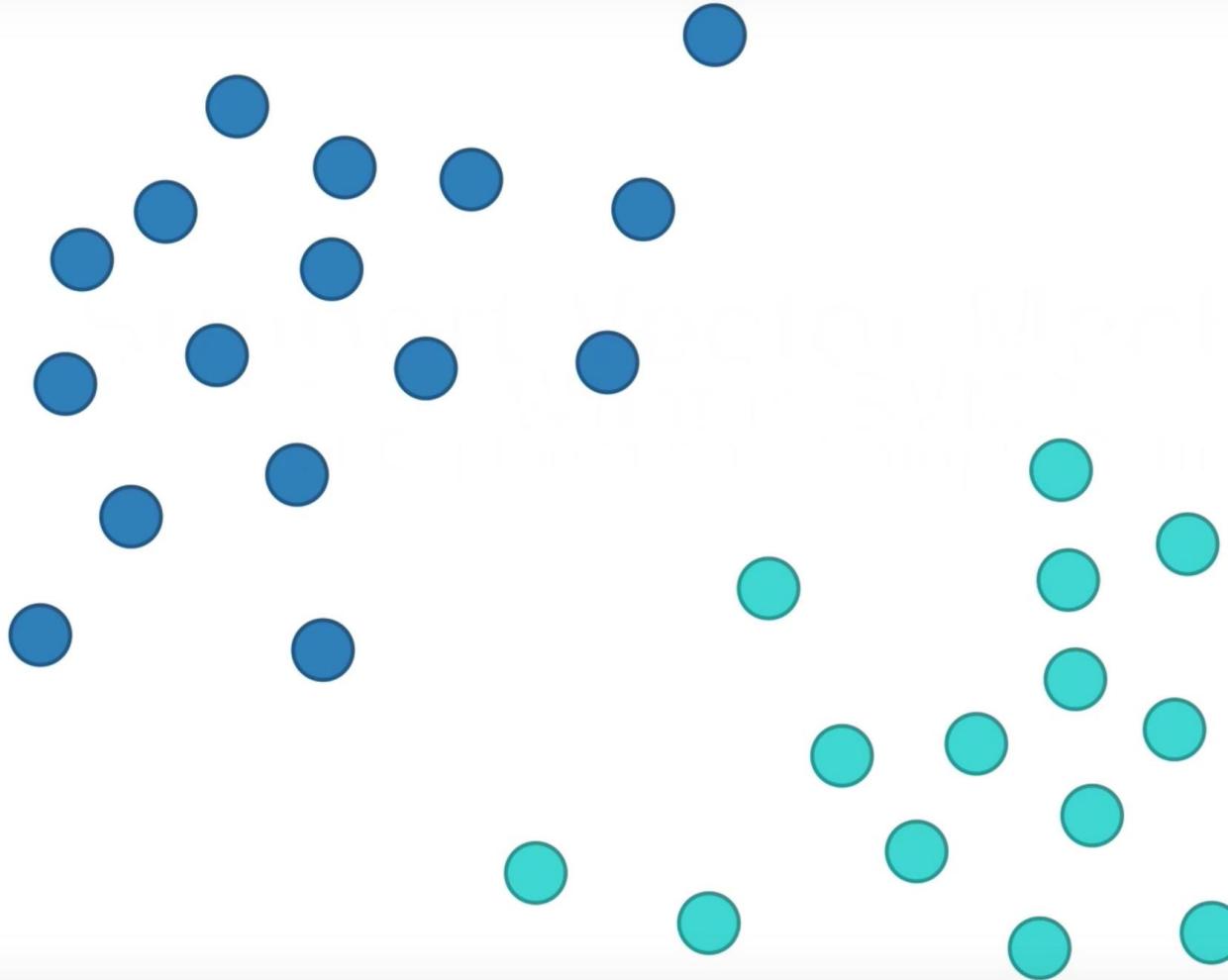
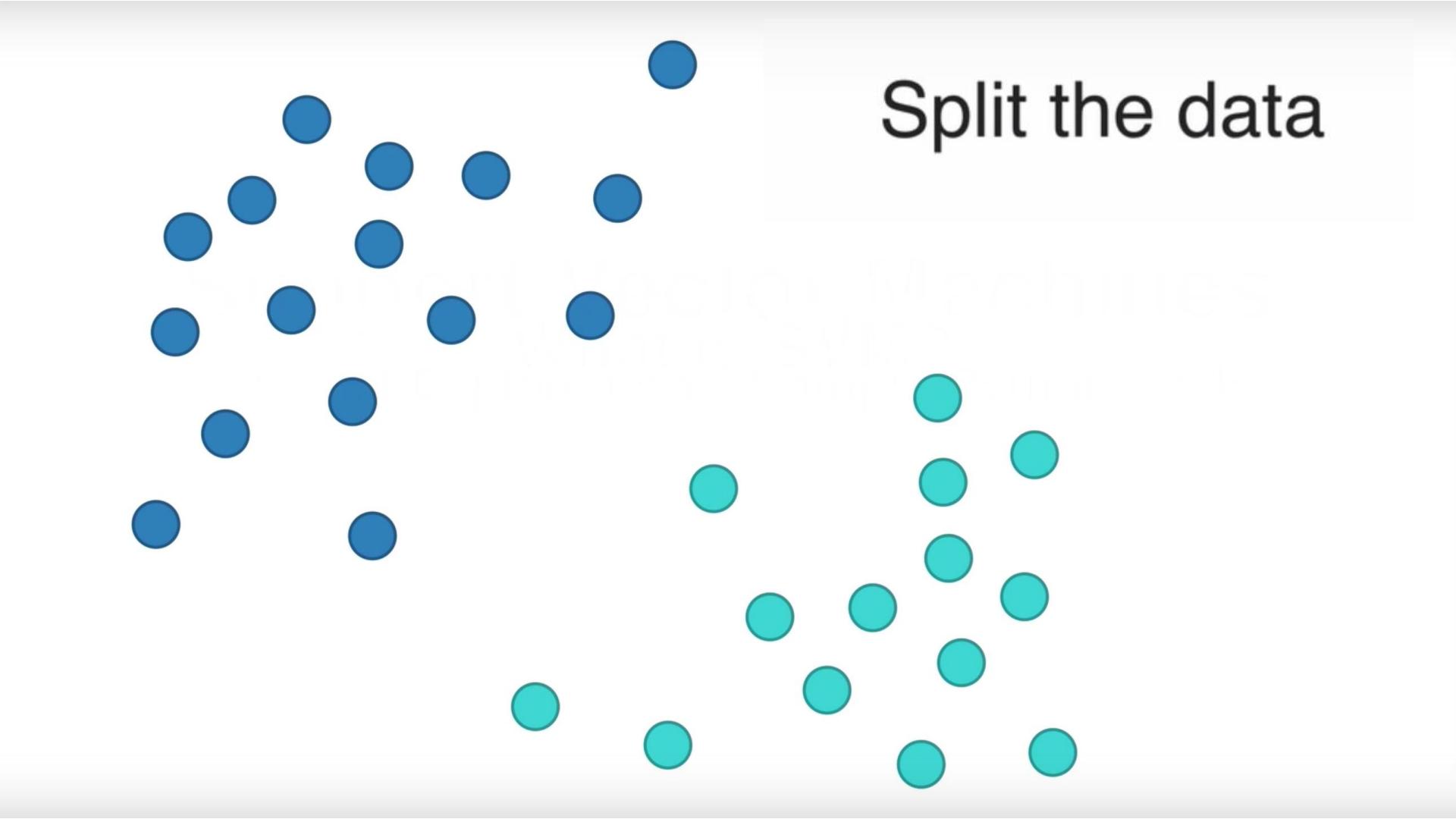


Support Vector Machine

What is SVM?

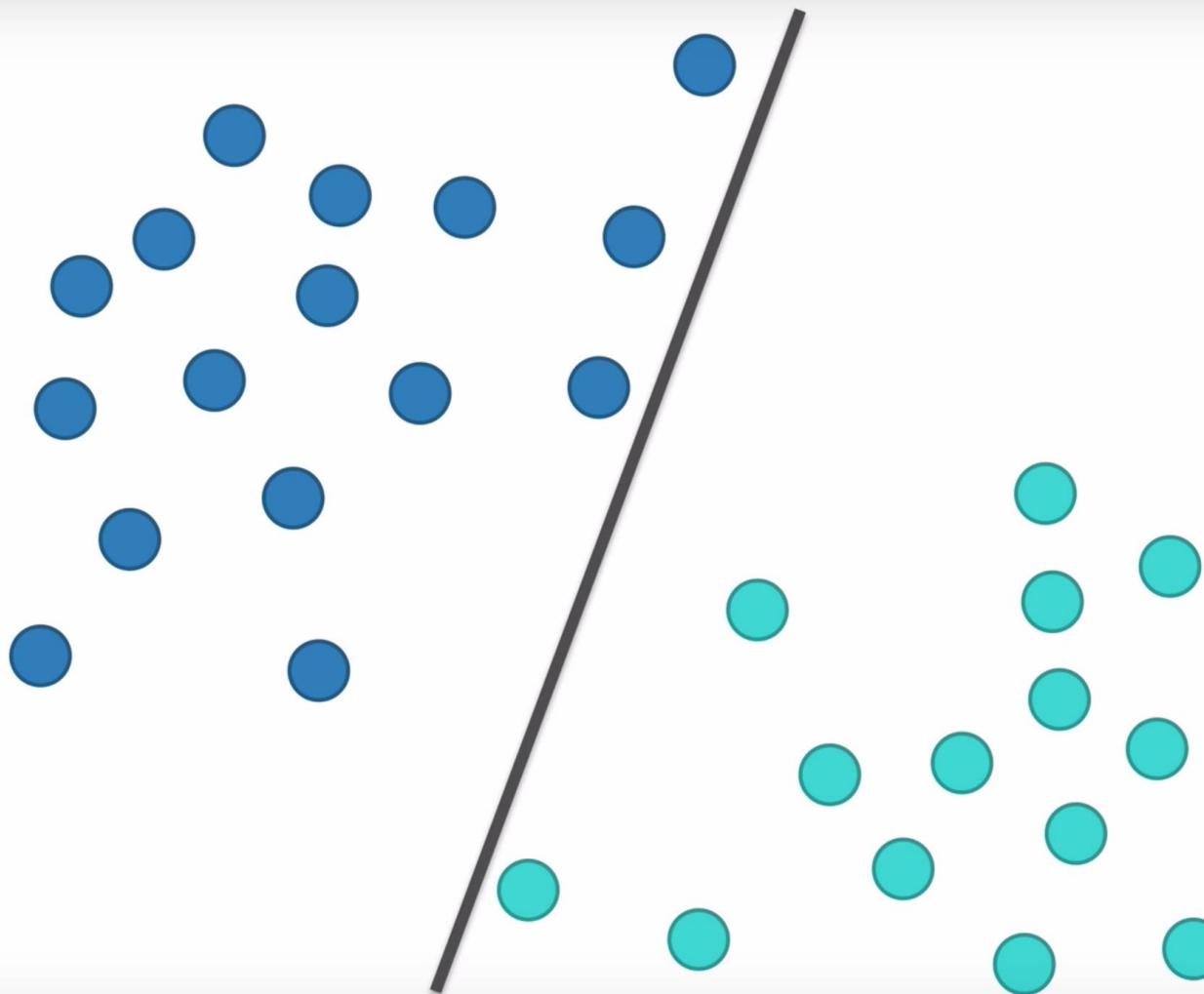
It's a classification technique.



