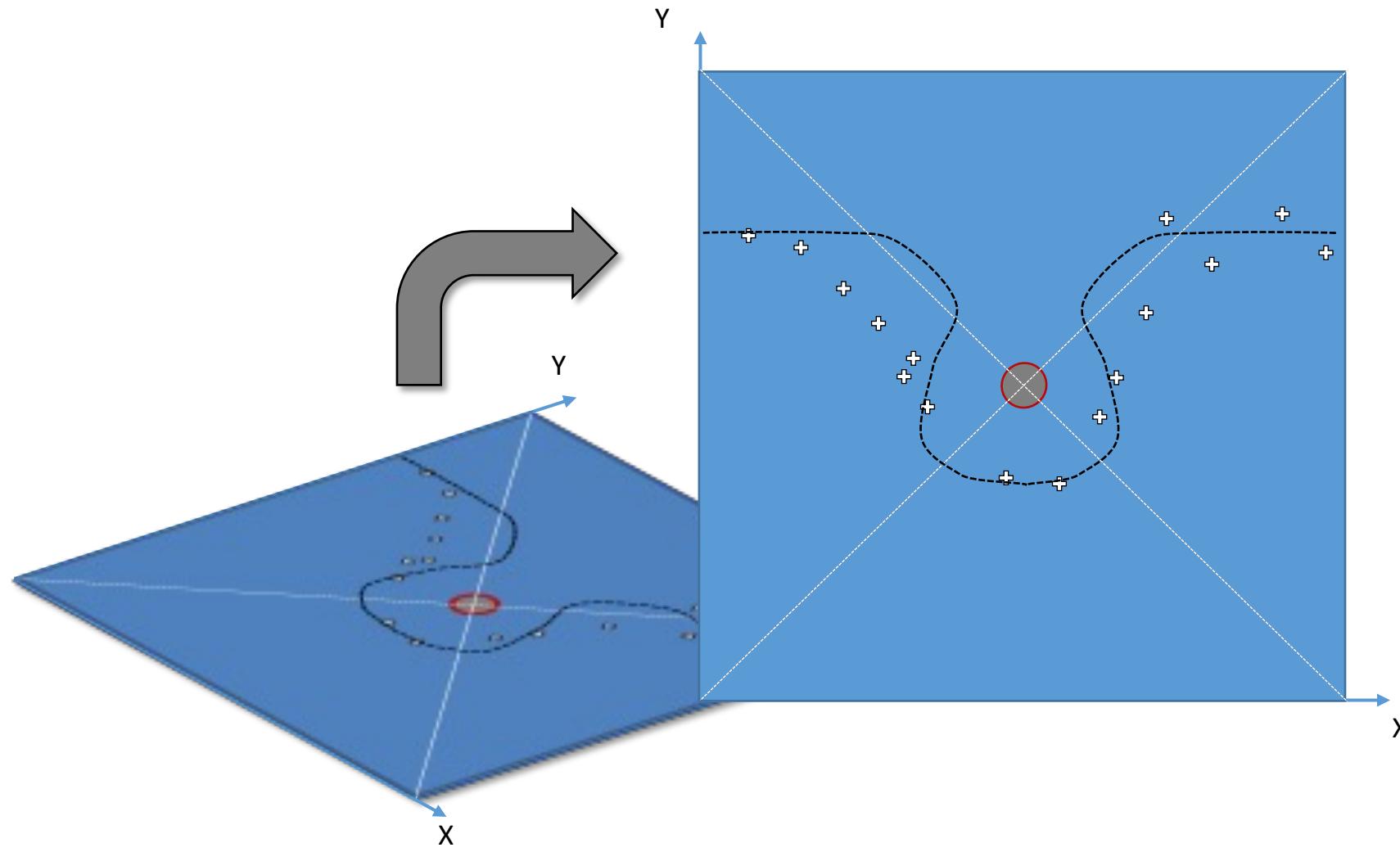
Non-Linear SVR



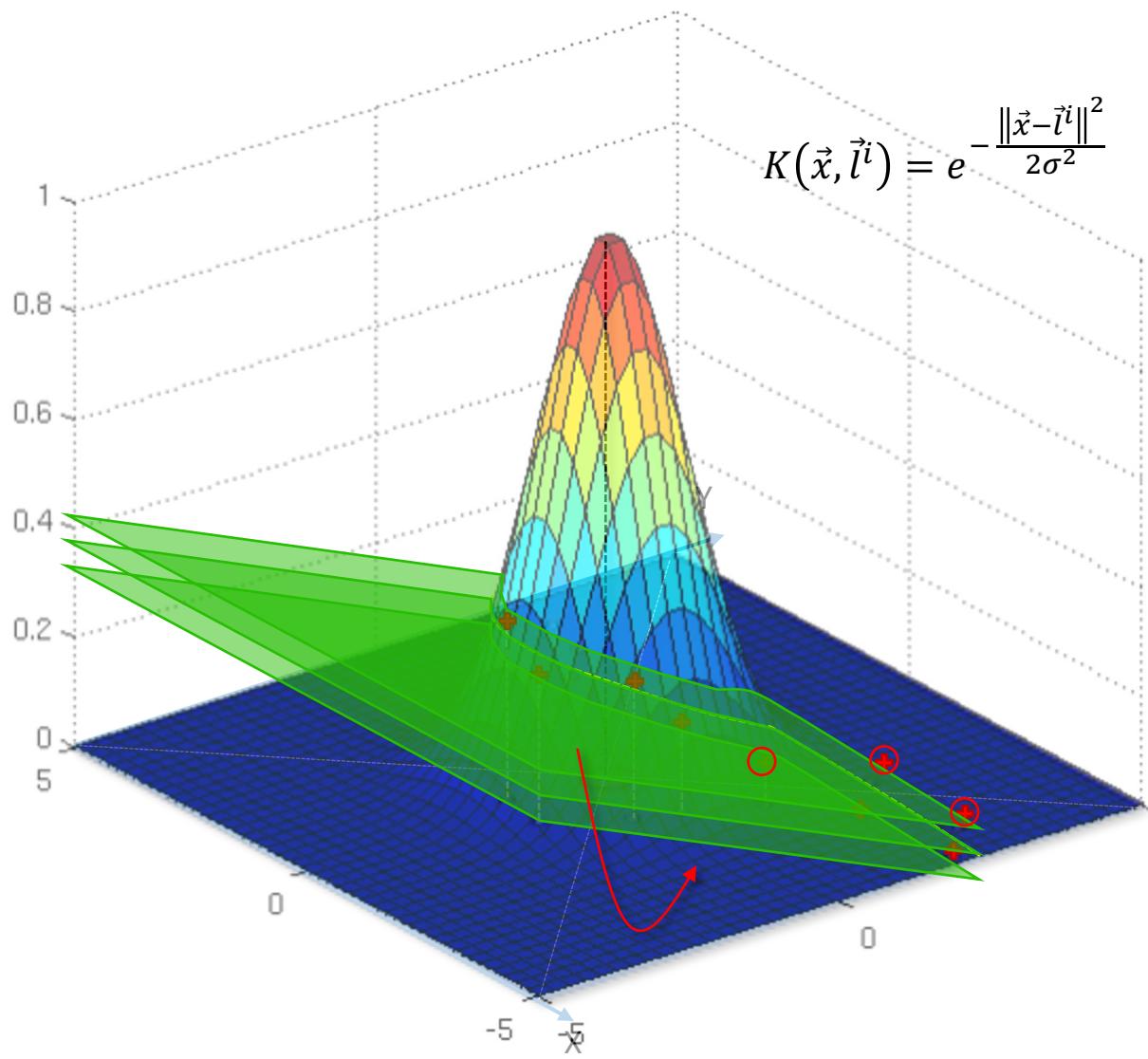
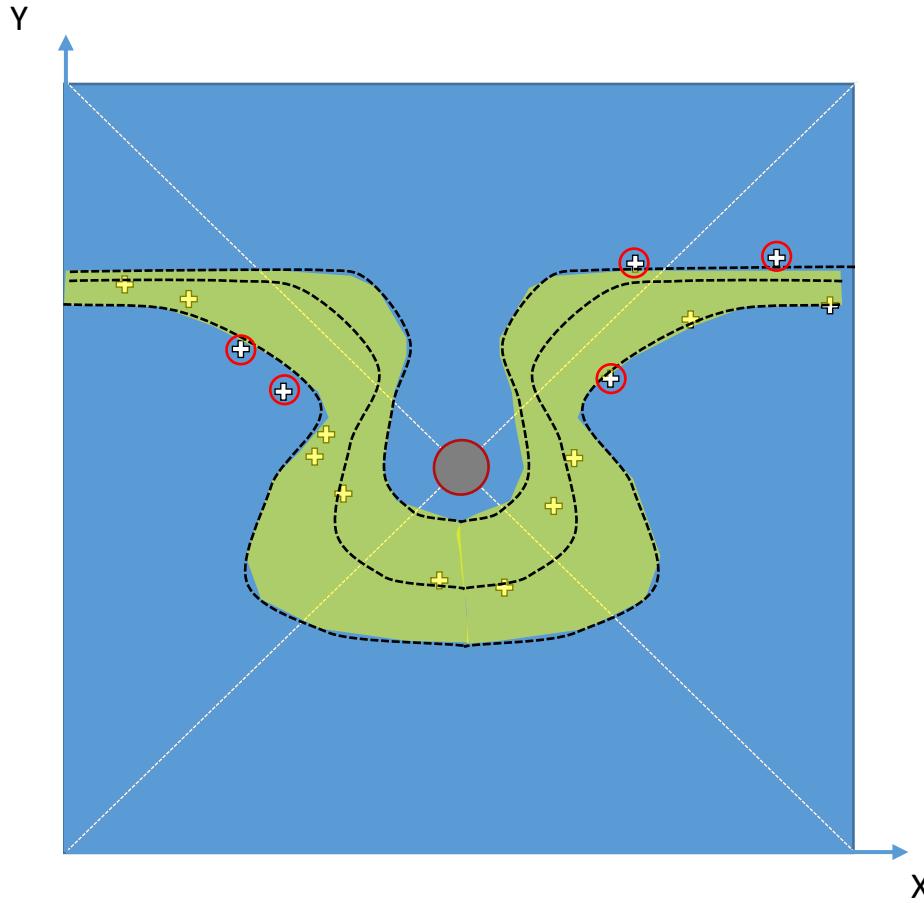
Non-Linear SVR