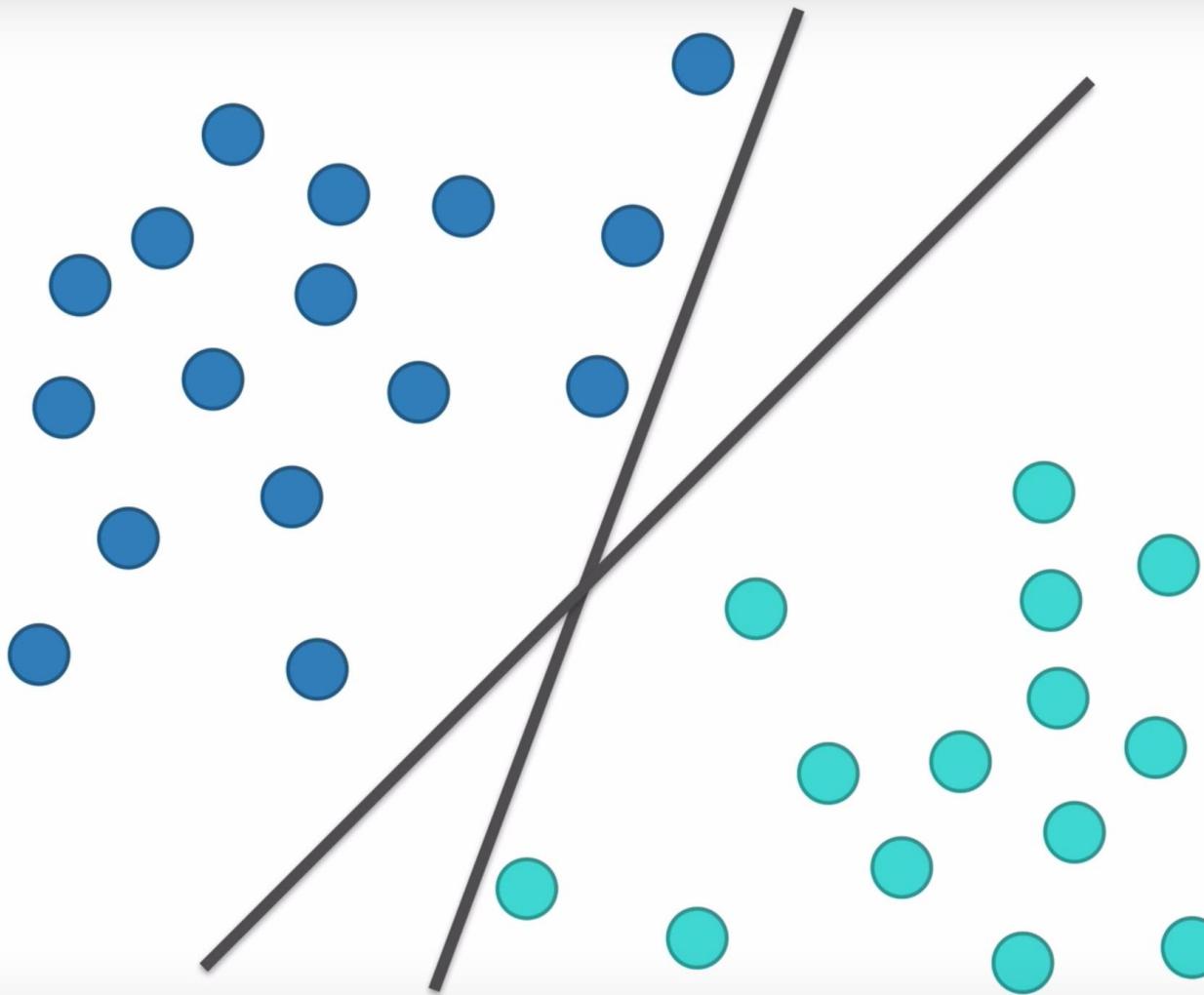


Split the data

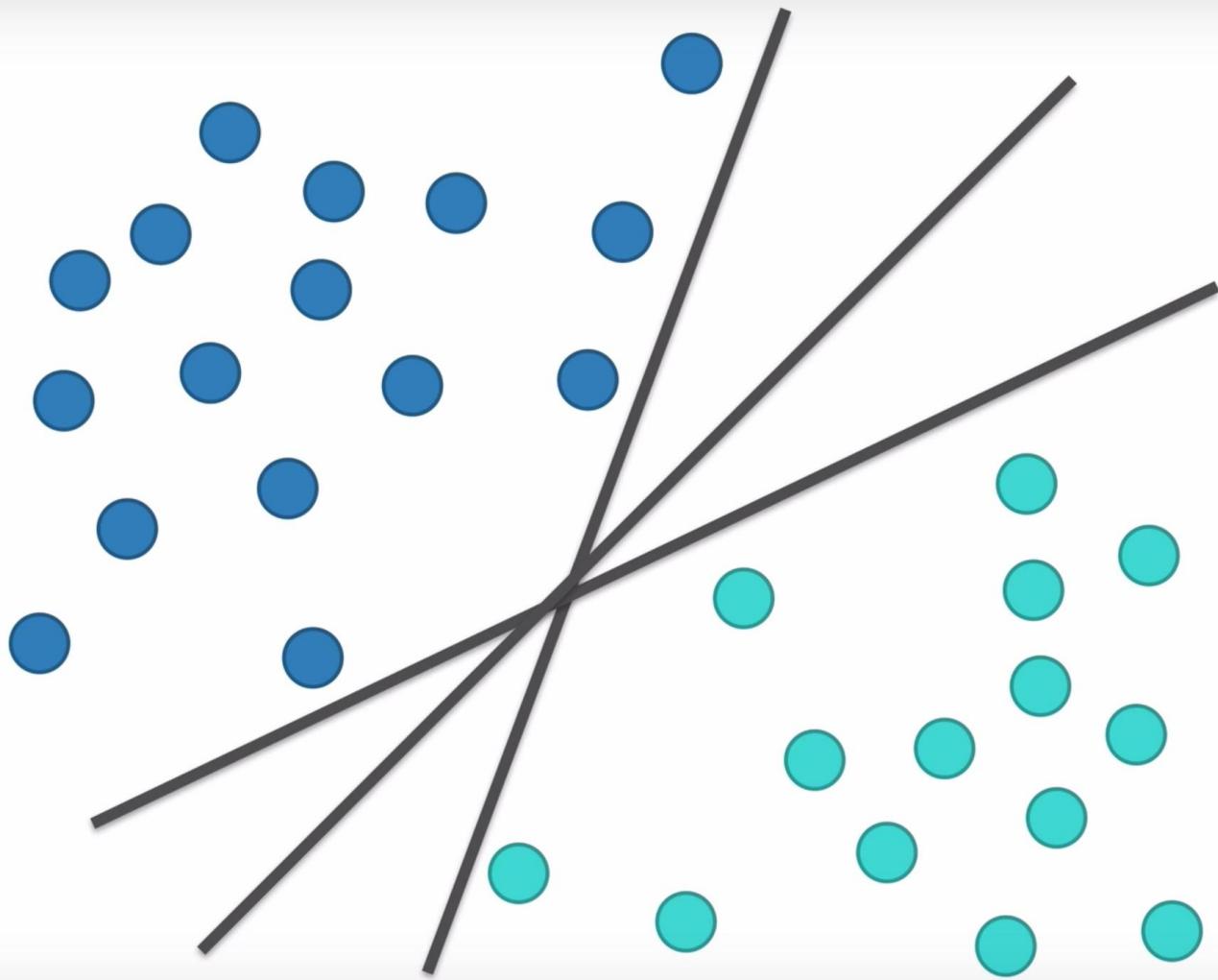
Split the data

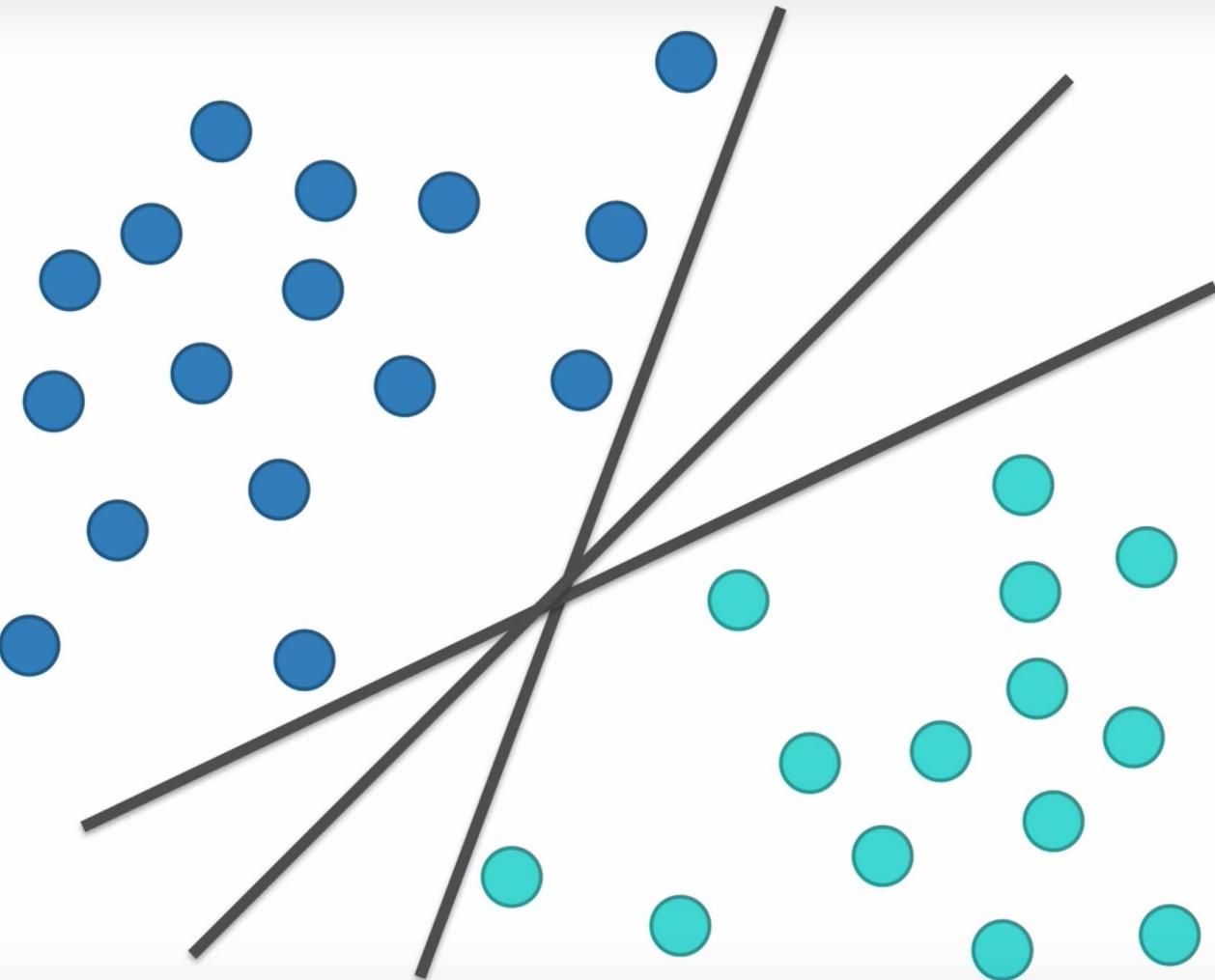


Split the data

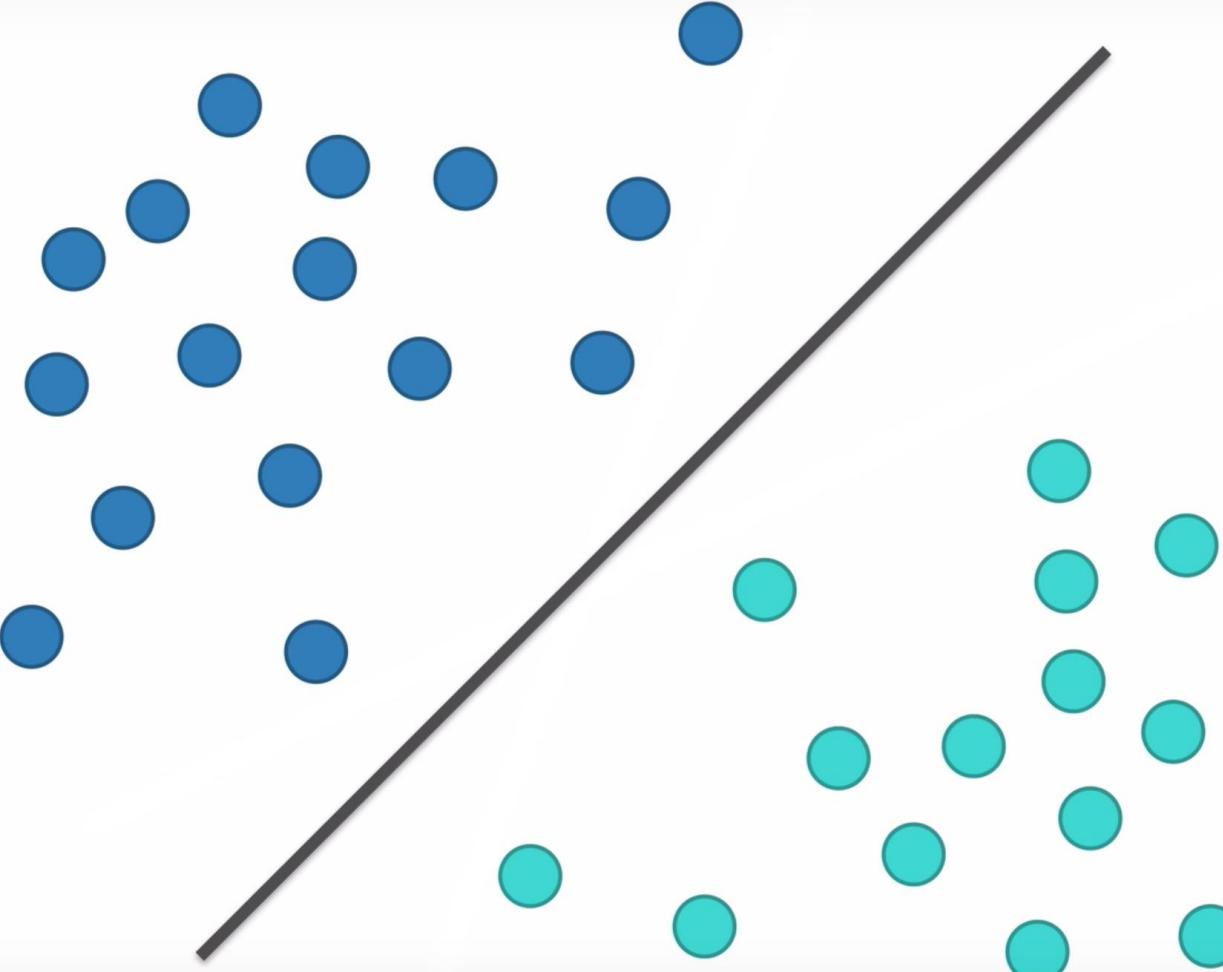


Split the data

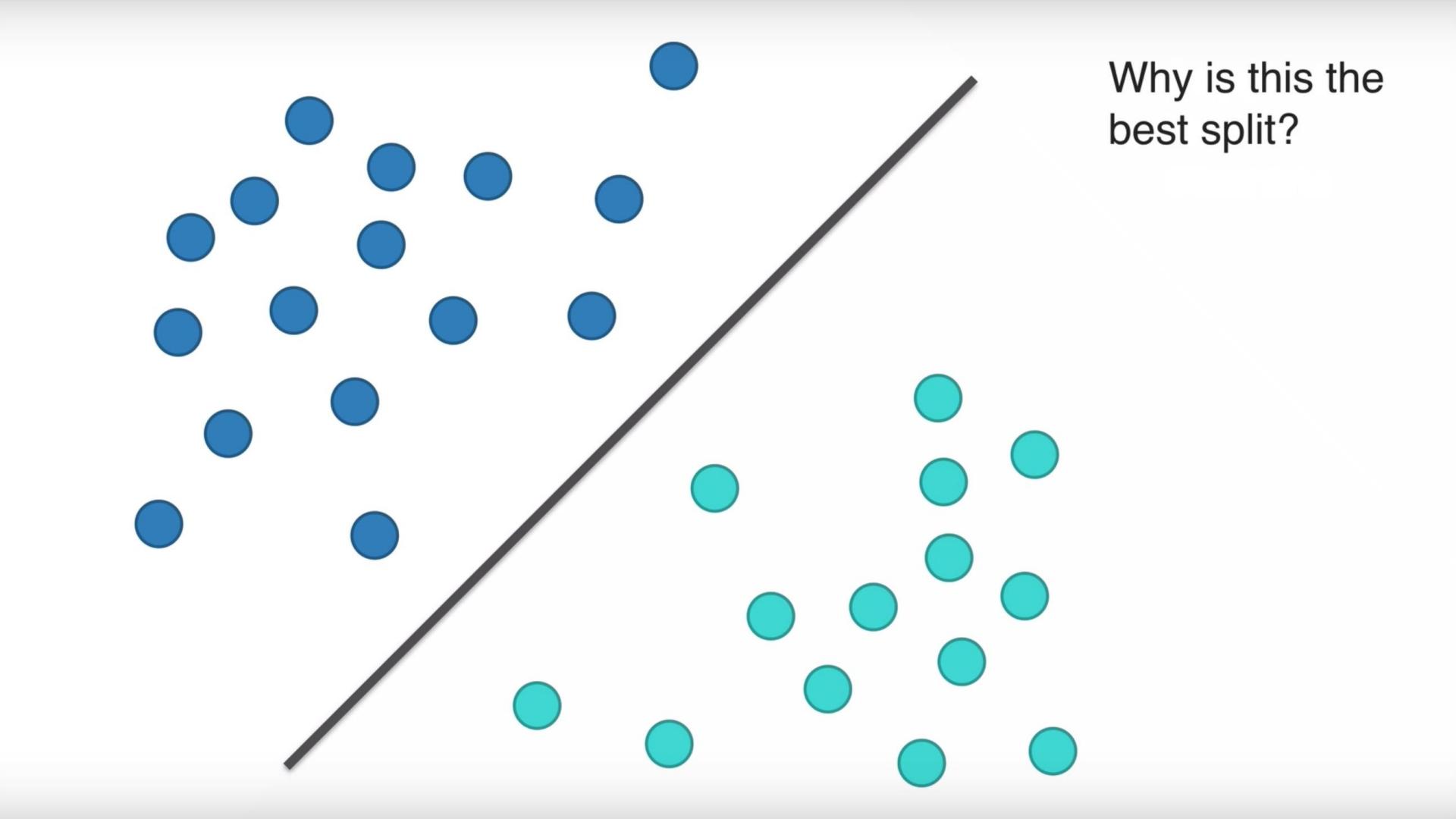




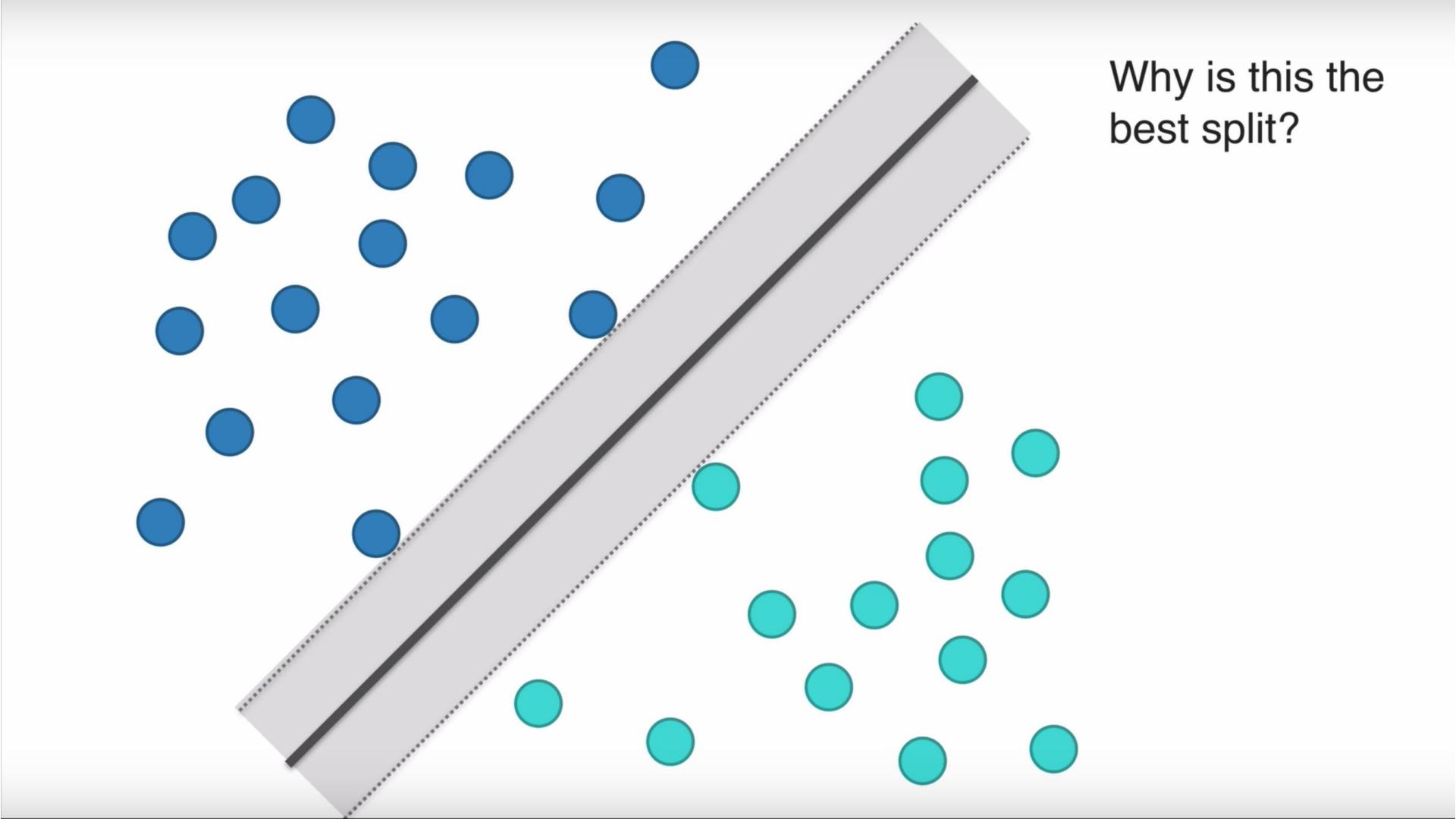
**Split the data
in the best
possible way**



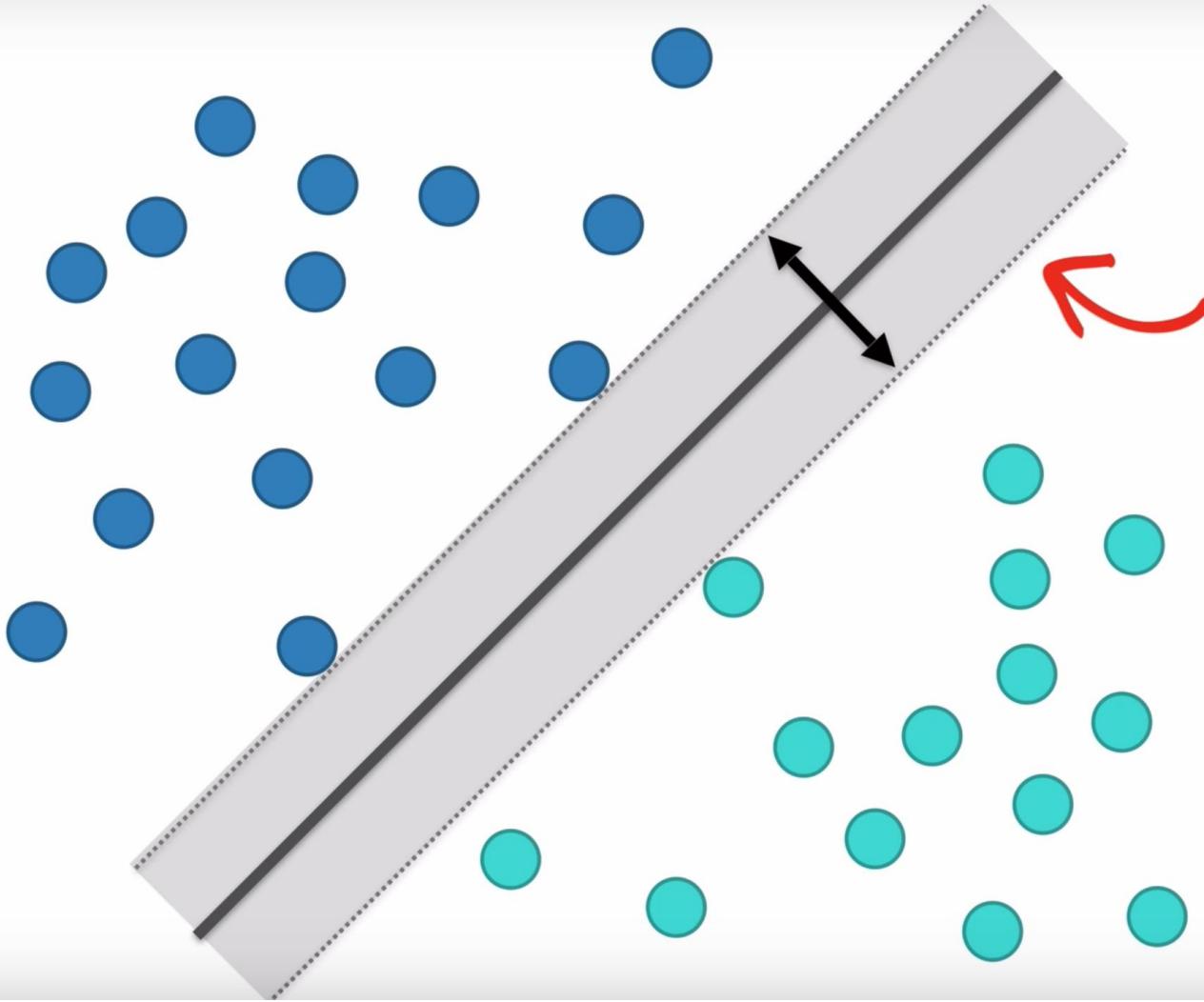
**Split the data
in the best
possible way**



Why is this the
best split?