Non-Linear SVR

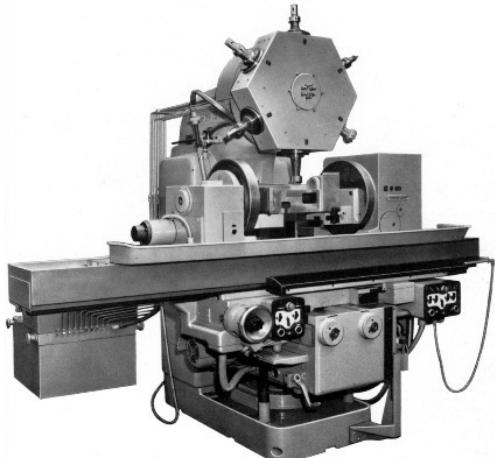


Non-Linear SVR

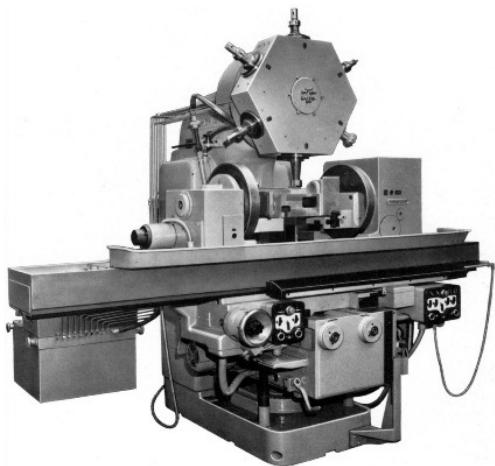


Bayes' Theorem

Bayes Theorem



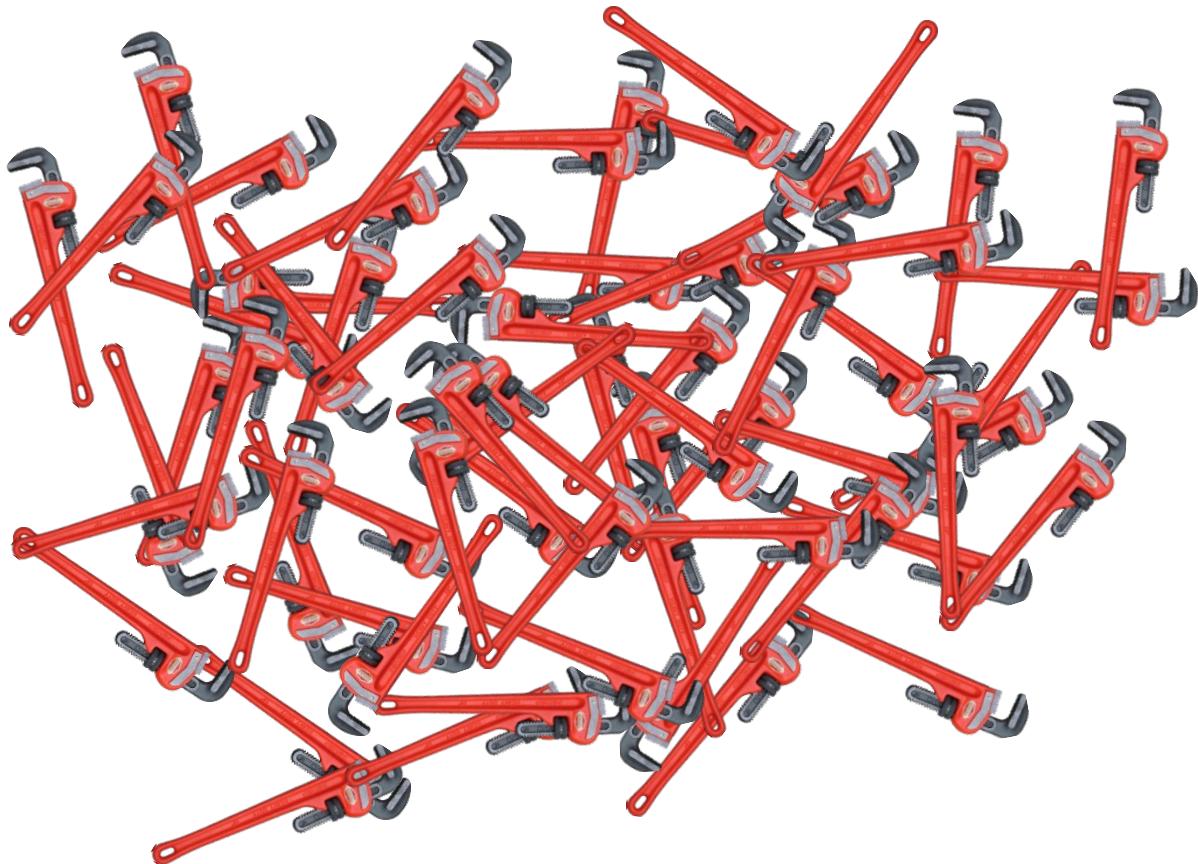
m1 m1



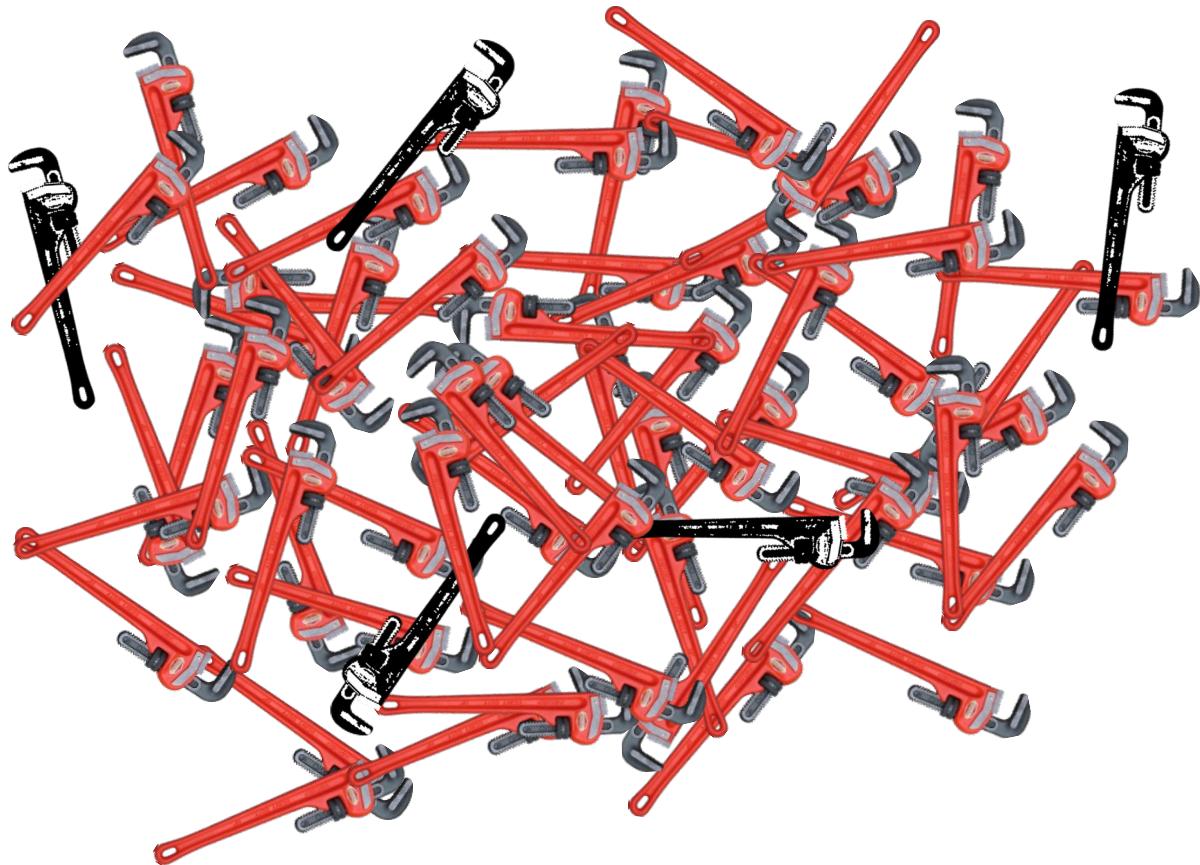
m2 m2 m2 m2 m2 m2 m2 m2 m2



Bayes Theorem

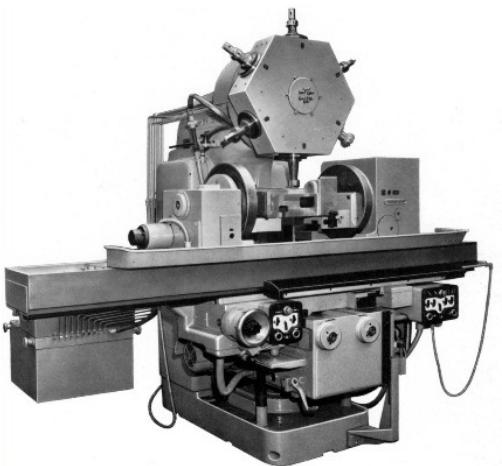


Bayes Theorem



Bayes Theorem

What's the probability?