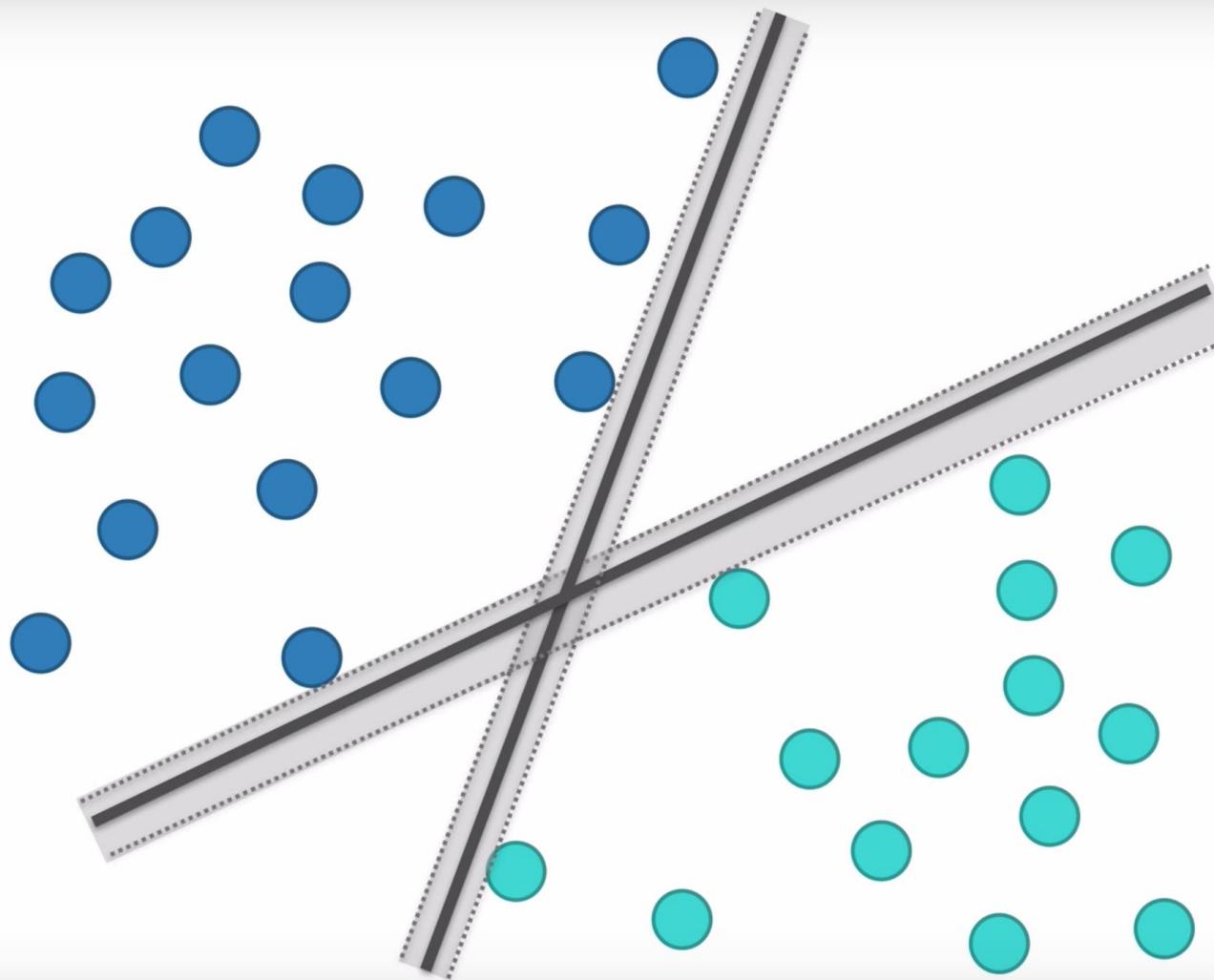
Why is this the
best split?

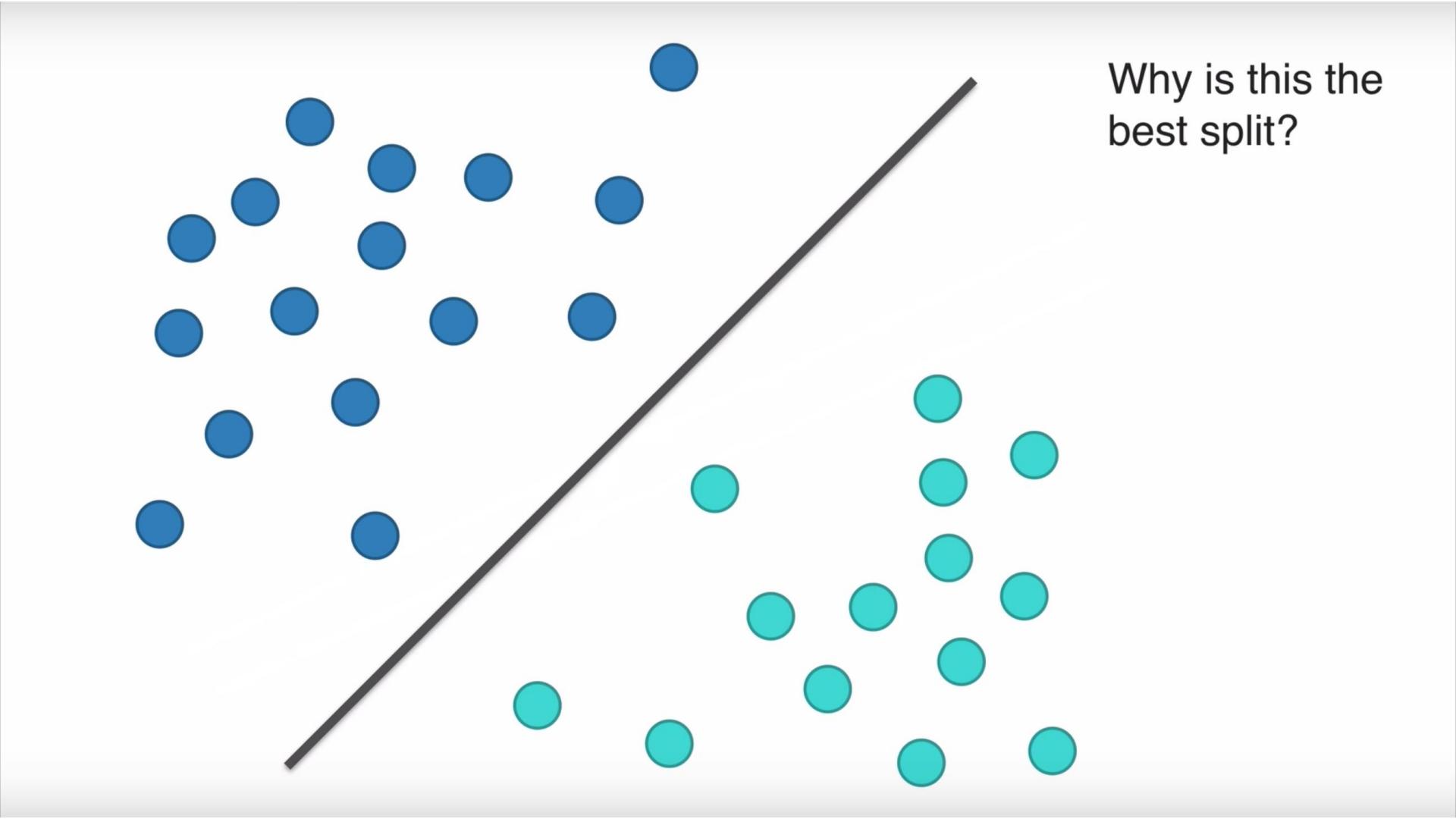


Why is this the
best split?

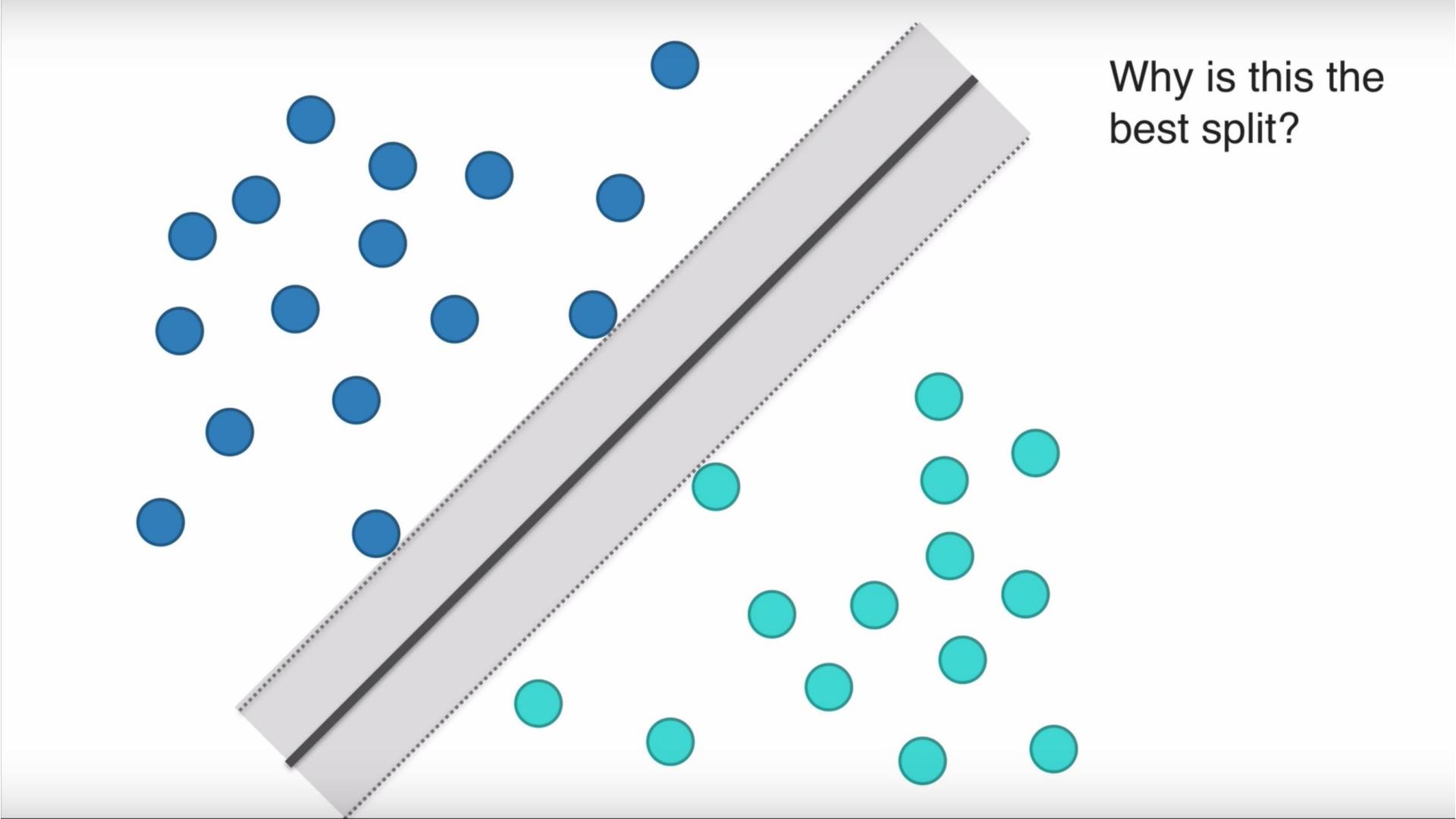
This is the
widest road that
separates the
two groups

Why is this the
best split?

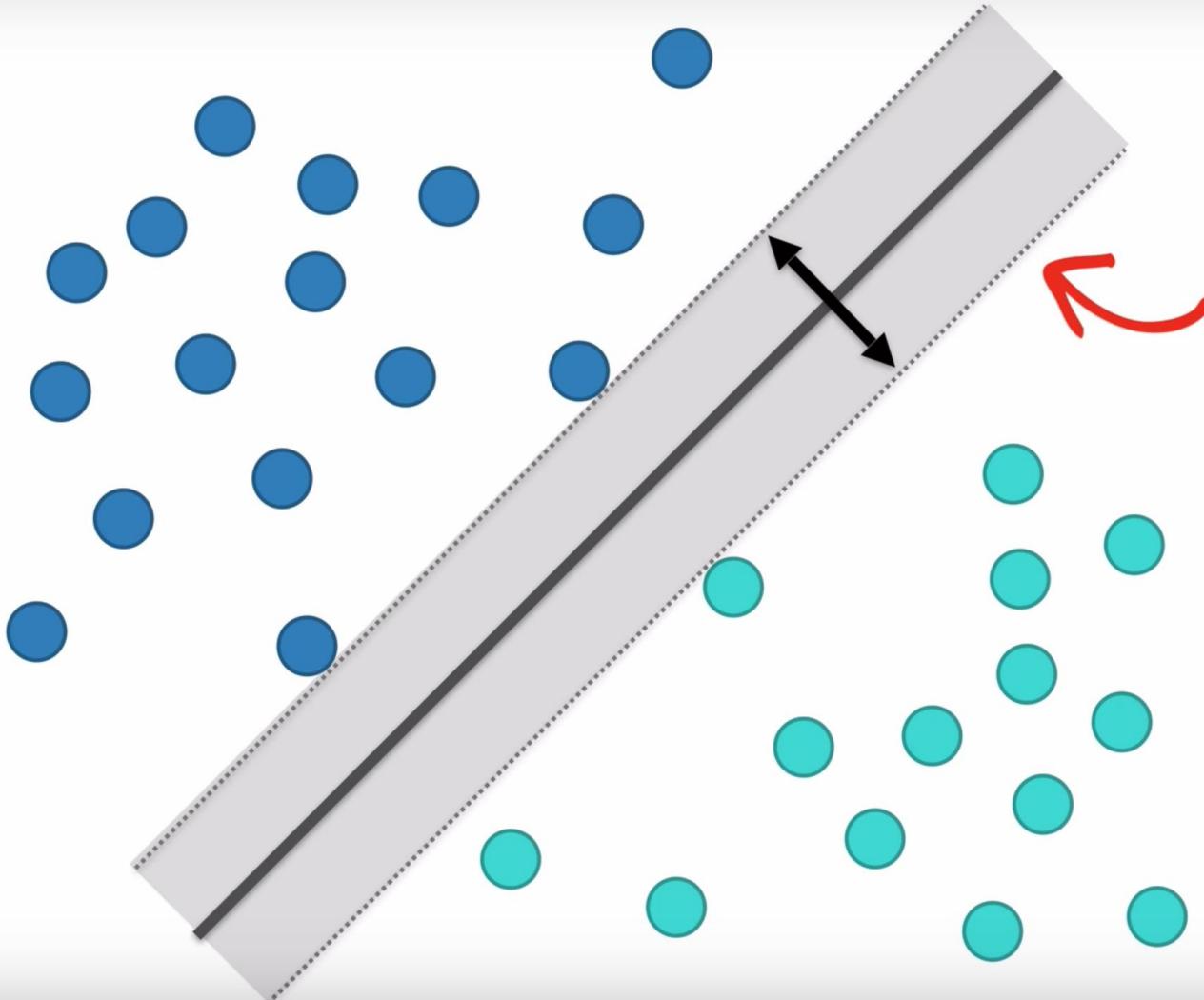




Why is this the
best split?

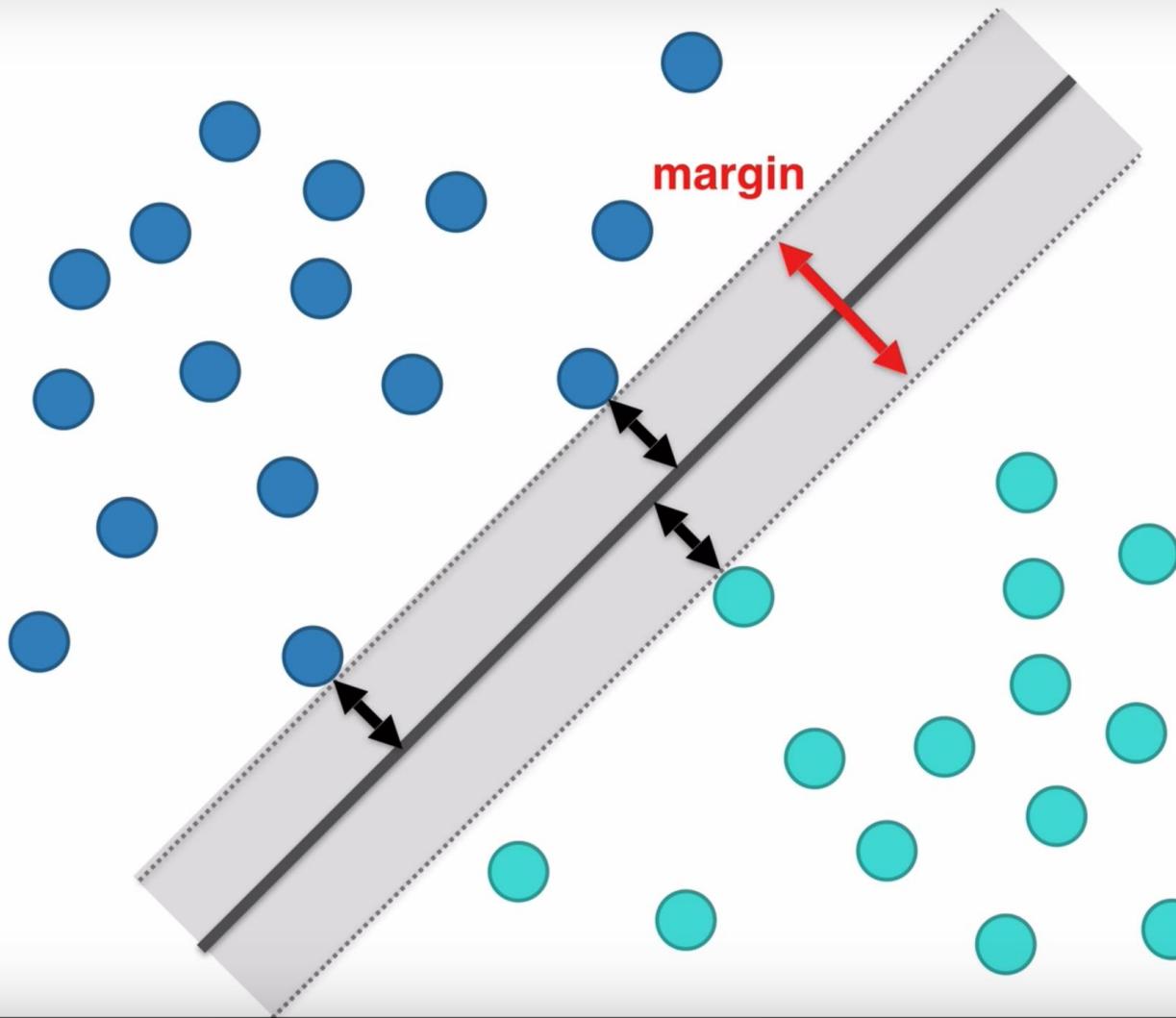


Why is this the
best split?