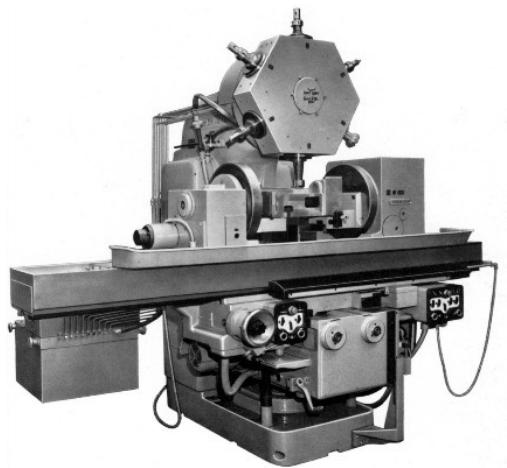


m2



Bayes Theorem

What's the probability?



m2



Bayes Theorem

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Bayes Theorem

Mach1: 30 wrenches / hr

Mach2: 20 wrenches / hr

Out of all produced parts:

We can SEE that 1% are defective

Out of all defective parts:

We can SEE that 50% came from mach1

And 50% came from mach2

Question:

What is the probability that a part
produced by mach2 is defective = ?

Bayes Theorem

Mach1: 30 wrenches / hr

Mach2: 20 wrenches / hr

$$\rightarrow P(\text{Mach1}) = 30/50 = 0.6$$

$$\rightarrow P(\text{Mach2}) = 20/50 = 0.4$$

Out of all produced parts:

We can SEE that 1% are defective

$$\rightarrow P(\text{Defect}) = 1\%$$

Out of all defective parts:

We can SEE that 50% came from mach1

And 50% came from mach2

$$\rightarrow P(\text{Mach1} | \text{Defect}) = 50\%$$

$$\rightarrow P(\text{Mach2} | \text{Defect}) = 50\%$$

Question:

What is the probability that a part
produced by mach2 is defective = ?

$$\rightarrow P(\text{Defect} | \text{Mach2}) = ?$$

Bayes Theorem

Mach1: 30 wrenches / hr

Mach2: 20 wrenches / hr

$$\cancel{\rightarrow P(\text{Mach1}) = 30/50 = 0.6}$$

$$\rightarrow P(\text{Mach2}) = 20/50 = 0.4$$

Out of all produced parts:

We can SEE that 1% are defective

$$\rightarrow P(\text{Defect}) = 1\%$$

Out of all defective parts:

We can SEE that 50% came from mach1

And 50% came from mach2

$$\cancel{\rightarrow P(\text{Mach1} \mid \text{Defect}) = 50\%}$$

$$\rightarrow P(\text{Mach2} \mid \text{Defect}) = 50\%$$

Question:

What is the probability that a part
produced by mach2 is defective = ?

$$\rightarrow P(\text{Defect} \mid \text{Mach2}) = ?$$

Bayes Theorem

Mach1: 30 wrenches / hr

Mach2: 20 wrenches / hr

Out of all produced parts:

We can SEE that 1% are defective

Out of all defective parts:

We can SEE that 50% came from mach1

And 50% came from mach2

Question:

What is the probability that a part
produced by mach2 is defective = ?

$$\rightarrow P(\text{Mach2}) = 20/50 = 0.4$$

$$\rightarrow P(\text{Defect}) = 1\%$$

$$\rightarrow P(\text{Mach2} | \text{Defect}) = 50\%$$

$$\rightarrow P(\text{Defect} | \text{Mach2}) = ?$$

Bayes Theorem

Mach1: 30 wrenches / hr

Mach2: 20 wrenches / hr

Out of all produced parts:

We can SEE that 1% are defective

Out of all defective parts:

We can SEE that 50% came from mach1

And 50% came from mach2

Question:

What is the probability that a part
produced by mach2 is defective = ?

$$\rightarrow P(\text{Mach2}) = 20/50 = 0.4$$

$$\rightarrow P(\text{Defect}) = 1\%$$

$$\rightarrow P(\text{Mach2} | \text{Defect}) = 50\%$$

$$\rightarrow P(\text{Defect} | \text{Mach2}) = ?$$

$$P(\text{Defect} | \text{Mach2}) = \frac{P(\text{Mach2} | \text{Defect}) * P(\text{Defect})}{P(\text{Mach2})}$$

Bayes Theorem

Mach1: 30 wrenches / hr

Mach2: 20 wrenches / hr

Out of all produced parts:

We can SEE that 1% are defective

Out of all defective parts:

We can SEE that 50% came from mach1

And 50% came from mach2

Question:

What is the probability that a part
produced by mach2 is defective = ?

$$\rightarrow P(\text{Mach2}) = 20/50 = 0.4$$

$$\rightarrow P(\text{Defect}) = 1\%$$

$$\rightarrow P(\text{Mach2} | \text{Defect}) = 50\%$$

$$\rightarrow P(\text{Defect} | \text{Mach2}) = ?$$

$$P(\text{Defect} | \text{Mach2}) = \frac{0.5 * 0.01}{0.4}$$

Bayes Theorem

Mach1: 30 wrenches / hr

Mach2: 20 wrenches / hr

Out of all produced parts:

We can SEE that 1% are defective

Out of all defective parts:

We can SEE that 50% came from mach1

And 50% came from mach2

Question:

What is the probability that a part
produced by mach2 is defective = ?

$$\rightarrow P(\text{Mach2}) = 20/50 = 0.4$$

$$\rightarrow P(\text{Defect}) = 1\%$$

$$\rightarrow P(\text{Mach2} | \text{Defect}) = 50\%$$

$$\rightarrow P(\text{Defect} | \text{Mach2}) = ?$$

$$P(\text{Defect} | \text{Mach2}) = \frac{0.5 * 0.01}{0.4} = 0.0125 = 1.25\%$$

It's intuitive!

$$P(\text{Defect} \mid \text{Mach2}) = \frac{P(\text{Mach2} \mid \text{Defect}) * P(\text{Defect})}{P(\text{Mach2})} = 1.25\%$$

Let's look at an example:

- 1000 wrenches
- 400 came from Mach2
- 1% have a defect = 10
- of them 50% came from Mach2 = 5
- % defective parts from Mach2 = $5/400 = 1.25\%$

It's intuitive!

Obvious question:

If the items are labeled, why couldn't we just count the number of defective wrenches that came from Mach2 and divide by the total number that came from Mach2?

Bayes Theorem

Quick exercise:

$$P(\text{Defect} \mid \text{Mach1}) = ?$$

Bayes Theorem

Mach1: 30 wrenches / hr

Mach2: 20 wrenches / hr

$$\rightarrow P(\text{Mach1}) = 30/50 = 0.6$$

$$\rightarrow P(\text{Mach2}) = 20/50 = 0.4$$

Out of all produced parts:

We can SEE that 1% are defective

$$\rightarrow P(\text{Defect}) = 1\%$$

Out of all defective parts:

We can SEE that 50% came from mach1

And 50% came from mach2

$$\rightarrow P(\text{Mach1} | \text{Defect}) = 50\%$$

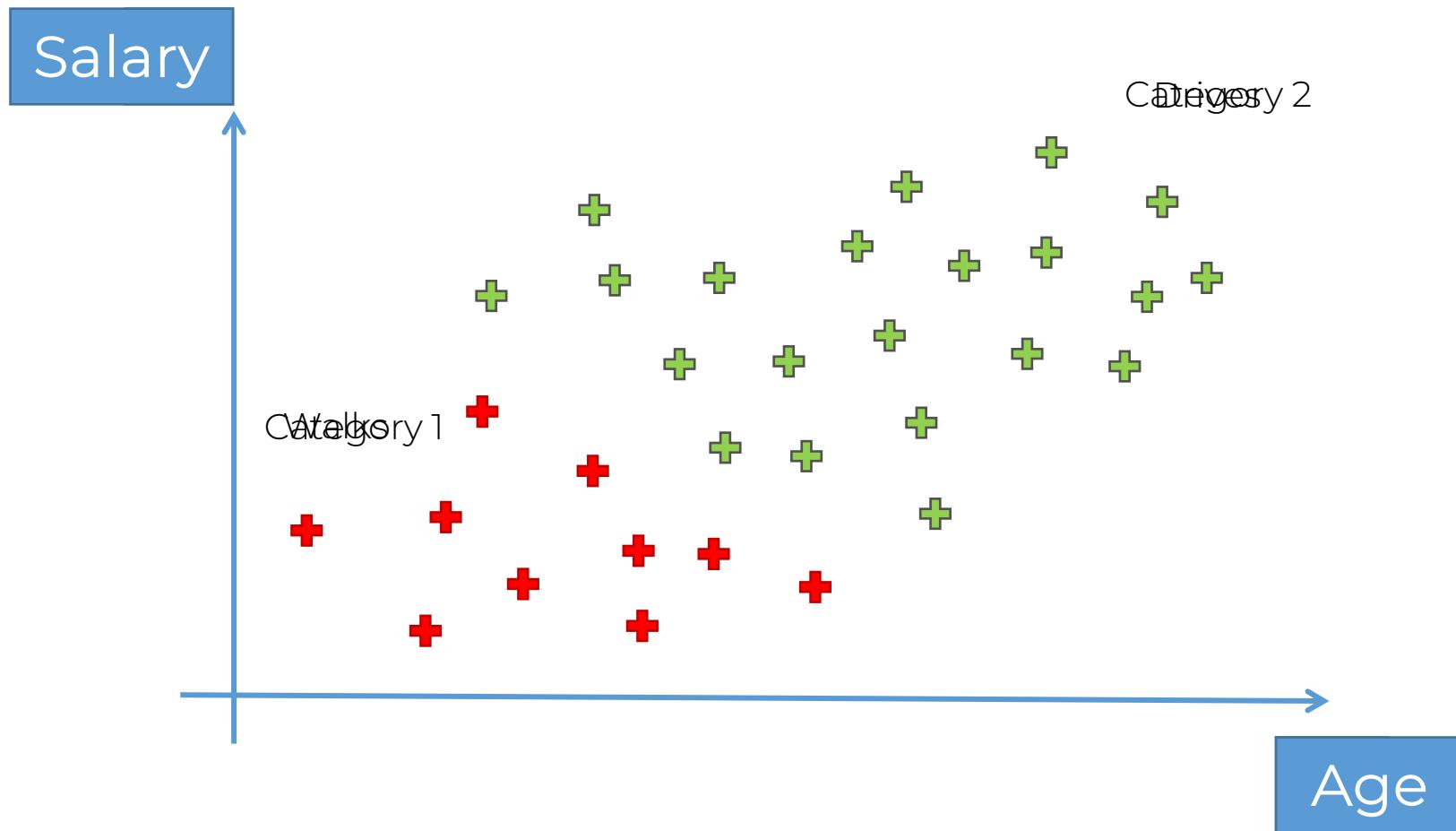
$$\rightarrow P(\text{Mach2} | \text{Defect}) = 50\%$$

Naïve Bayes Classifier Intuition

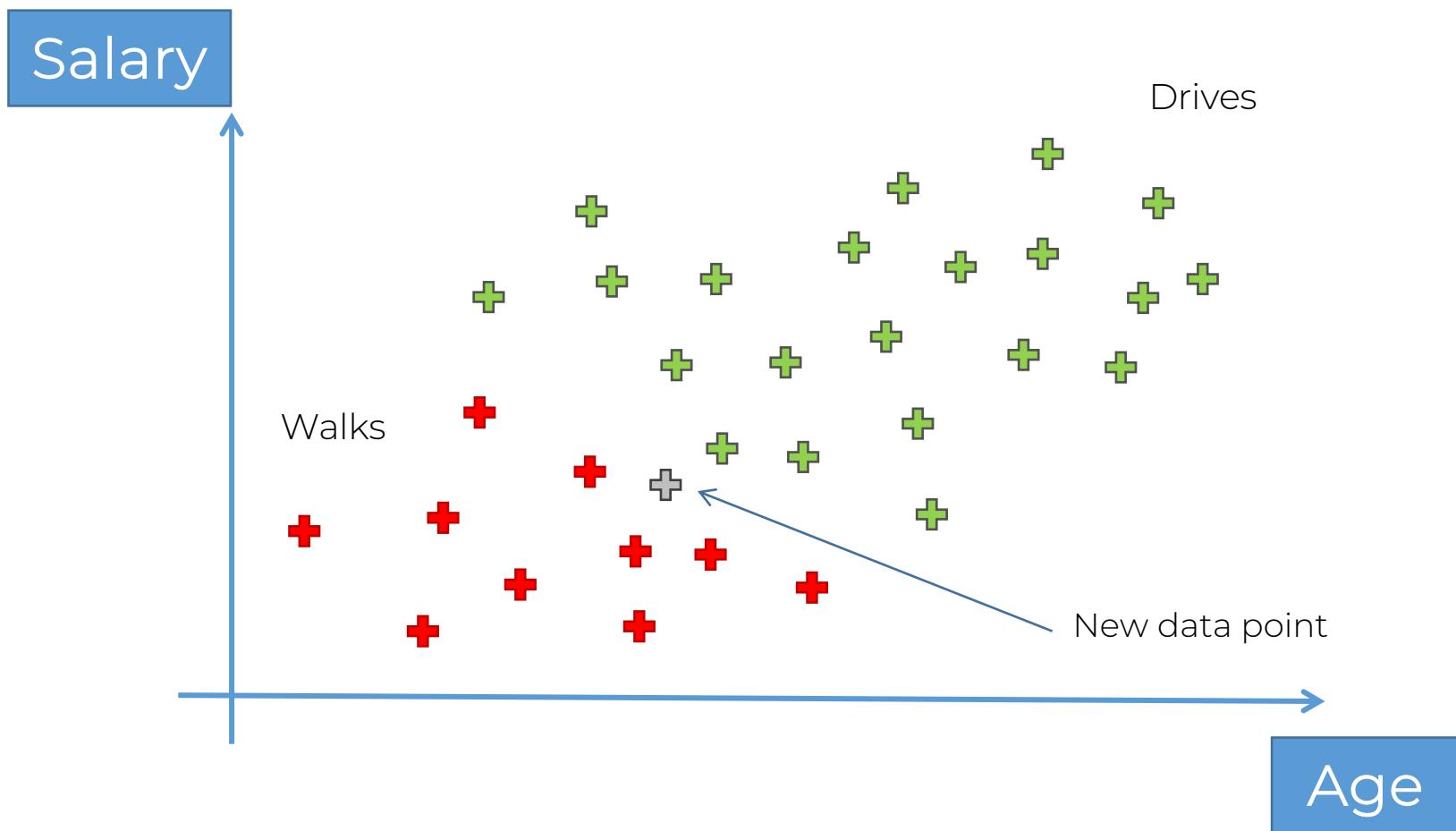
Naïve Bayes

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Naïve Bayes



Naïve Bayes



Naïve Bayes

Plan of Attack

Naïve Bayes

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Step 1

$$P(Walks|X) = \frac{P(X|Walks) * P(Walks)}{P(X)}$$

#4 Posterior Probability

#3 Likelihood

#1 Prior Probability

#2 Marginal Likelihood

```
graph TD; A["#4 Posterior Probability"] --> B["P(Walks|X)"]; C["#3 Likelihood"] --> B; D["#1 Prior Probability"] --> B; E["#2 Marginal Likelihood"] --> B;
```