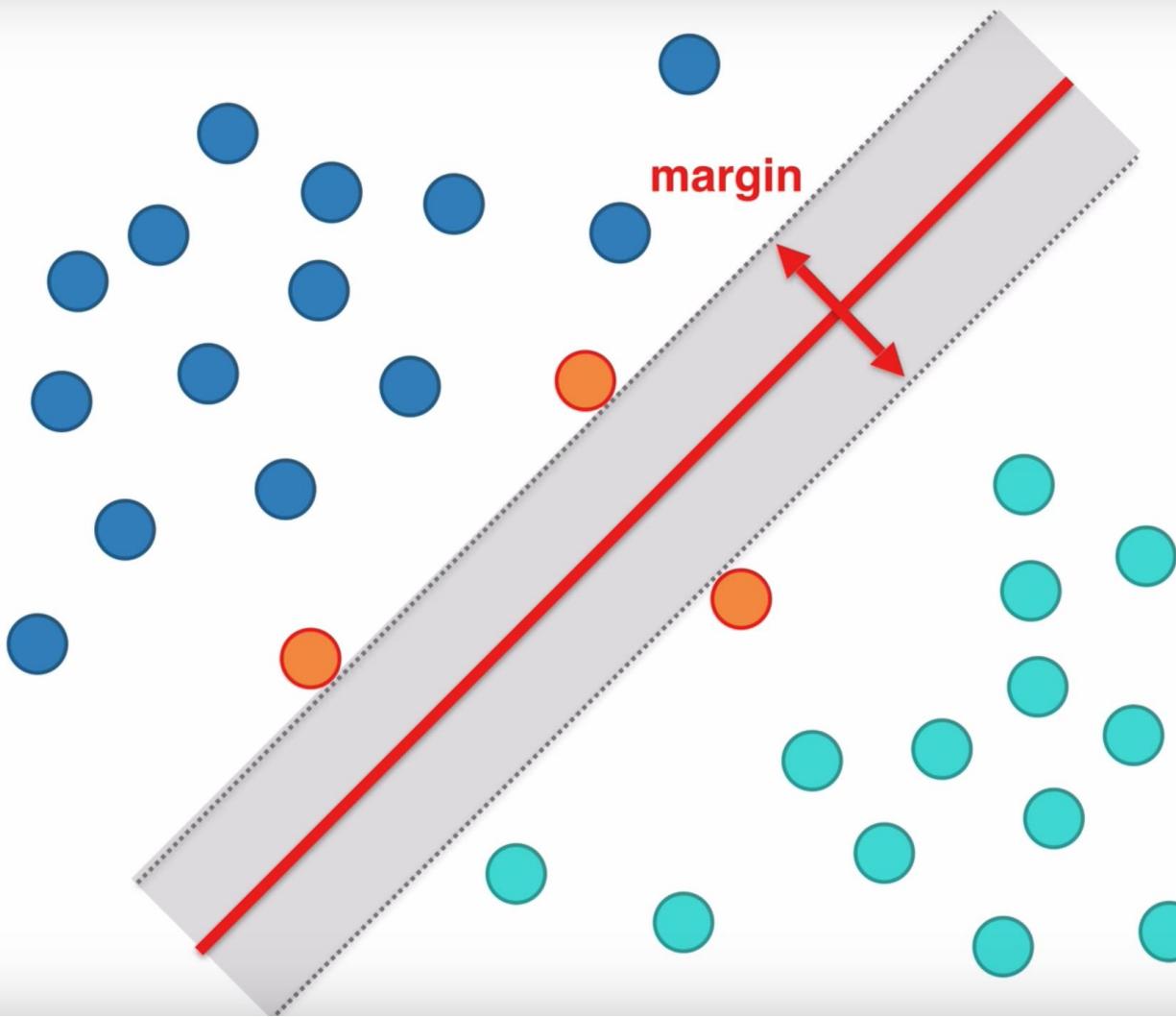
Why is this the
best split?

This is the
widest road that
separates the
two groups



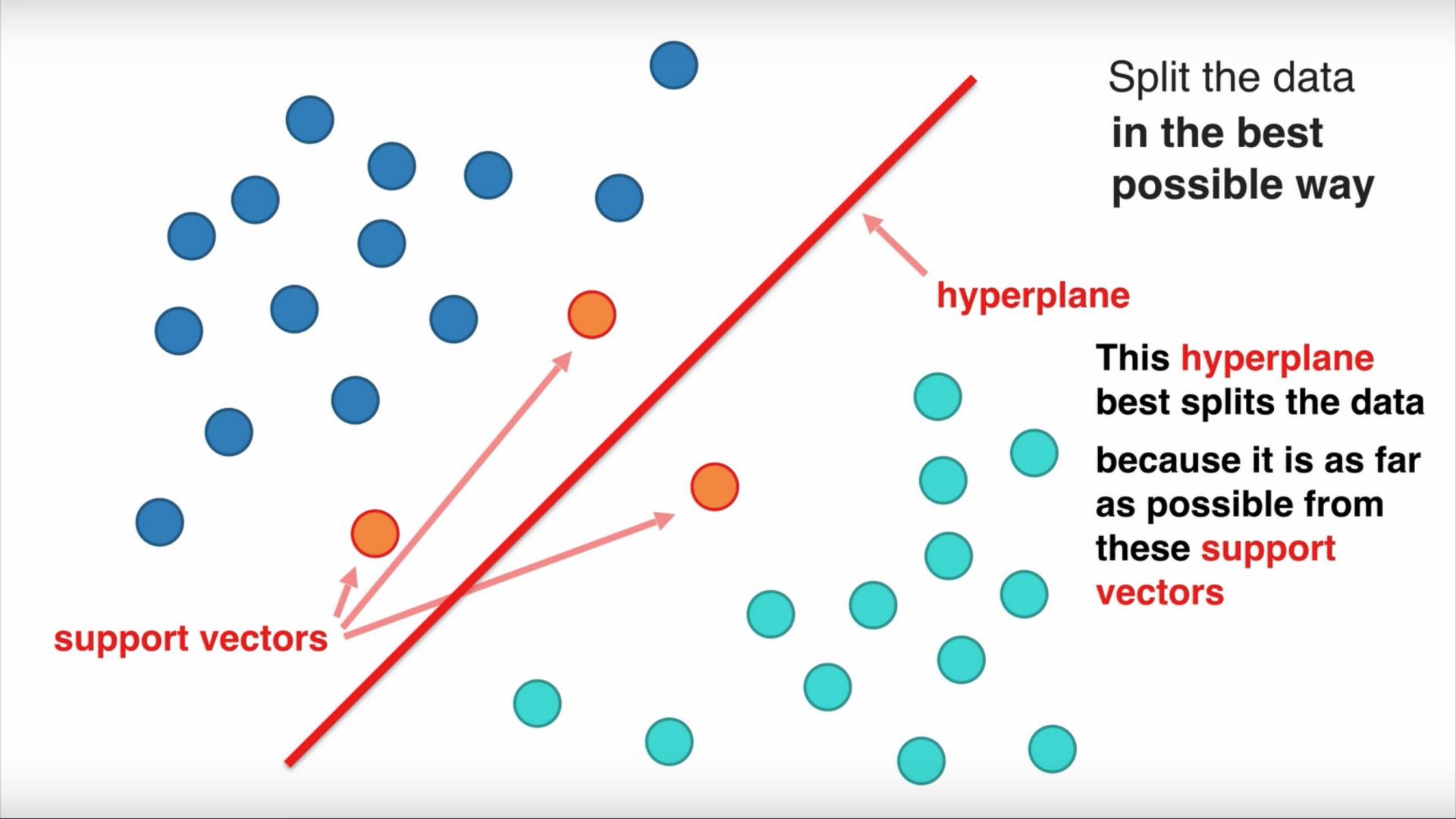
Why is this the best split?

The distance between the points and the line are as far as possible



Why is this the best split?

The distance between the support vectors and the hyperplane are as far as possible

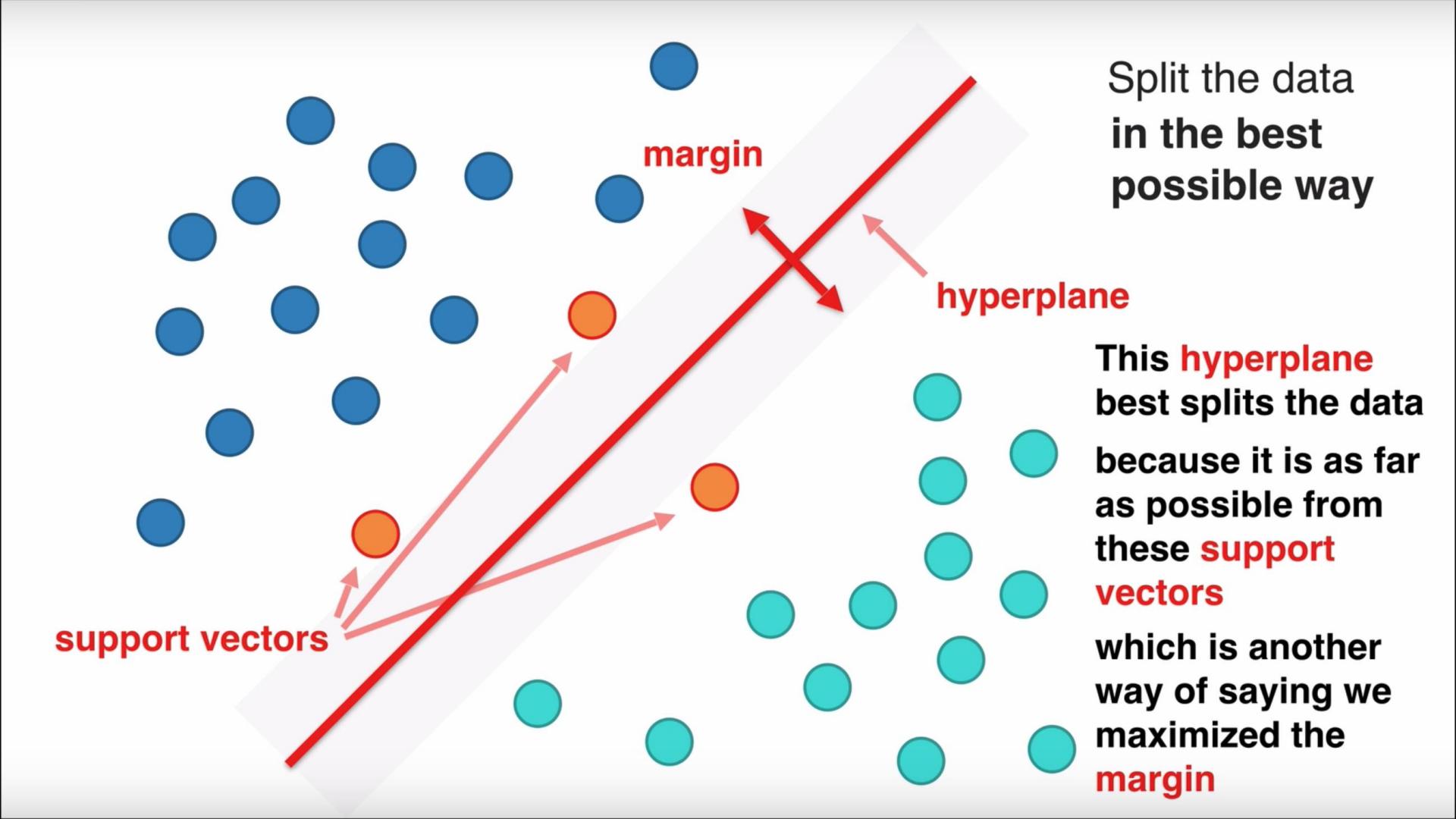


Split the data
in the best
possible way

hyperplane

support vectors

This **hyperplane**
best splits the data
because it is as far
as possible from
these **support
vectors**



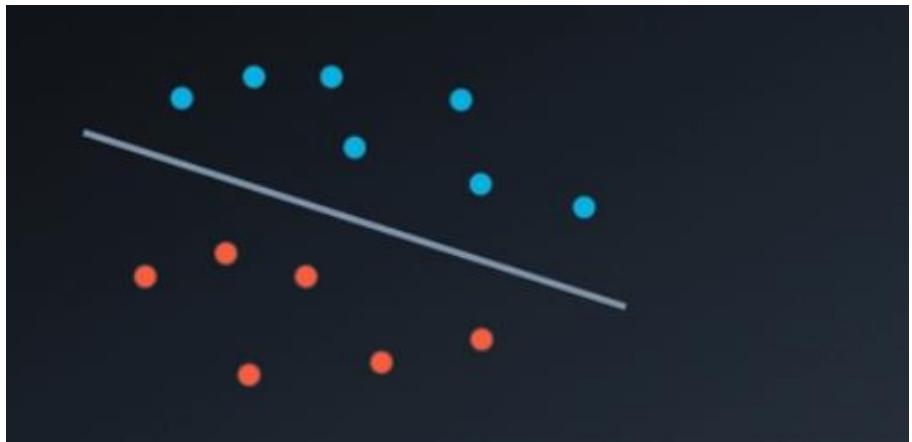
Split the data
in the best
possible way

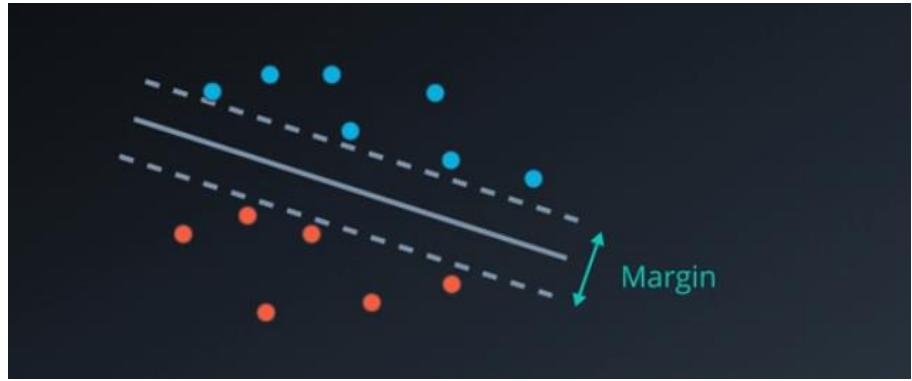
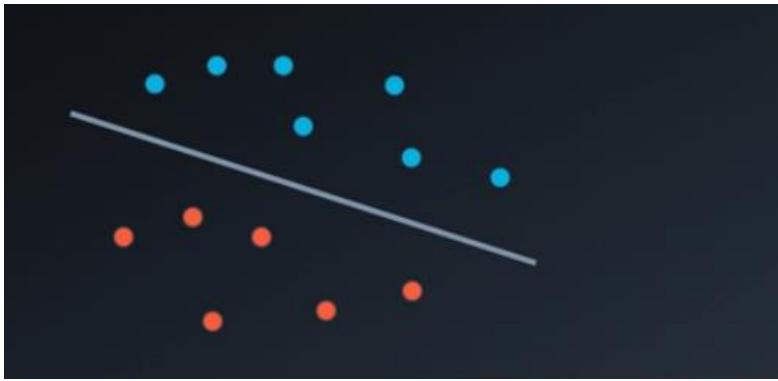
hyperplane

support vectors

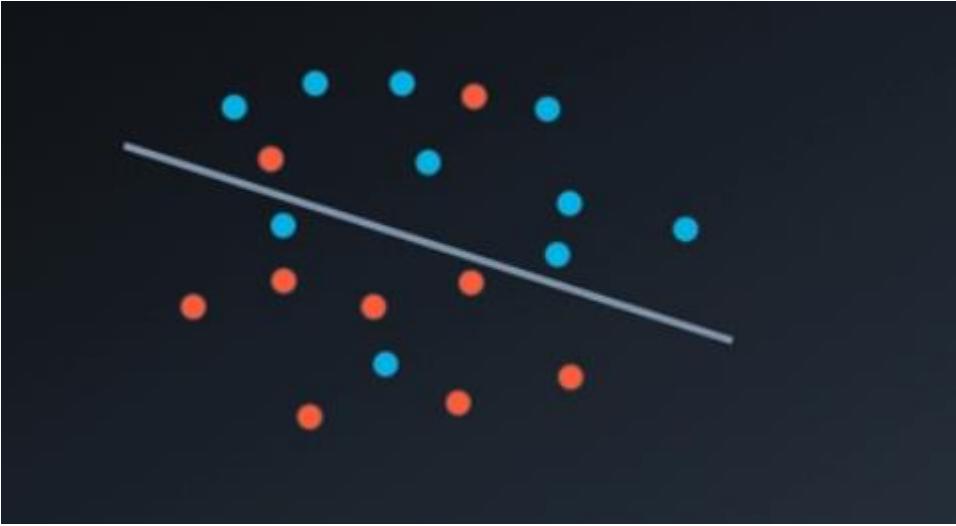
margin

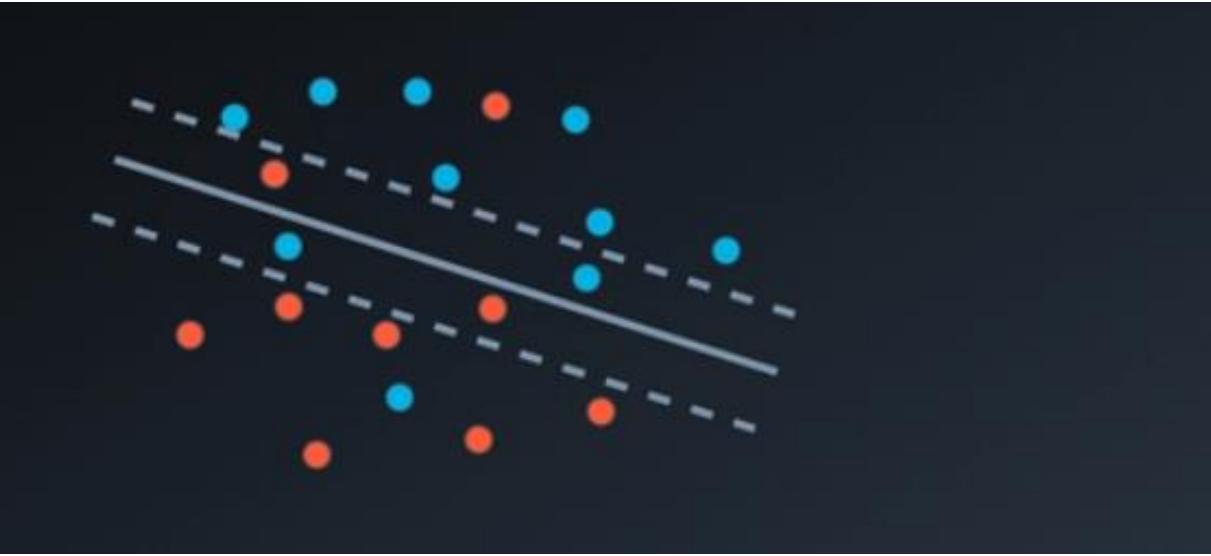
This **hyperplane**
best splits the data
because it is as far
as possible from
these **support
vectors**
which is another
way of saying we
maximized the
margin