Step 2

$$P(\\text{Drives}|X) = \\frac{P(X|\\text{Drives}) * P(\\text{Drives})}{P(X)}$$

#4 Posterior Probability

#3 Likelihood

#1 Prior Probability

#2 Marginal Likelihood

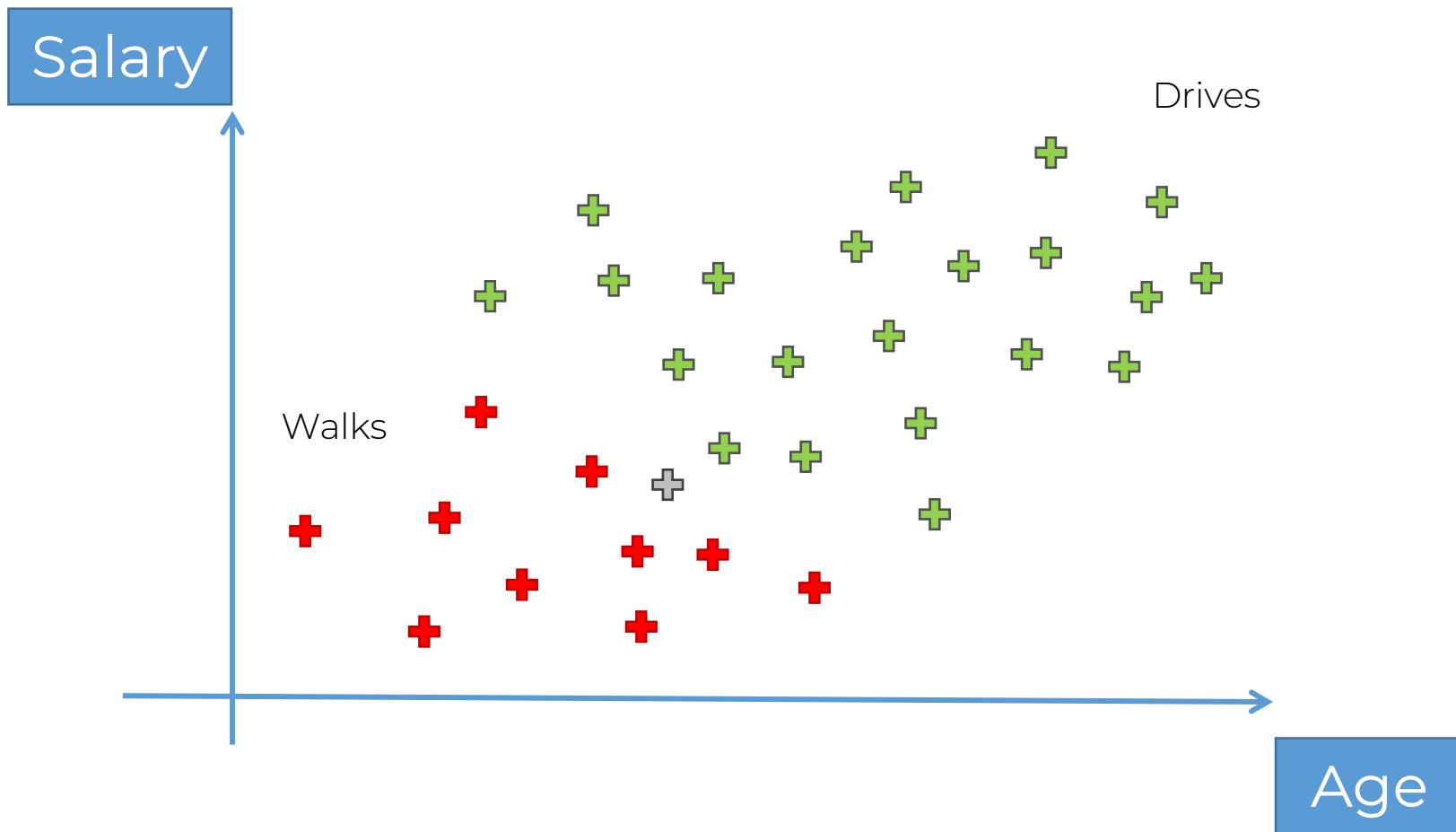
Step 3

$P(\text{Walks}|X)$ v.s. $P(\text{Drives}|X)$

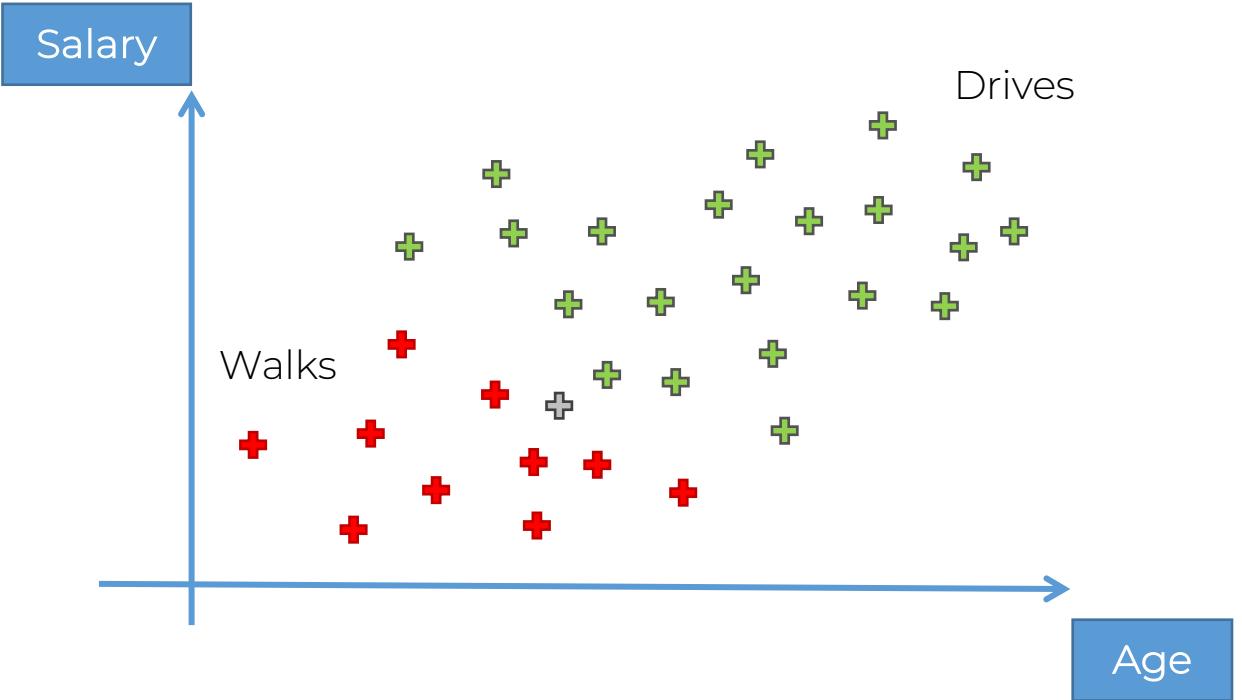
Naïve Bayes

Ready?

Naïve Bayes: Step 1



Naïve Bayes: Step 1



#1. $P(\text{Walks})$

$$P(\text{Walks}) = \frac{\text{Number of Walkers}}{\text{Total Observations}}$$

$$P(\text{Walks}) = \frac{10}{30}$$

Naïve Bayes: Step 1

$$P(Walks|X) = \frac{P(X|Walks) * P(Walks)}{P(X)}$$

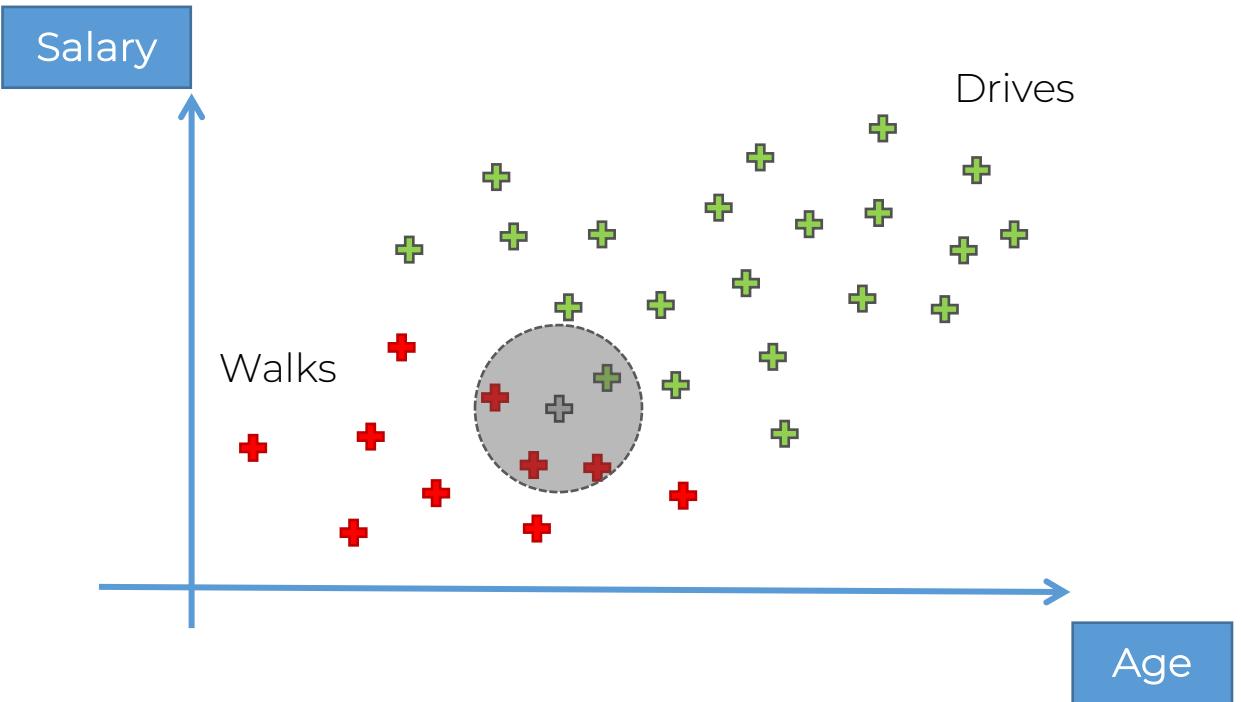
The diagram illustrates the components of the Naïve Bayes formula:

- #4 Posterior Probability
- #3 Likelihood
- #1 Prior Probability (marked with a green checkmark)
- #2 Marginal Likelihood (circled in red)

Arrows point from each component to its corresponding term in the formula:

- #4 Posterior Probability points to $P(Walks|X)$.
- #3 Likelihood points to $P(X|Walks)$.
- #1 Prior Probability points to $P(Walks)$.
- #2 Marginal Likelihood points to $P(X)$.

Naïve Bayes: Step 1

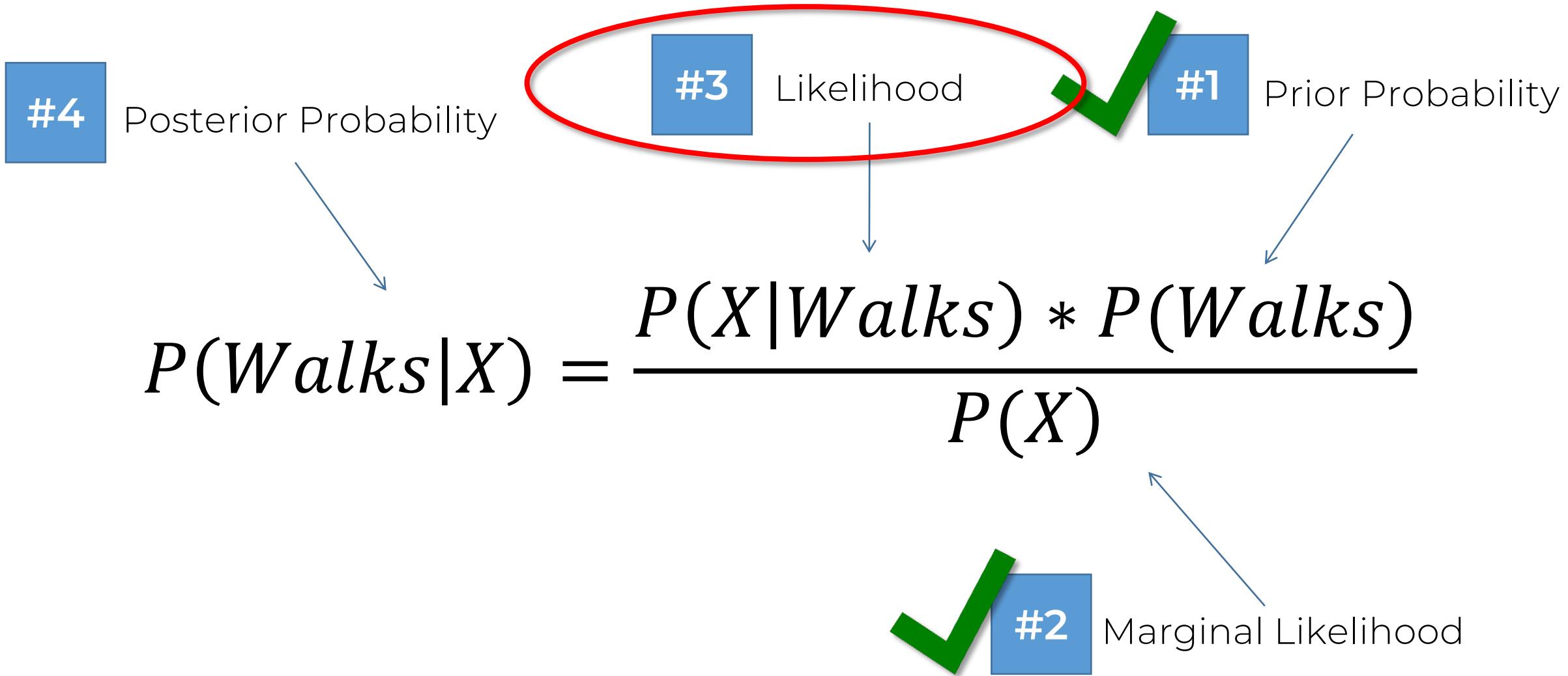


#2. $P(X)$

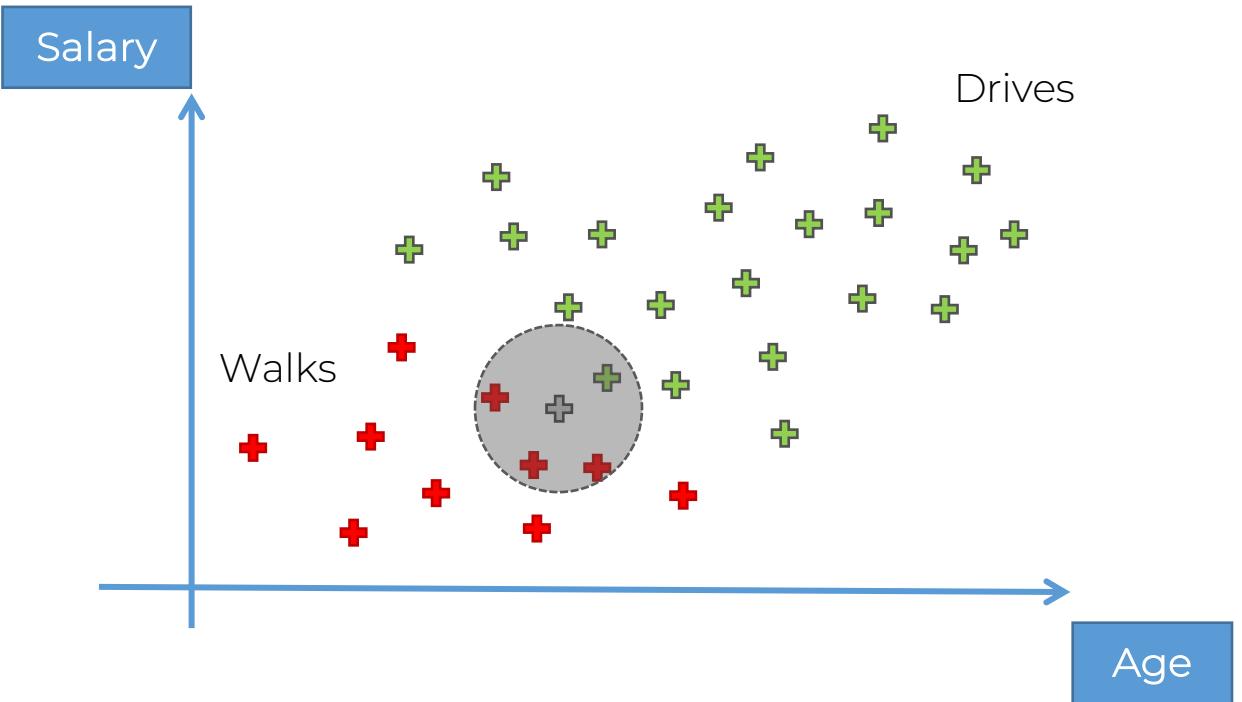
$$P(X) = \frac{\text{Number of Similar Observations}}{\text{Total Observations}}$$

$$P(X) = \frac{4}{30}$$

Naïve Bayes: Step 1



Naïve Bayes: Step 1

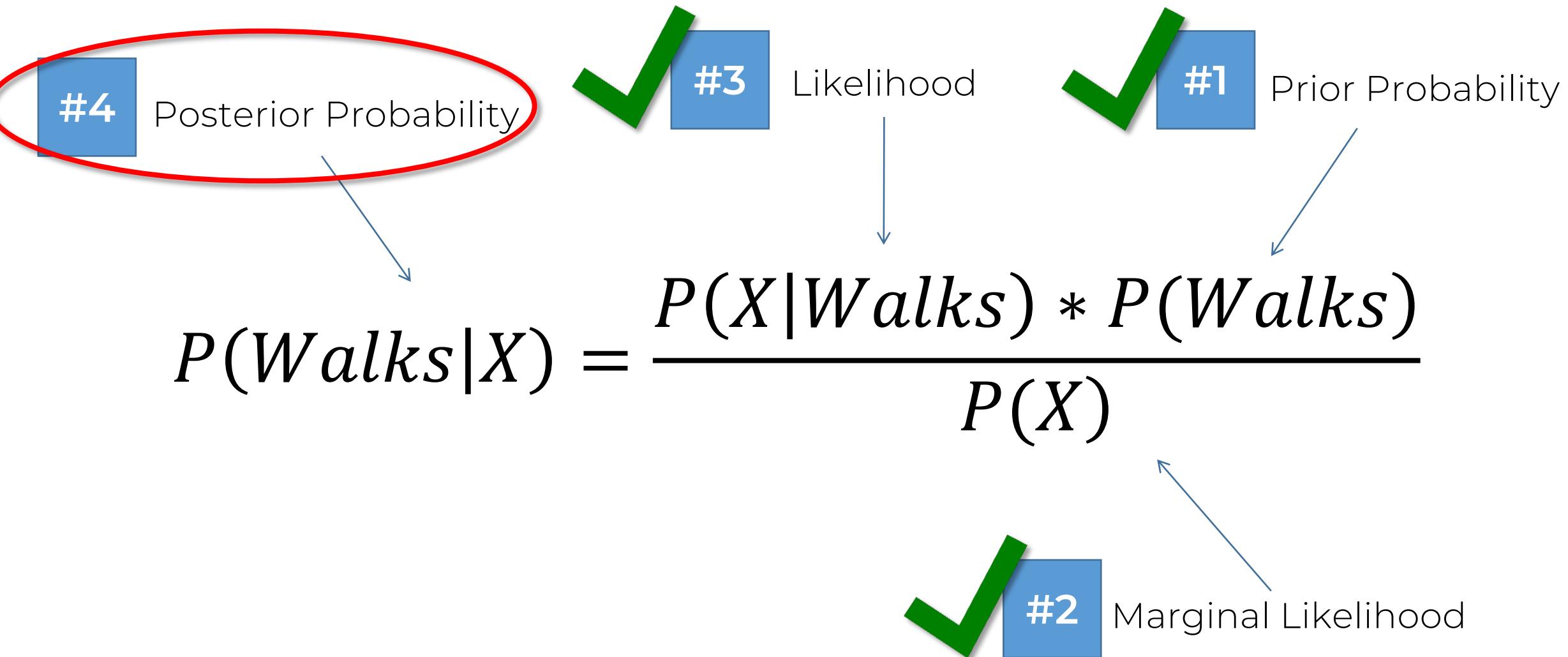


#3. $P(X|Walks)$

Number of Similar Observations

$$P(X|Walks) = \frac{\text{Among those who Walk}}{\text{Total number of Walkers}}$$
$$P(X|Walks) = \frac{3}{10}$$

Naïve Bayes: Step 1



Naïve Bayes: Step 1

#4

Posterior Probability

#3

Likelihood

#1

Prior Probability

$$P(Walks|X) = \frac{\frac{3}{10} * \frac{10}{30}}{\frac{4}{30}} = 0.75$$

#2

Marginal Likelihood

Naïve Bayes

Step 1 - Done.

Step 2

$$P(\\text{Drives}|X) = \\frac{P(X|\\text{Drives}) * P(\\text{Drives})}{P(X)}$$

#4 Posterior Probability

#3 Likelihood

#1 Prior Probability

#2 Marginal Likelihood

Naïve Bayes: Step 2

#4

Posterior Probability

#3

Likelihood

#1

Prior Probability

$$P(\text{Drives}|X) = \frac{\frac{1}{20} * \frac{20}{30}}{\frac{4}{30}} = 0.25$$

#2

Marginal Likelihood

Naïve Bayes

Step 2 - Done.