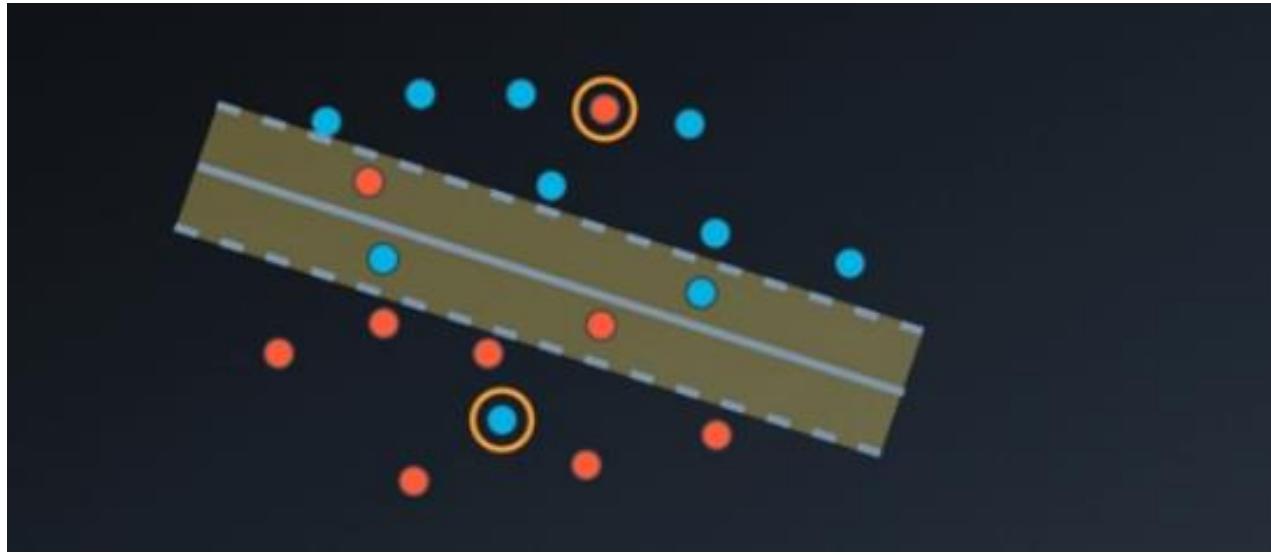


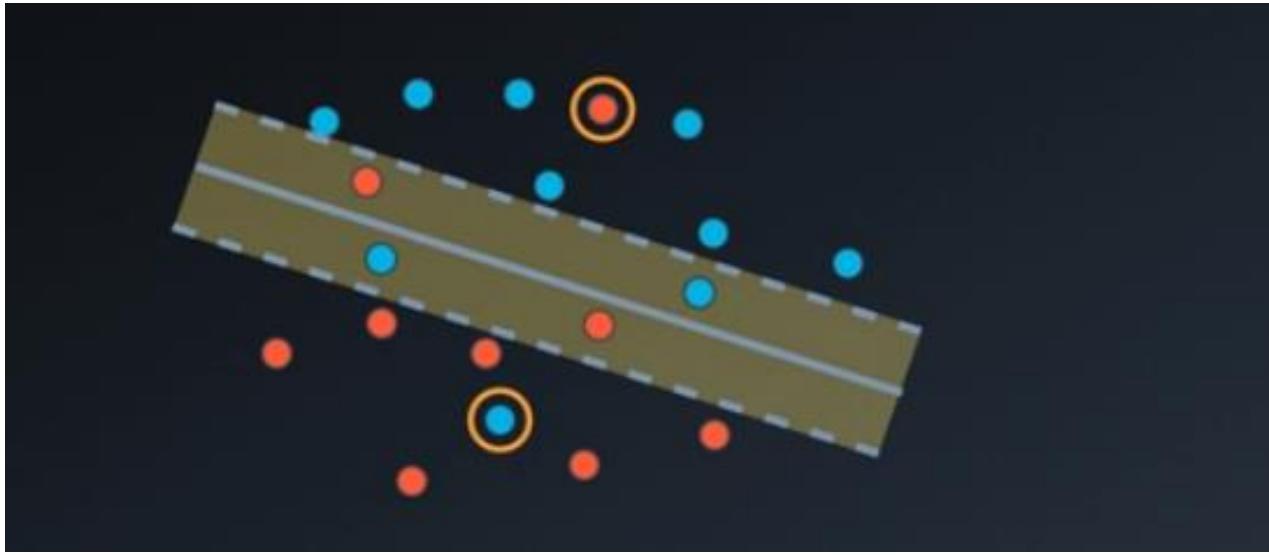










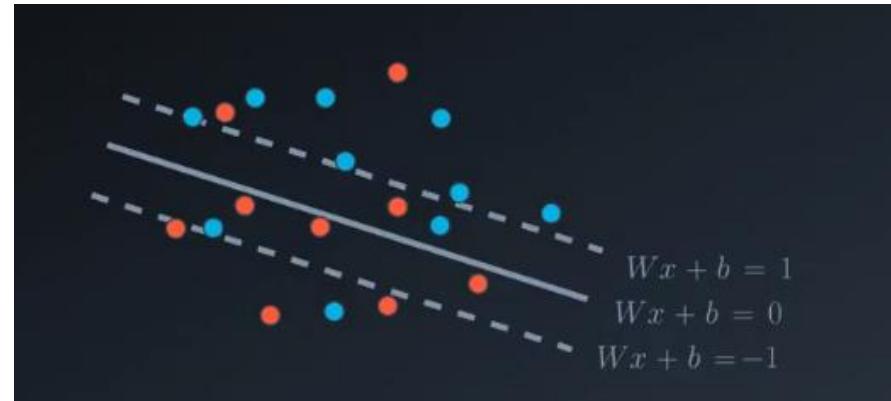


ERROR = CLASSIFICATION ERROR + MARGIN ERROR

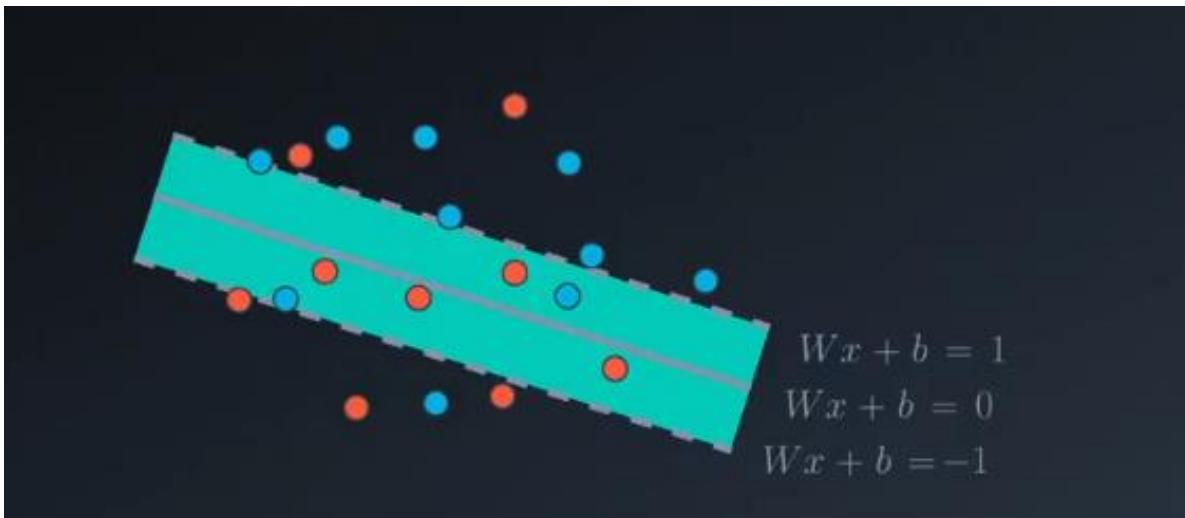
Classification Error

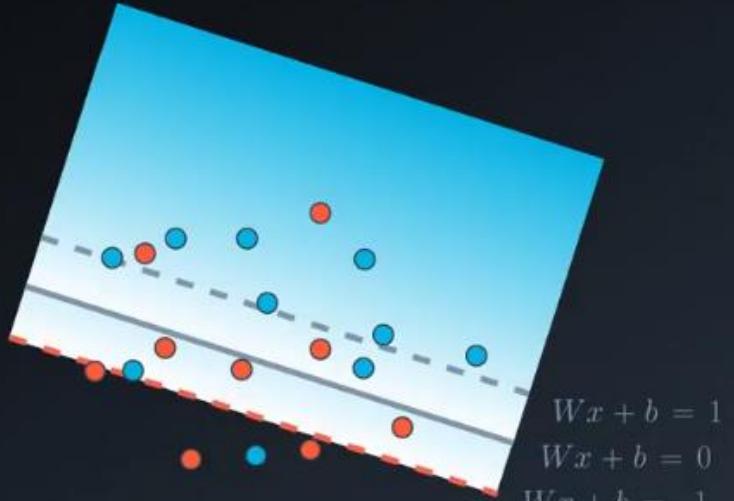


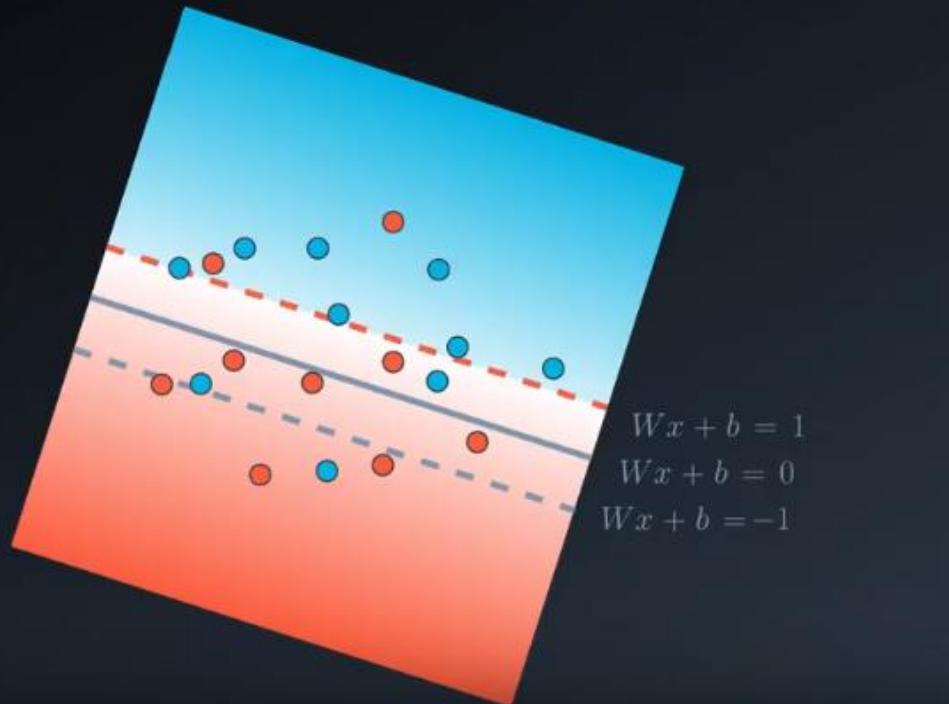
Classification Error

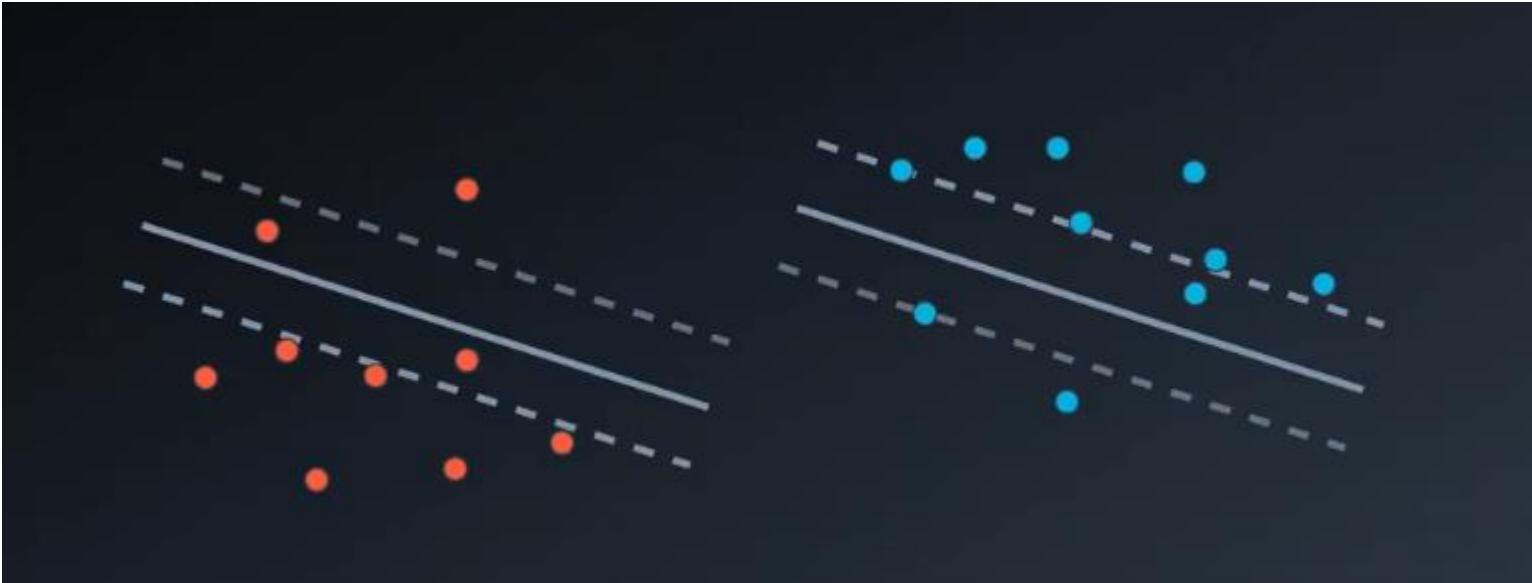


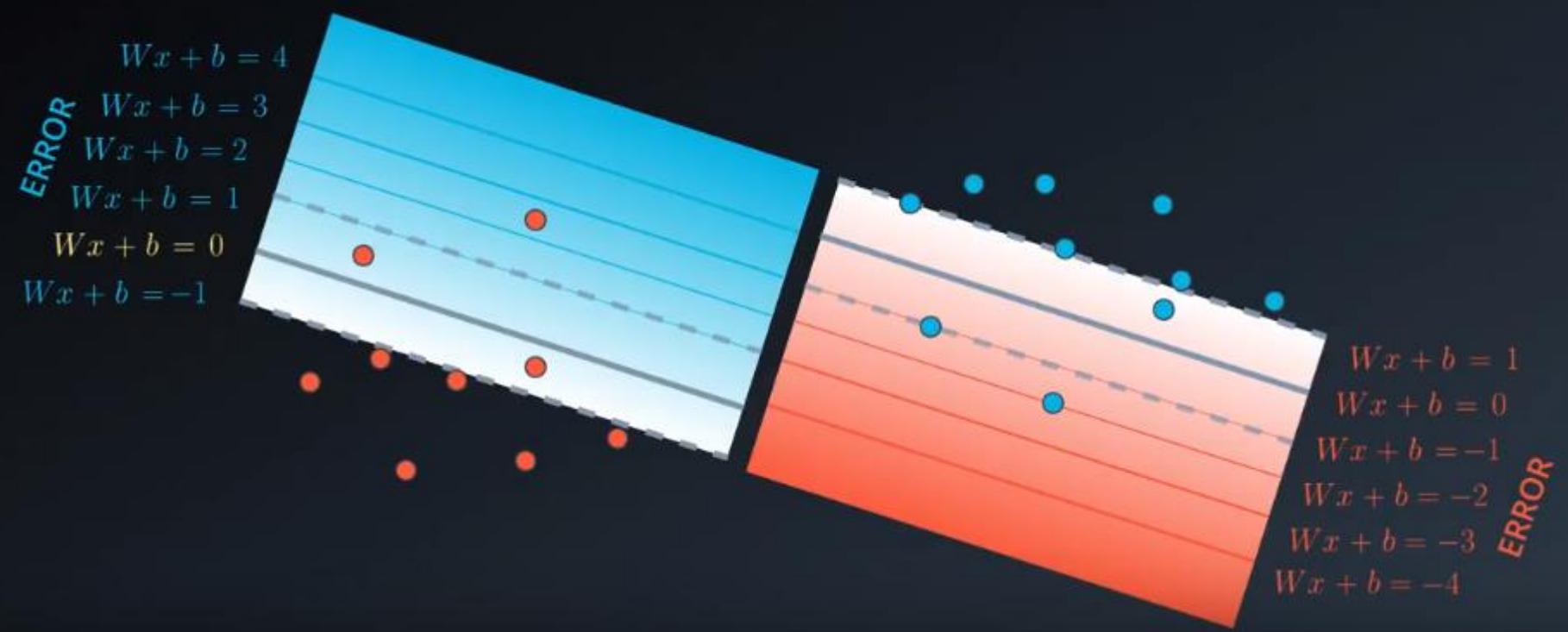
Classification Error

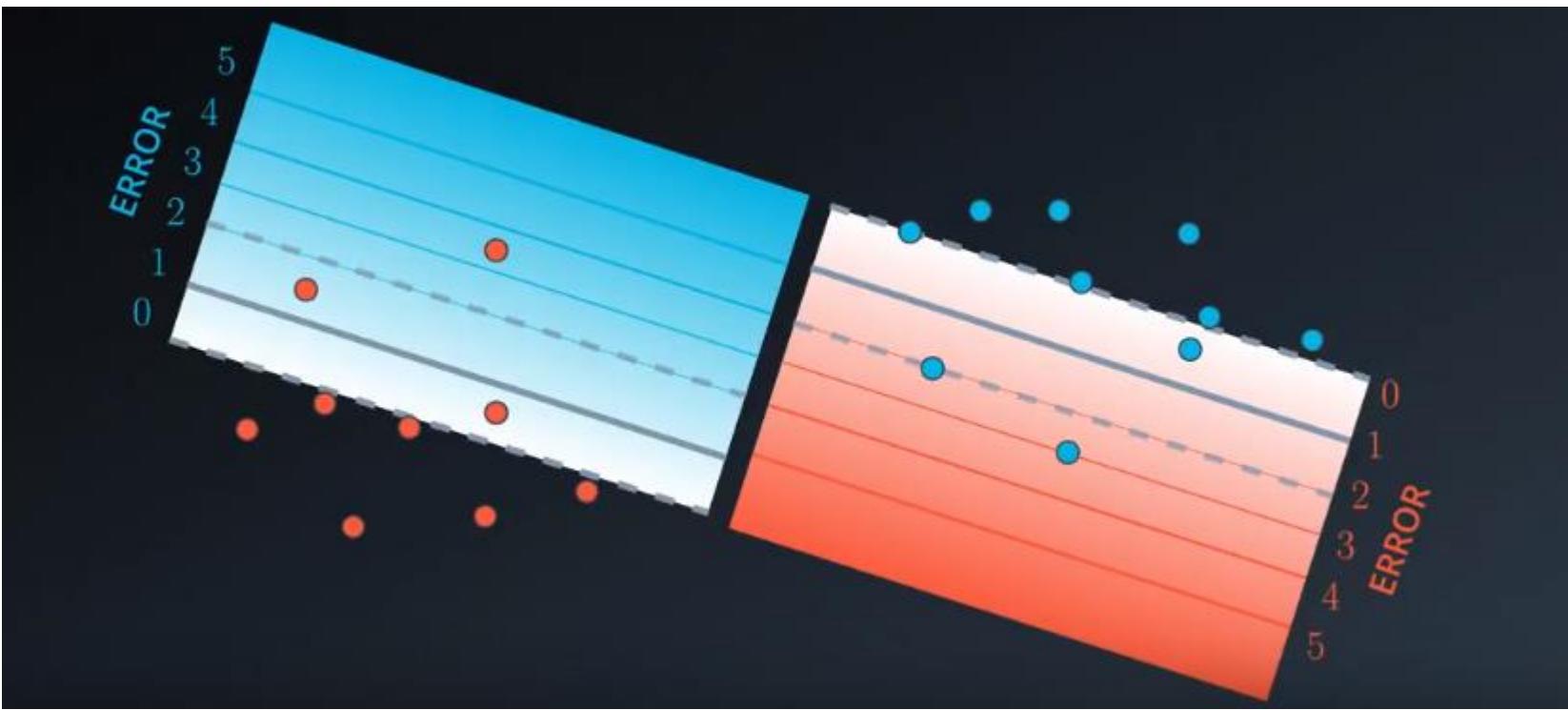






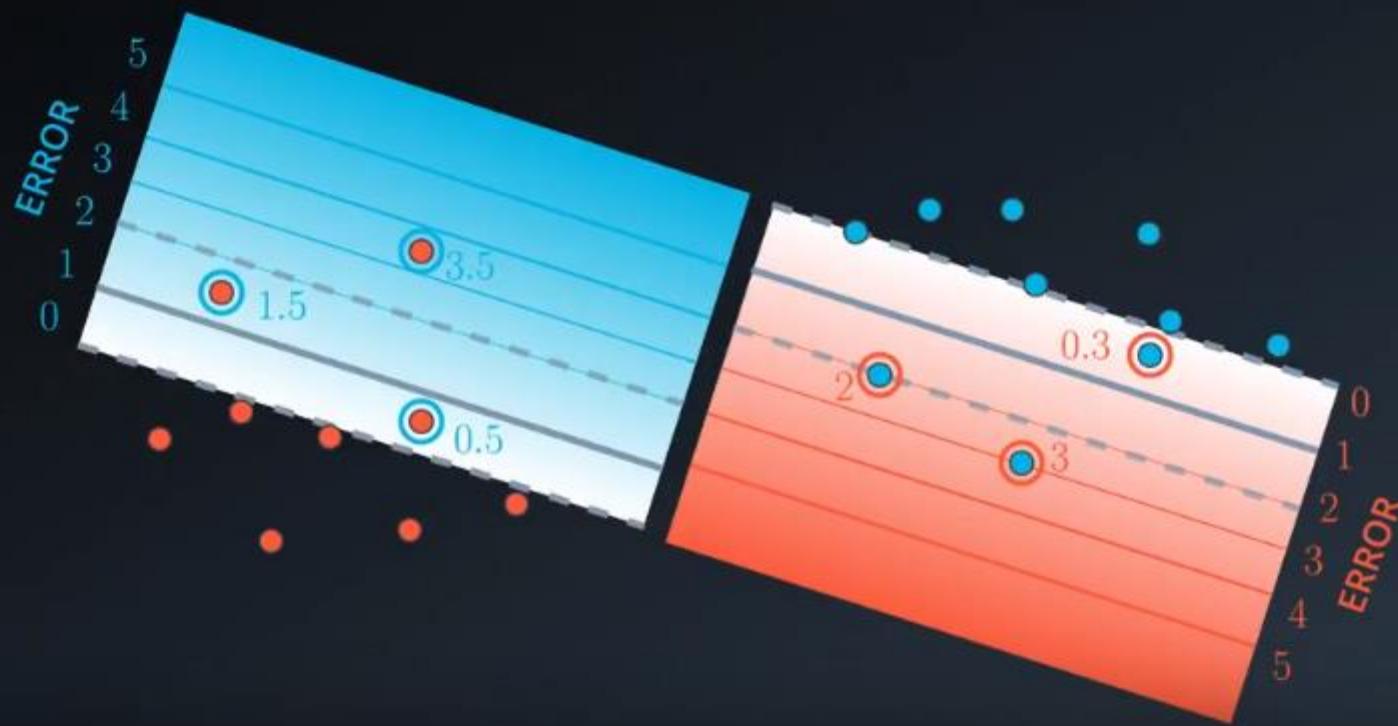




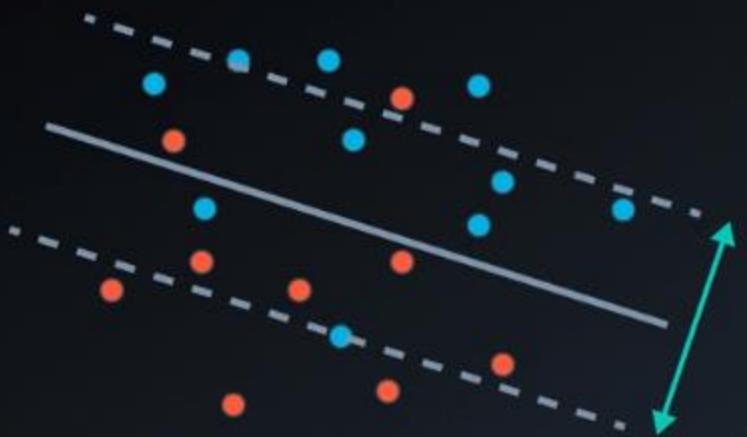


Classification Error

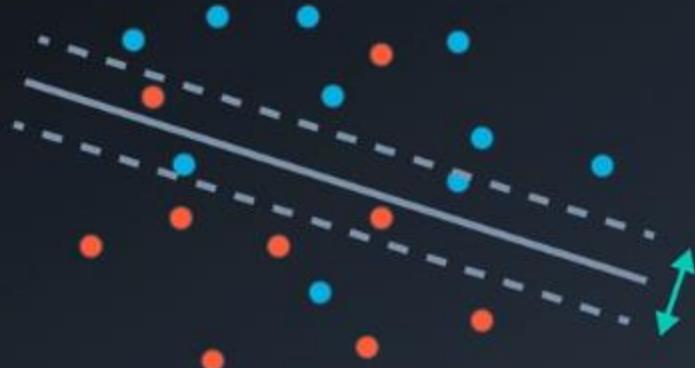
$$\text{Error} = 1.5 + 3.5 + 0.5 + 2 + 3 + 0.3 = 10.8$$



Margin Error



LARGE MARGIN
SMALL ERROR



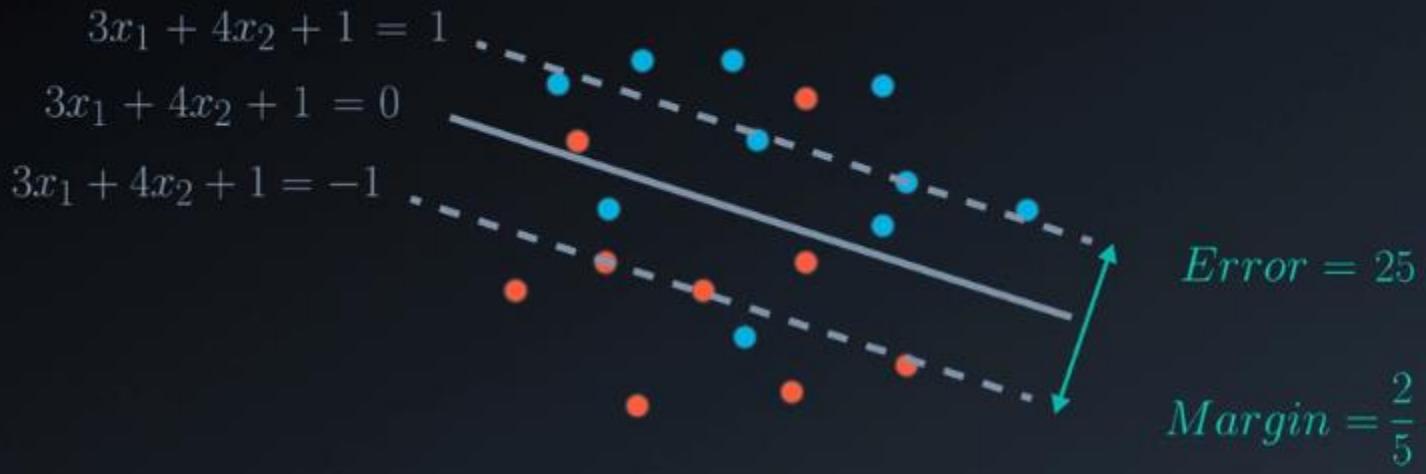
SMALL MARGIN
LARGE ERROR

$$\text{Error} = |W|^2$$

$$\text{Margin} = \frac{2}{|W|}$$

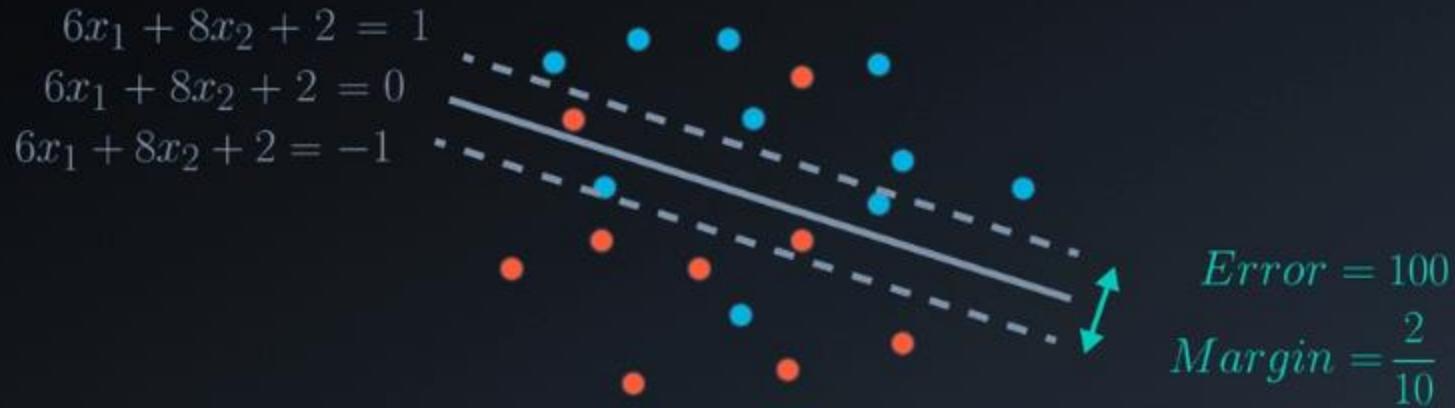
LARGE MARGIN \longrightarrow SMALL ERROR

SMALL MARGIN \longrightarrow LARGE ERROR



$$W = (3, 4) \quad w_1 x_1 + w_2 x_2 + b = 0 \quad |W|^2 = 3^2 + 4^2 = 25$$

$$b = 1 \quad 3x_1 + 4x_2 + 1 = 0 \quad \frac{2}{|W|} = \frac{2}{5}$$



$$W = (6, 8)$$

$$w_1x_1 + w_2x_2 + b = 0$$

$$|W|^2 = 6^2 + 8^2 = 100$$

$$b = 2$$

$$6x_1 + 8x_2 + 2 = 0$$

$$\frac{2}{|W|} = \frac{2}{10}$$

The C Parameter



Error = C Classification Error + Margin Error

Cheat Code

Small C

Large margin

May make classification errors

Large C

Classifies points well

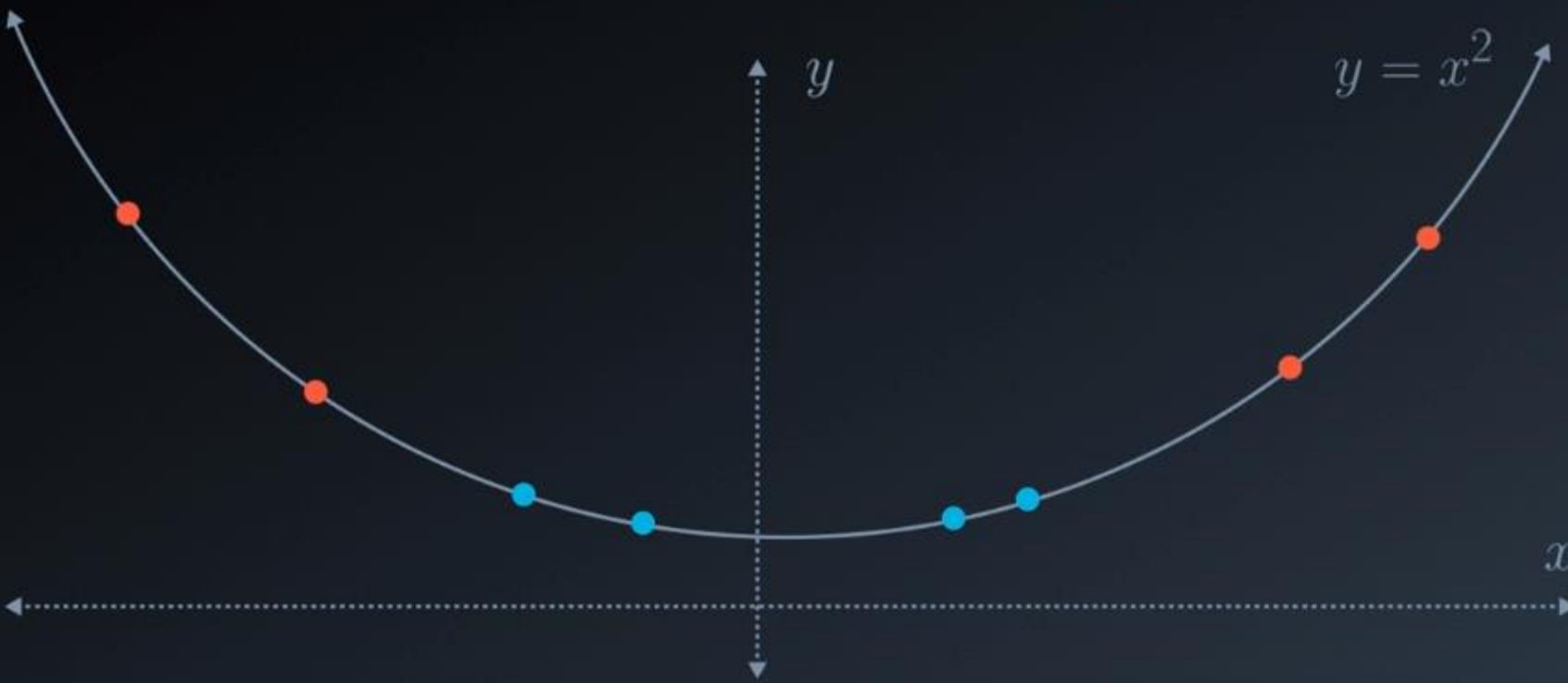
May have a small margin

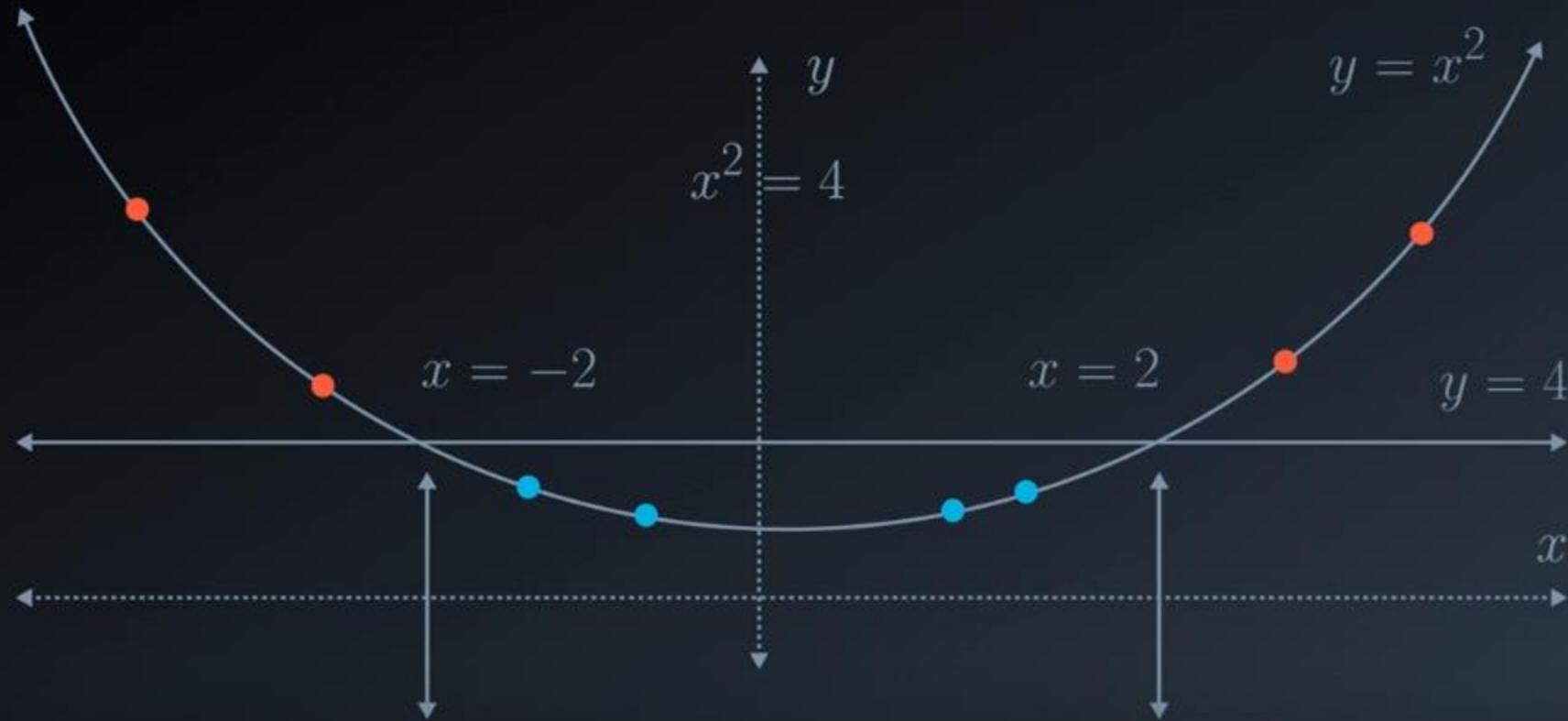
Kernel Trick



Kernel Trick





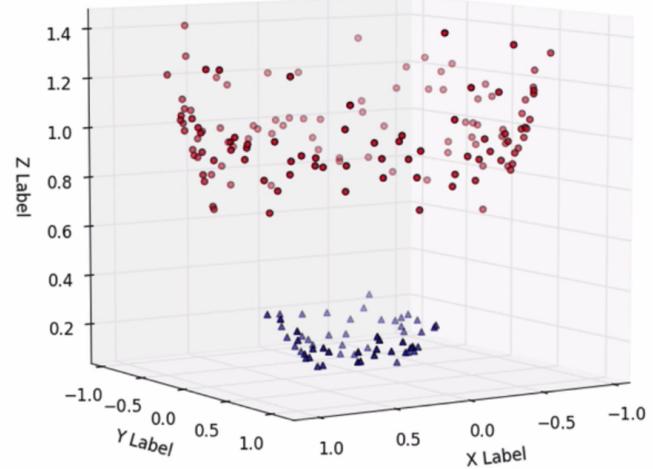
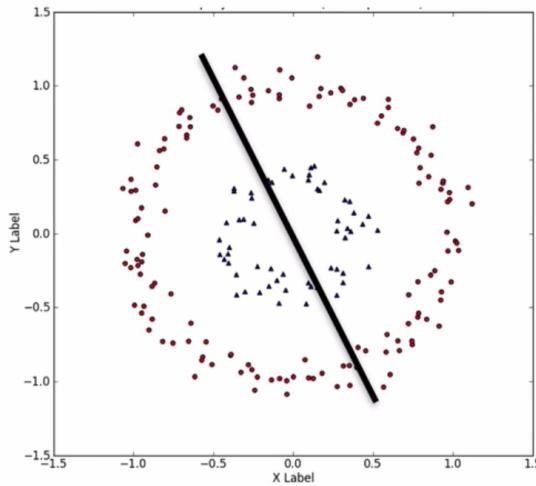