Step 3

$P(\text{Walks}|X)$ v.s. $P(\text{Drives}|X)$

Step 3

0.75 v. s. 0.25

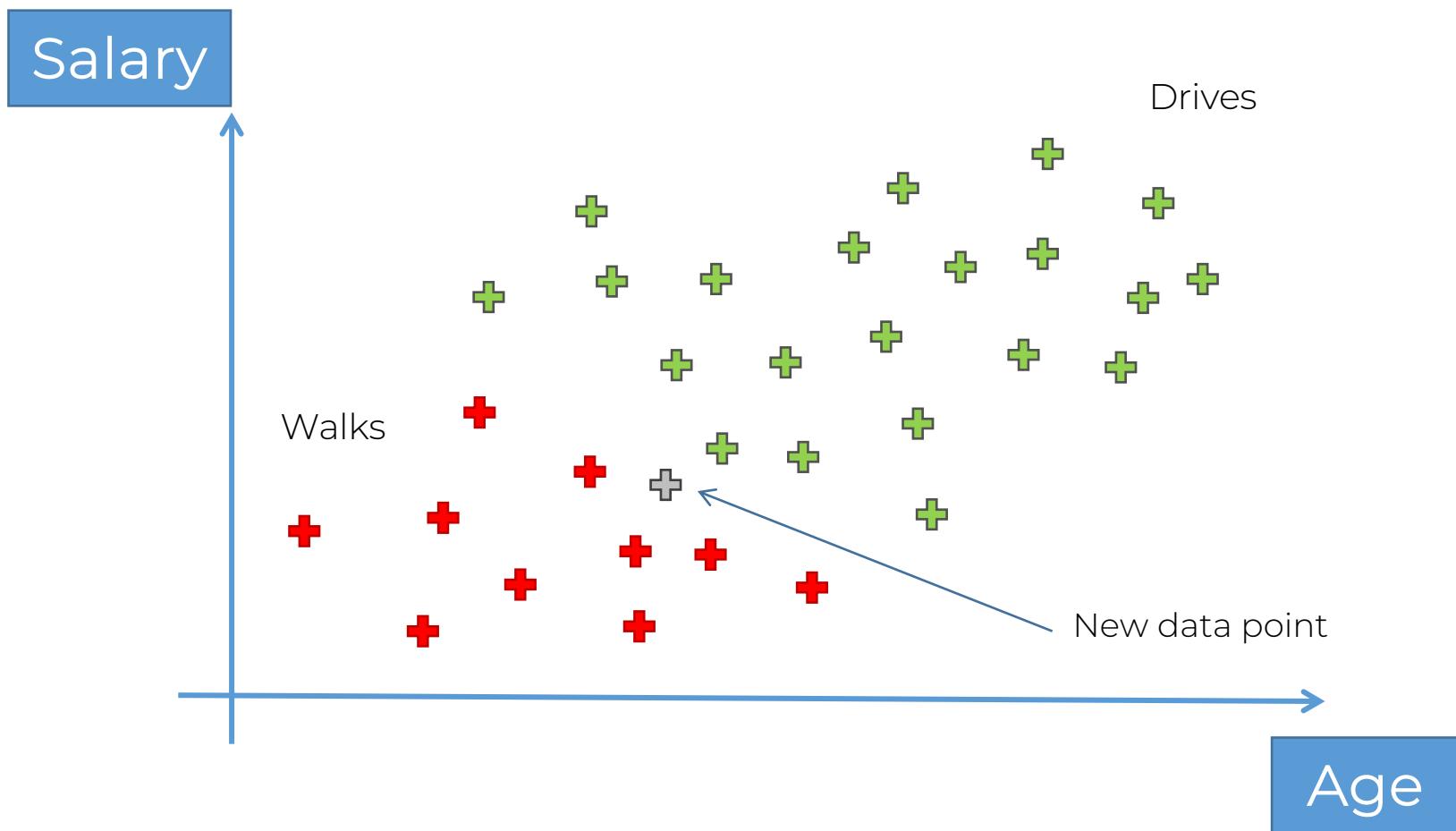
Step 3

$$0.75 > 0.25$$

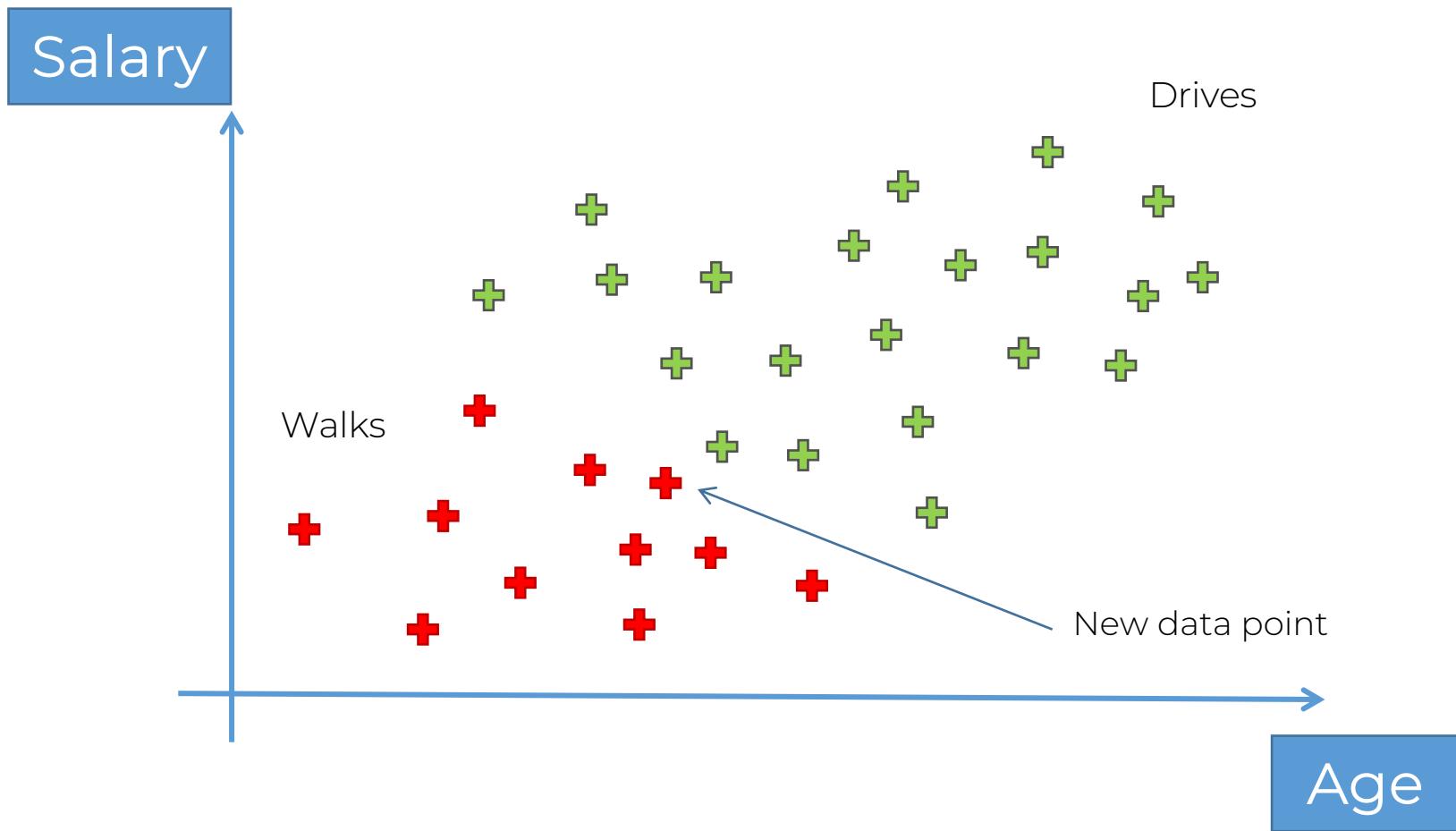
Step 3

$$P(Walks|X) > P(Drives|X)$$

Naïve Bayes



Naïve Bayes



Naïve Bayes Classifier Intuition (Challenge Reveal)