$$x = -2$$



$$x = 2$$

Kernel Trick: Visual



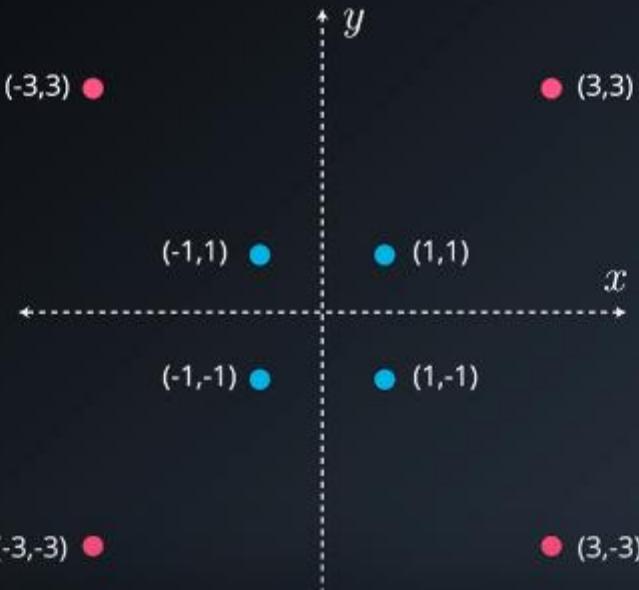
Higher Dimensions

C Parameter

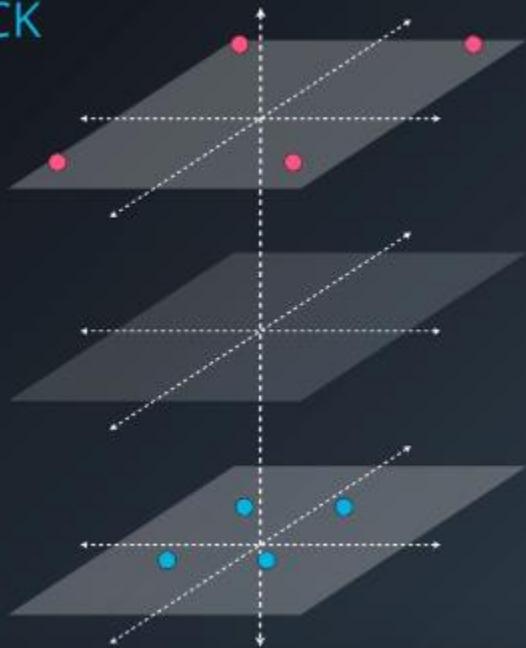
Multiple Classes

Kernel Trick

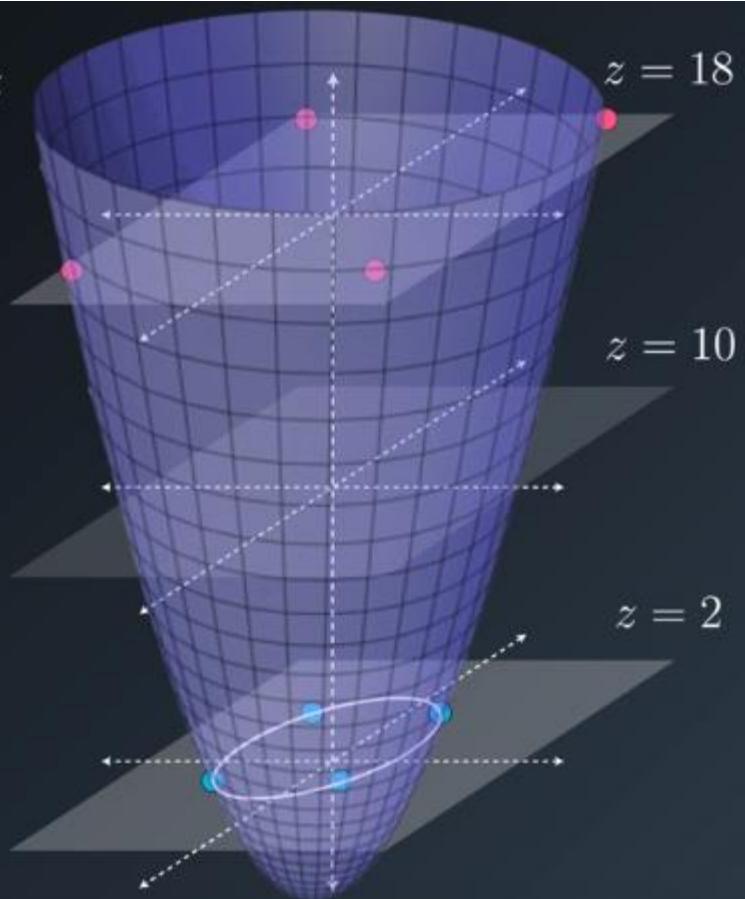
KERNEL TRICK



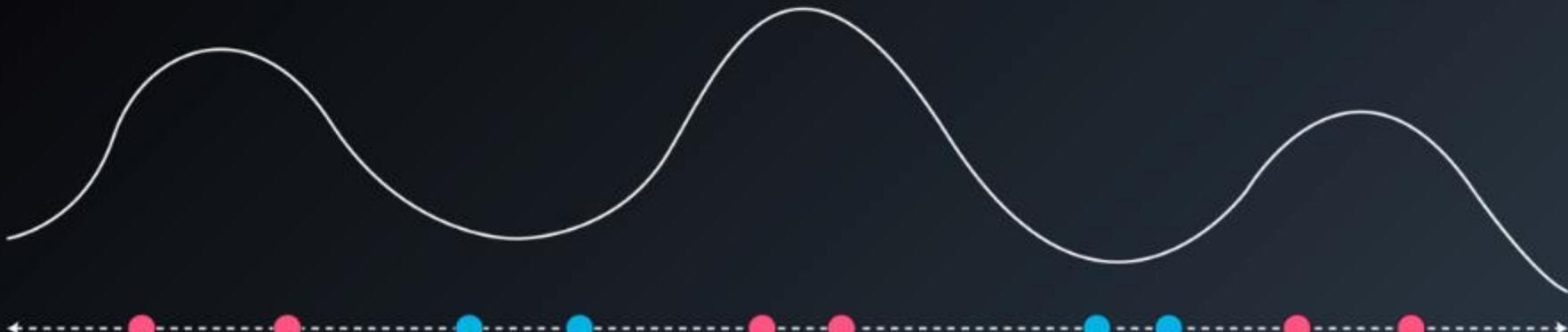
KERNEL TRICK



$$x^2 + y^2 = z$$



RBF KERNEL



RBF KERNEL



...

So now what?

**Let's apply this to a real world
problem.**

Cupcakes



Muffins



versus

“a cupcake is just a muffin with frosting”

“a muffin is just a cupcake with random bits of stuff in it”

Cupcakes vs Muffins

The Challenge

Classify recipes as cupcakes or muffins. When given a new recipe, determine if it's a cupcake or a muffin.

The Steps

1. Find the data
2. Apply a data science model
3. Review the results

Cupcakes vs Muffins

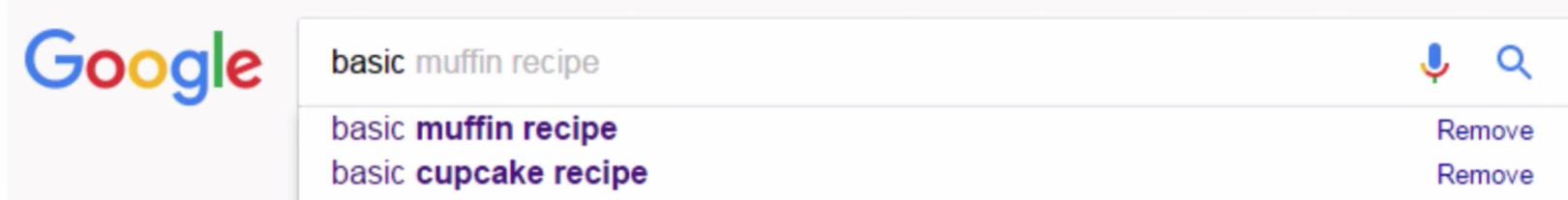
The Challenge

Classify recipes as cupcakes or muffins. When given a new recipe, determine if it's a cupcake or a muffin.

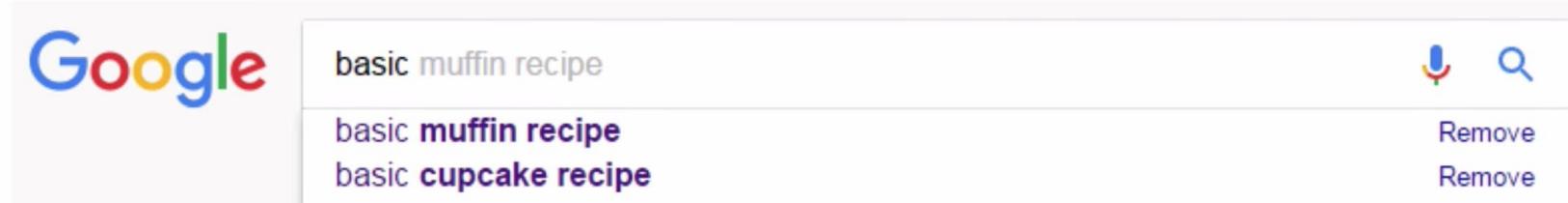
The Steps

1. Find the data
2. Apply a data science model
3. Review the results

1. Find the data



1. Find the data



Recorded the top 10 muffin and top 10 cupcake recipes

1. Find the data

Problem

Each recipe yields different amounts of batter

1. Find the data

Problem

Each recipe yields different amounts of batter

Solution

Amount-based

Recipe	Flour	Sugar	Other
Muffin1	2 cups	1/2 cup	...
Cupcake1	2 cups	3/4 cup	...



Percent-based

Recipe	Flour	Sugar	Other	Total Volume
Muffin1	47%	24%	...	100%
Cupcake1	42%	21%	...	100%

1. Find the data

Type	Flour	Milk	Sugar	Butter	Egg	Baking Powder	Vanilla	Salt
Muffin	55	28	3	7	5	2	0	0
Muffin	47	24	12	6	9	1	0	0
Muffin	47	23	18	6	4	1	0	0
Muffin	50	25	12	6	5	2	1	0
Muffin	55	27	3	7	5	2	1	0
Muffin	54	27	7	5	5	2	0	0
Muffin	47	26	10	10	4	1	0	0
Muffin	50	17	17	8	6	1	0	0
Muffin	50	17	17	11	4	1	0	0
Cupcake	39	0	26	19	14	1	1	0
Cupcake	34	17	20	20	5	2	1	0
Cupcake	39	13	17	19	10	1	1	0
Cupcake	38	15	23	15	8	0	1	0
Cupcake	42	18	25	9	5	1	0	0
Cupcake	36	14	21	14	11	2	1	0
Cupcake	38	15	31	8	6	1	1	0
Cupcake	36	16	24	12	9	1	1	0
Cupcake	34	17	23	11	13	0	1	0

Cupcakes vs Muffins

The Challenge

Classify recipes as cupcakes or muffins. When given a new recipe, determine if it's a cupcake or a muffin.

The Steps

1. Find the data ✓
2. Apply a data science model
3. Review the results

Cupcakes vs Muffins

The Challenge

Classify recipes as cupcakes or muffins. When given a new recipe, determine if it's a cupcake or a muffin.

The Steps

1. Find the data ✓
2. Apply a data science model
3. Review the results

Python Script

jupyter SVM Demo Lab (changed 22 minutes ago (autosaved))

File Edit View Insert Cell Kernel Help

Cell CellEditor

Cupcakes vs Muffins with SVM

Step 1: Import Packages

```
In [1]: # Allow sheets to appear in the notebook
%matplotlib inline

# Libraries for analysis
import numpy as np
import pandas as pd
from sklearn import svm
```

Libraries for visualizations
import matplotlib.pyplot as plt
import seaborn as sns; sns.set(style="darkgrid")

Step 2: Import Data

```
In [2]: # Load in muffin and cupcake dependent data
muffins = pd.read_csv('https://raw.githubusercontent.com/justmarkham/DAT8/master/data/muffins.csv')
muffins.head()
```

Type	Fat	With Sugar	Ghee	Cup	Gelling Powder	Vanilla	Skin
Muffin	22	2	7	2	0	0	0
Muffin	21	12	4	9	1	0	0
Muffin	27	14	6	6	1	0	0
Muffin	22	12	4	9	2	1	0
Muffin	22	17	2	9	2	1	0

Step 3: Preprocess Data

```
In [3]: # Clean the ingredients
sns.lmplot("Type", "Vanilla", data=muffins, hue="Type", fit_reg=False, scatter_kws={"s": 100})
sns.despine(left=True)
```

```
In [4]: # Specify inputs for the model
muffin_features = muffins[['Type', 'Fat', 'With Sugar', 'Ghee']]
muffin_labels = np.where(muffins['Type'] == 'Muffin', 0, 1)
```

Step 4: Fit the Model

```
In [5]: # Fit the SVM model
model = svm.SVC(kernel='linear')
model.fit(muffin_features, muffin_labels)
```

```
In [6]: # SVC(0.1, muffin_features, muffin_labels).decision_function(muffin_features).mean()
# SVC(0.1, muffin_features, muffin_labels).support_vectors_.dot(muffin_features, axis=0).mean() + b
# SVC(0.1, muffin_features, muffin_labels).support_vectors_.dot(muffin_features, axis=0).mean() + 0.1
# b = -0.1
```

Step 5: Visualize Results

```
In [7]: # Get one separating hyperplane
w = np.array([0, 1])
w1 = w / np.linalg.norm(w)
w2 = -w / np.linalg.norm(w)

# Get the possible two separating hyperplanes that pass through the support vectors
SVC(0.1, muffin_features, muffin_labels).decision_function(muffin_features) / np.linalg.norm(w)
```

```
In [8]: # Get the Aggregates
sns.lmplot("Type", "Vanilla", data=muffins, hue="Type", palette="Set1", fit_reg=False, scatter_kws={"s": 100})
plt.plot(muffin_features, w1 * np.array([-10, 10]) + b)
plt.plot(muffin_features, w2 * np.array([-10, 10]) + b)
```

```
In [9]: # Get the Aggregates
sns.lmplot("Type", "Vanilla", data=muffins, hue="Type", palette="Set1", fit_reg=False, scatter_kws={"s": 100})
plt.plot(muffin_features, w1 * np.array([-10, 10]) + b)
plt.plot(muffin_features, w2 * np.array([-10, 10]) + b)
```



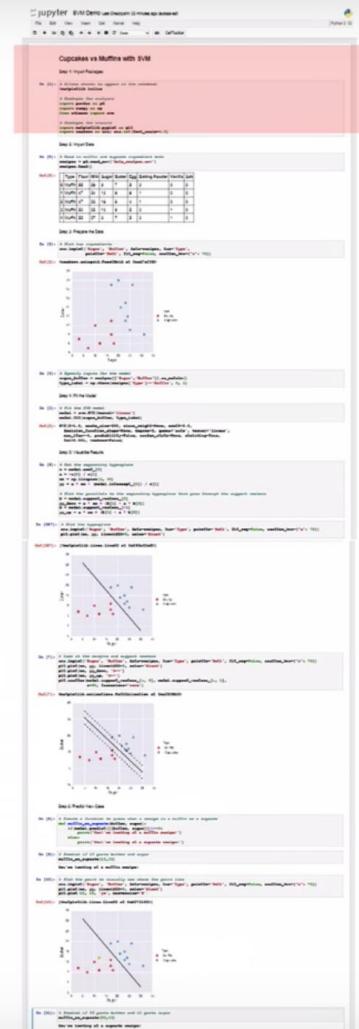
- a. Import Libraries
- b. Import Data
- c. Prepare the Data
- d. Fit the Model
- e. Visualize Results
- f. Predict New Case

a. Import Libraries

Cupcakes vs Muffins with SVM

Step 1: Import Libraries

```
In [13]: # Allows charts to appear in the notebook  
%matplotlib inline  
  
# Libraries for analysis  
import pandas as pd  
import numpy as np  
from sklearn import svm  
  
# Libraries for visuals  
import matplotlib.pyplot as plt  
import seaborn as sns; sns.set(font_scale=1.2)
```



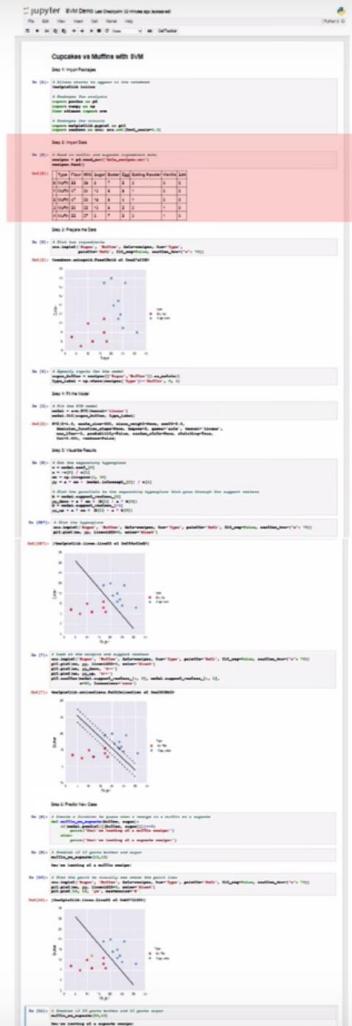
b. Import Data

Step 2: Import Data

```
In [2]: # Read in muffin and cupcake ingredient data  
recipes = pd.read_csv('data_recipes.csv')  
recipes.head()
```

Out [2]:

	Type	Flour	Milk	Sugar	Butter	Egg	Baking Powder	Vanilla	Salt
0	Muffin	55	28	3	7	5	2	0	0
1	Muffin	47	24	12	6	9	1	0	0
2	Muffin	47	23	18	6	4	1	0	0
3	Muffin	50	25	12	6	5	2	1	0
4	Muffin	55	27	3	7	5	2	1	0

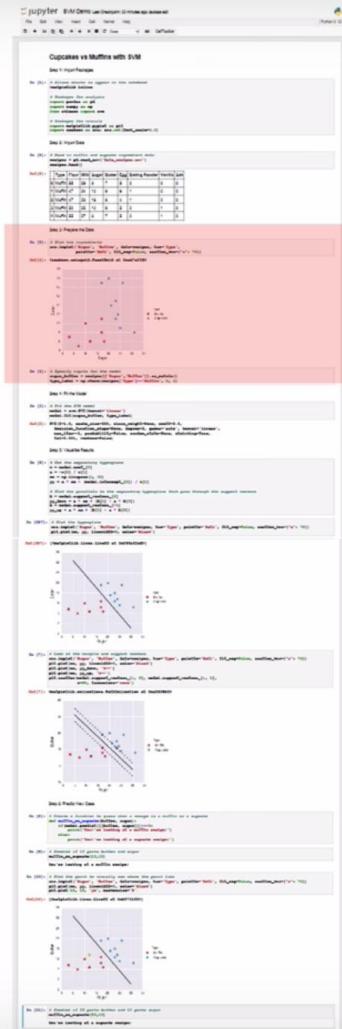
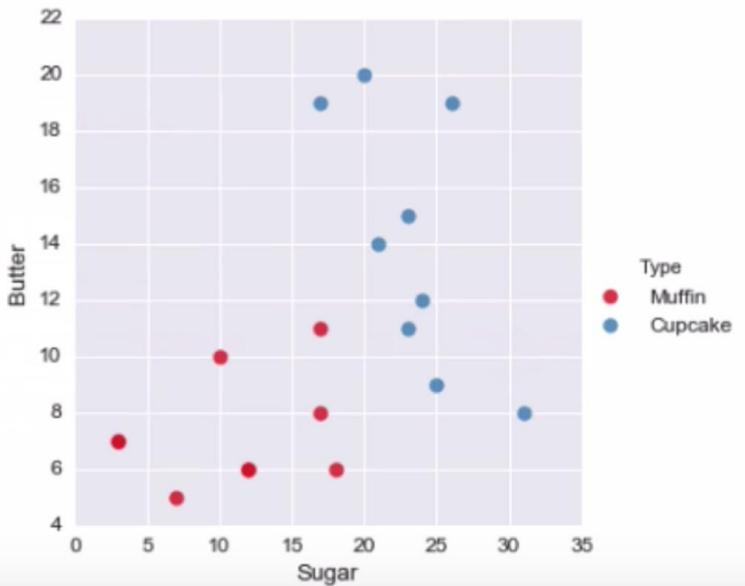


c. Prepare the Data

Step 3: Prepare the Data

```
In [3]: # Plot two ingredients  
sns.lmplot('Sugar', 'Butter', data=recipes, hue='Type',  
           palette='Set1', fit_reg=False, scatter_kws={"s": 70})
```

```
Out[3]: <seaborn.axisgrid.FacetGrid at 0xad7af28>
```



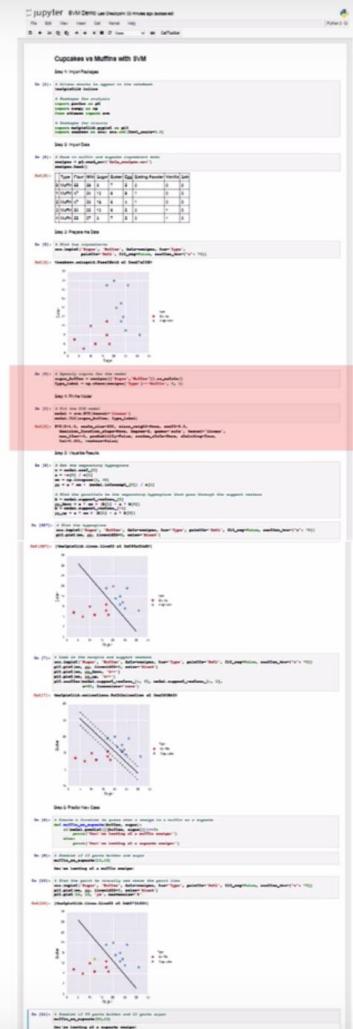
d. Fit the Model

Step 4: Fit the Model

```
In [4]: # Specify inputs for the model  
sugar_butter = recipes[['Sugar','Butter']].as_matrix()  
type_label = np.where(recipes['Type']=='Muffin', 0, 1)
```

```
In [5]: # Fit the SVM model  
model = svm.SVC(kernel='linear')  
model.fit(sugar_butter, type_label)
```

```
Out[5]: SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,  
           decision_function_shape=None, degree=3, gamma='auto', kernel='linear',  
           max_iter=-1, probability=False, random_state=None, shrinking=True,  
           tol=0.001, verbose=False)
```



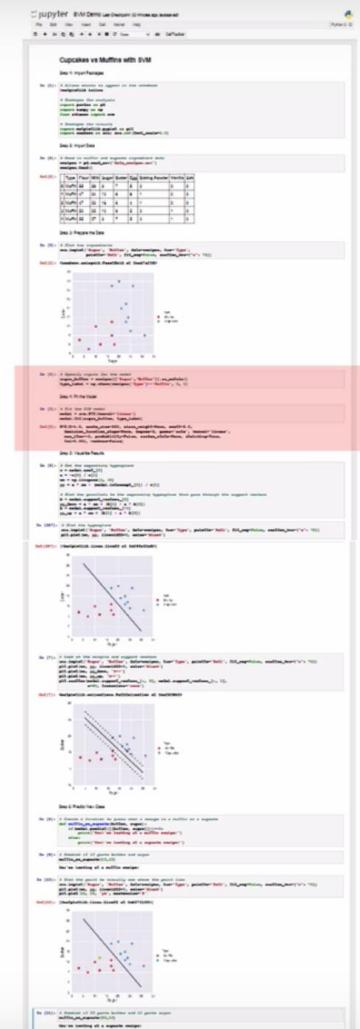
d. Fit the Model

Step 4: Fit the Model

```
In [4]: # Specify inputs for the model  
sugar_butter = recipes[['Sugar', 'Butter']].as_matrix()  
type_label = np.where(recipes['Type']=='Muffin', 0, 1)
```

```
In [5]: # Fit the SVM model  
model = svm.SVC(kernel='linear')  
model.fit(sugar_butter, type_label)
```

```
Out[5]: SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,  
           decision_function_shape=None, degree=3, gamma='auto', kernel='linear',  
           max_iter=-1, probability=False, random_state=None, shrinking=True,  
           tol=0.001, verbose=False)
```

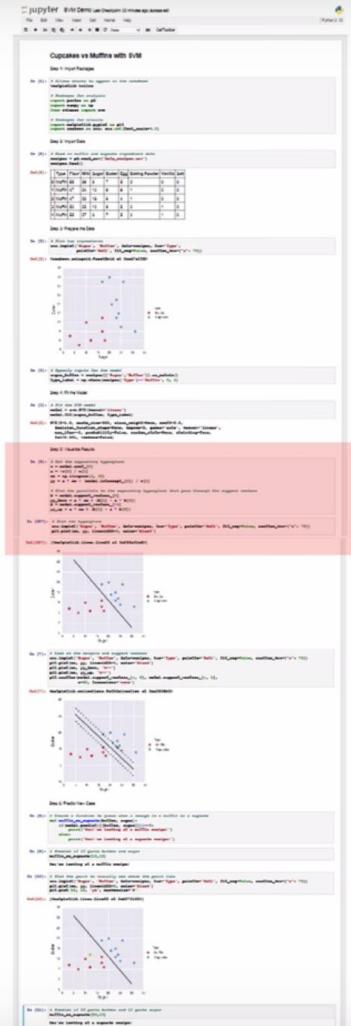


e. Visualize Results

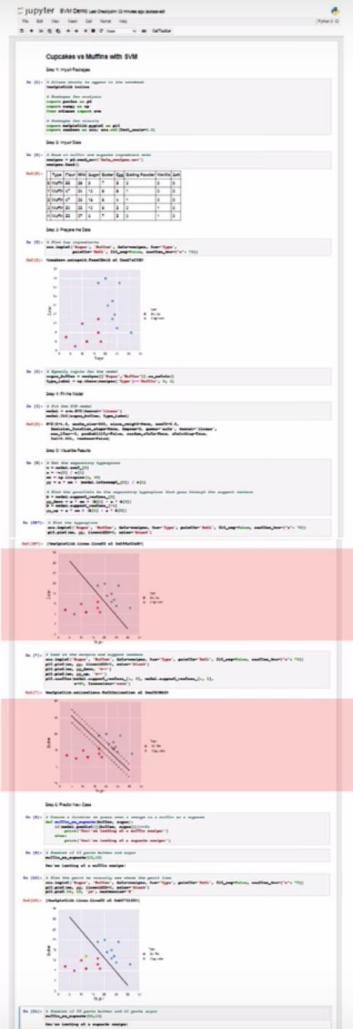
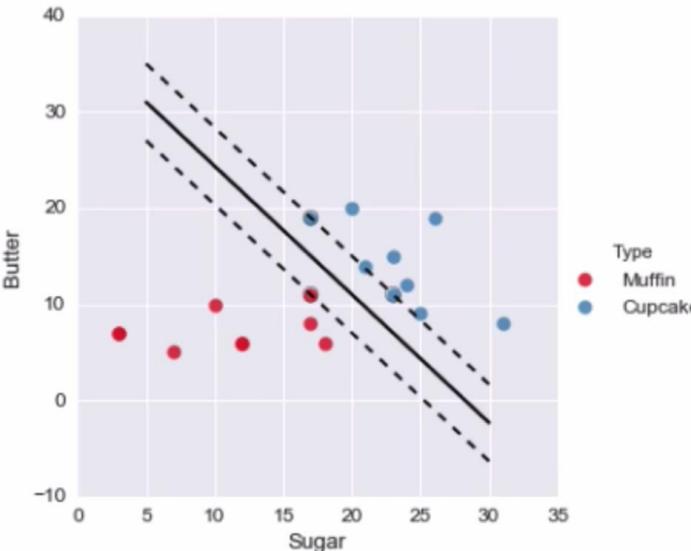
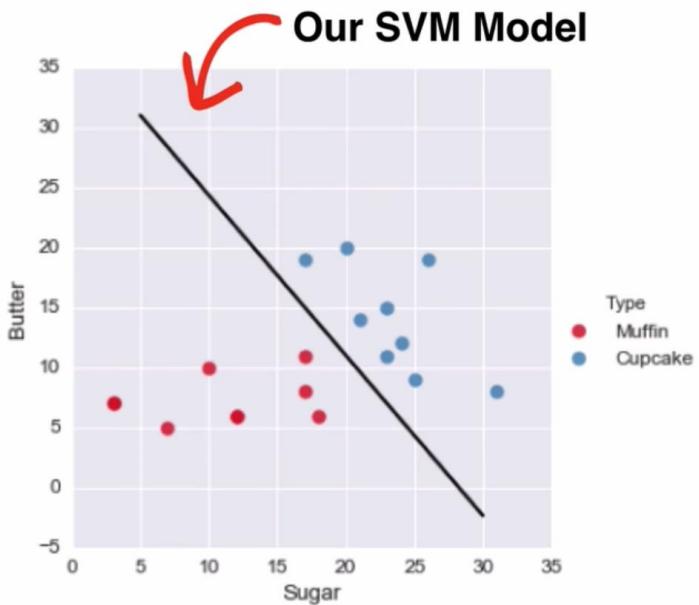
Step 5: Visualize Results

```
In [6]: # Get the separating hyperplane
w = model.coef_[0]
a = -w[0] / w[1]
xx = np.linspace(5, 30)
yy = a * xx - (model.intercept_[0]) / w[1]

# Plot the parallels to the separating hyperplane
# that pass through the support vectors
b = model.support_vectors_[0]
yy_down = a * xx + (b[1] - a * b[0])
b = model.support_vectors_[-1]
yy_up = a * xx + (b[1] - a * b[0])
```



e. Visualize Results



Cupcakes vs Muffins

The Challenge

Classify recipes as cupcakes or muffins. When given a new recipe, determine if it's a cupcake or a muffin.

The Steps

1. Find the data ✓
2. Apply a data science model ✓
3. Review the results

3. Review the Results

The Challenge

Classify recipes as cupcakes or muffins. ✓

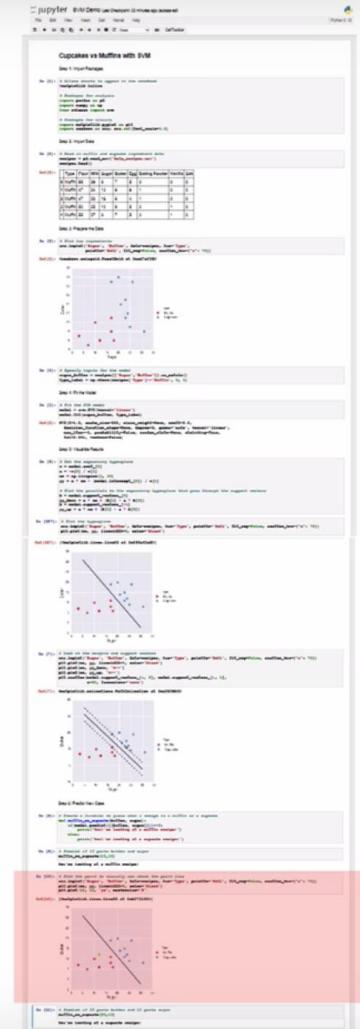
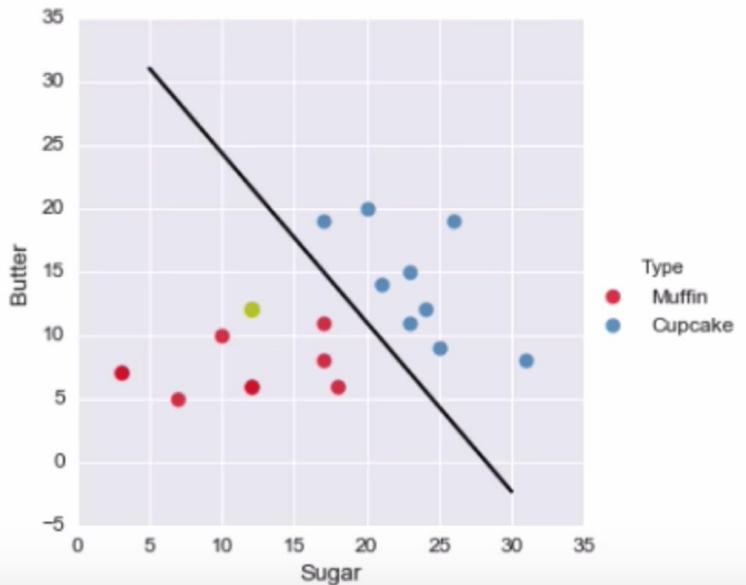
When given a new recipe, determine if it's a cupcake or a muffin.



f. Predict New Case

```
In [10]: # Plot the point to visually see where the point lies
sns.lmplot('Sugar', 'Butter', data=recipes, hue='Type',
            palette='Set1', fit_reg=False, scatter_kws={"s": 70})
plt.plot(xx, yy, linewidth=2, color='black')
plt.plot(12, 12, 'yo', markersize='9')
```

```
Out[10]: [
```



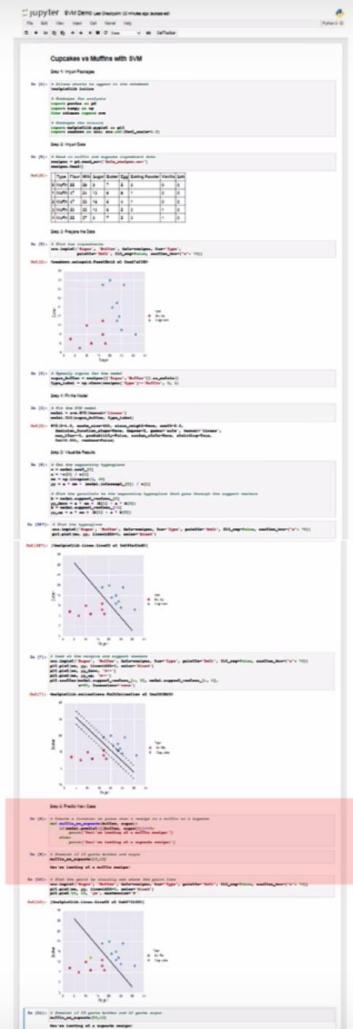
f. Predict New Case

Step 6: Predict New Case

```
In [8]: # Create a function to guess when a recipe is a muffin  
# or a cupcake using the SVM model we created  
def muffin_or_cupcake(butter, sugar):  
    if(model.predict([[butter, sugar]])==0):  
        print('You\'re looking at a muffin recipe!')  
    else:  
        print('You\'re looking at a cupcake recipe!')
```

```
In [9]: # Predict if 12 parts butter and sugar  
muffin_or_cupcake(12,12)
```

You're looking at a muffin recipe!



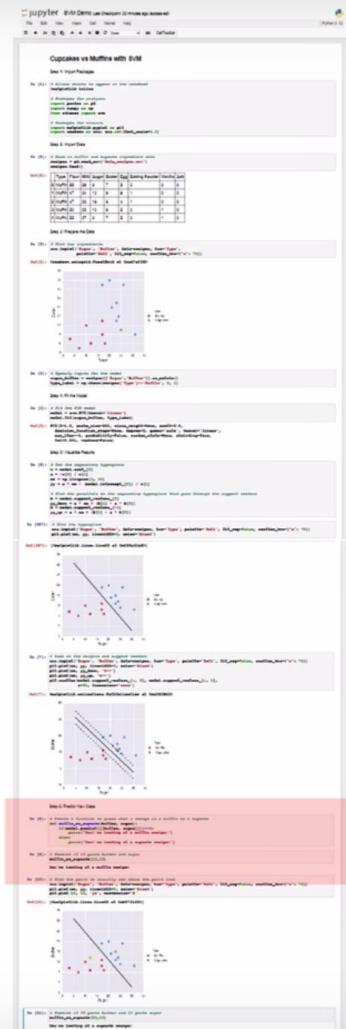
f. Predict New Case

Step 6: Predict New Case

```
In [8]: # Create a function to guess when a recipe is a muffin  
# or a cupcake using the SVM model we created  
def muffin_or_cupcake(butter, sugar):  
    if(model.predict([[butter, sugar]])) == 0:  
        print('You\'re looking at a muffin recipe!')  
    else:  
        print('You\'re looking at a cupcake recipe!')
```

```
In [9]: # Predict if 12 parts butter and sugar  
muffin_or_cupcake(12,12)
```

You're looking at a muffin recipe!



Cupcakes



Muffins



versus

Cupcakes and muffins can be classified using Support Vector Machines!

THANK YOU