Step 2

$$P(\\text{Drives}|X) = \\frac{P(X|\\text{Drives}) * P(\\text{Drives})}{P(X)}$$

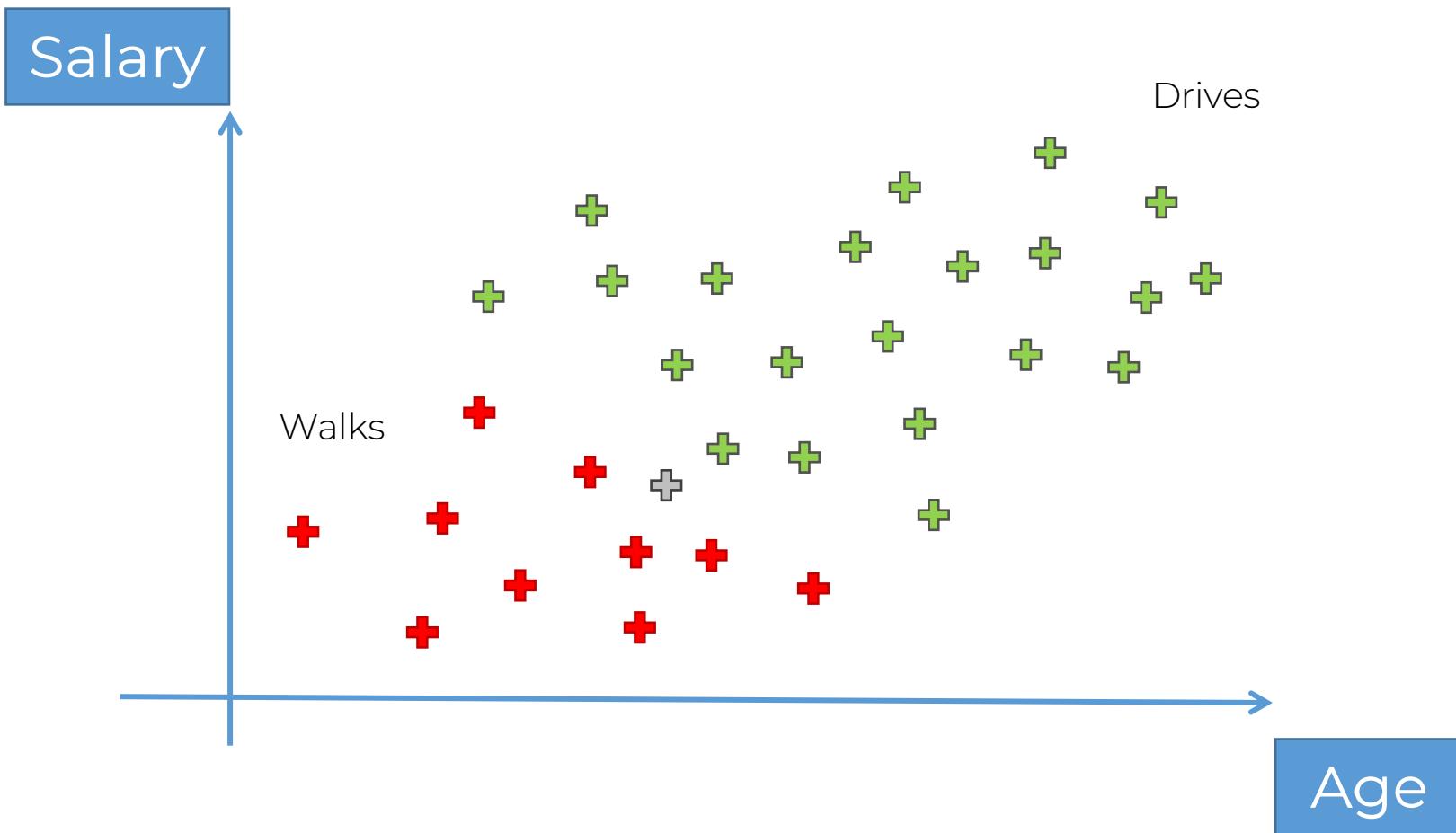
#4 Posterior Probability

#3 Likelihood

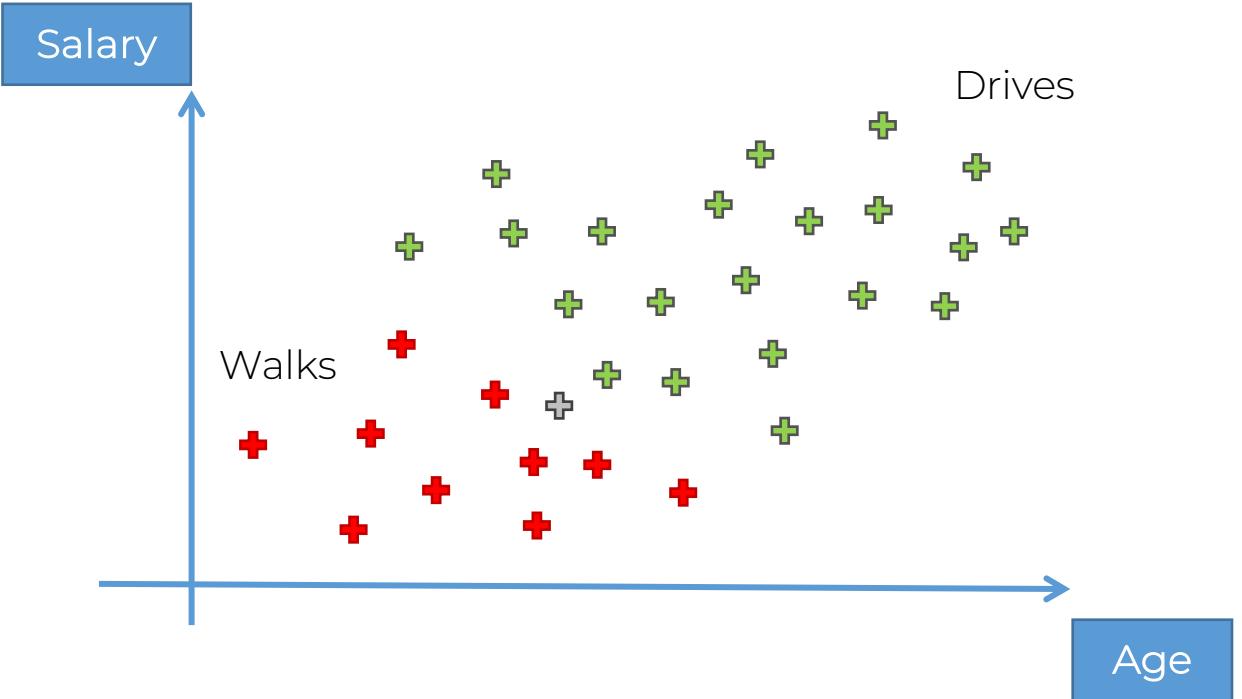
#1 Prior Probability

#2 Marginal Likelihood

Naïve Bayes: Step 2



Naïve Bayes: Step 2



#1. $P(\text{Drives})$

$$P(\text{Drives}) = \frac{\text{Number of Drivers}}{\text{Total Observations}}$$

$$P(\text{Drives}) = \frac{20}{30}$$

Naïve Bayes: Step 2

$$P(\\text{Drives}|X) = \\frac{P(X|\\text{Drives}) * P(\\text{Drives})}{P(X)}$$

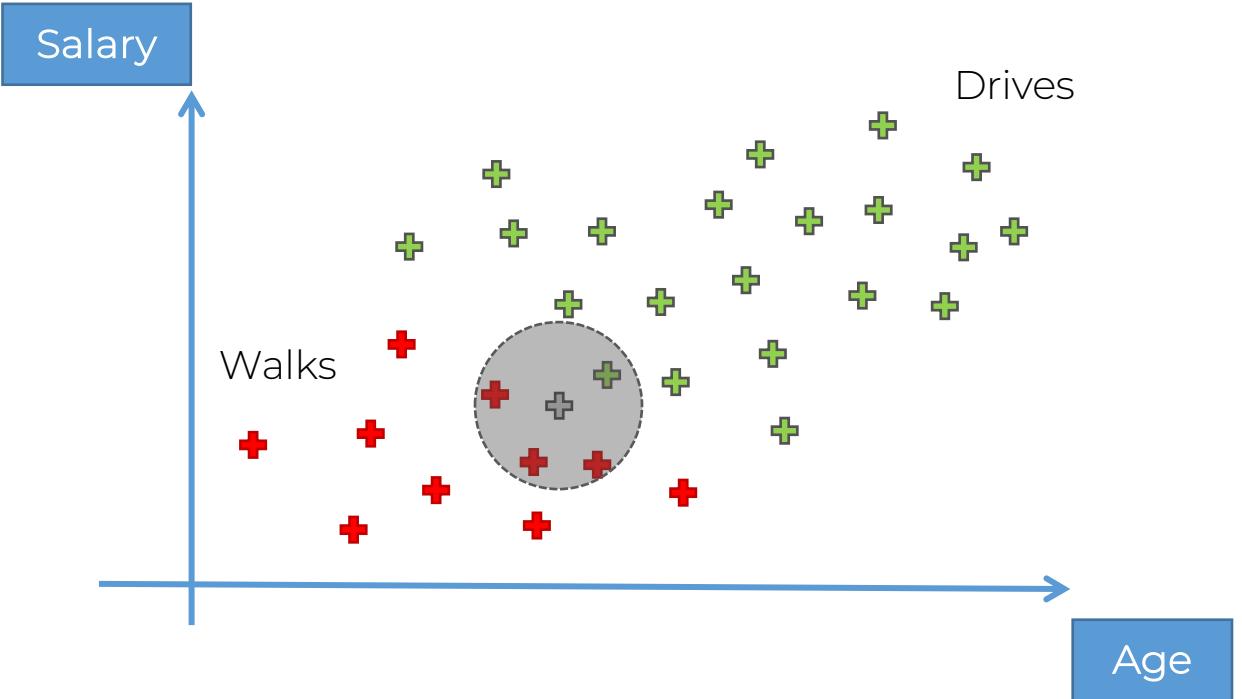
Diagram illustrating the components of the Naïve Bayes formula:

- #4 Posterior Probability
- #3 Likelihood
- #1 Prior Probability
- #2 Marginal Likelihood (circled in red)

Arrows point from the labels to their corresponding terms in the formula:

- #4 Posterior Probability points to $P(\\text{Drives}|X)$
- #3 Likelihood points to $P(X|\\text{Drives})$
- #1 Prior Probability points to $P(\\text{Drives})$
- #2 Marginal Likelihood points to $P(X)$

Naïve Bayes: Step 2



#2. $P(X)$

$$P(X) = \frac{\text{Number of Similar Observations}}{\text{Total Observations}}$$

$$P(X) = \frac{4}{30}$$