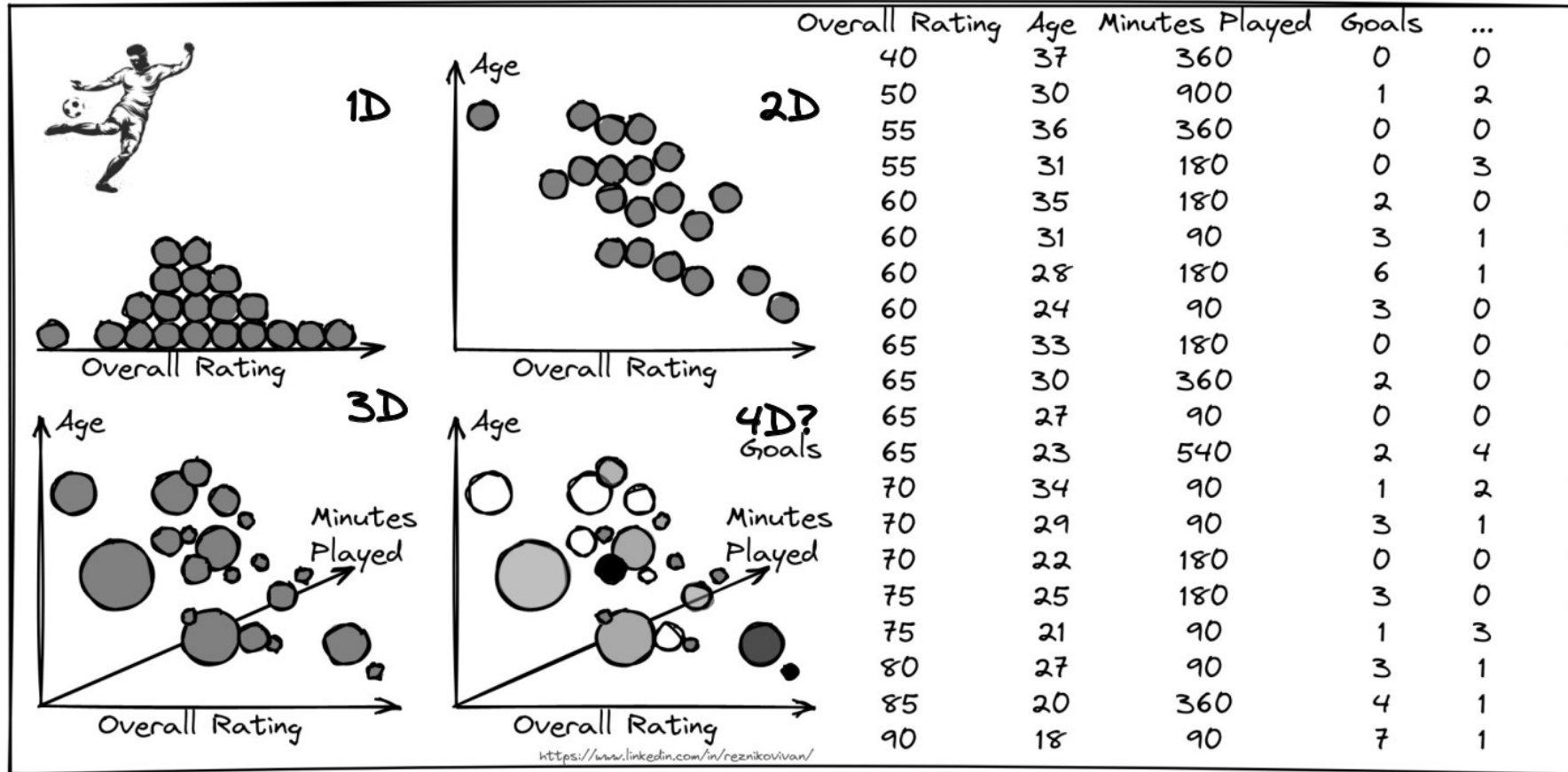


What is a Dimension?



Dimensions = features, attributes, variables, etc.

How many Dimension is a lot?

1. Generate random train data:

Size = 51, dimensions = 3, range (0,1)

```
In [3]: np_arr = np.random.rand(size,3)
np_arr
```

```
Out[3]: array([[0.69646919, 0.28613933, 0.22685145],
 [0.55131477, 0.71946897, 0.42310646],
 [0.9807642 , 0.68482974, 0.4809319 ],
 [0.39211752, 0.34317802, 0.72904971],
 [0.43857224, 0.0596779 , 0.39804426],
 [0.73799541, 0.18249173, 0.17545176],
 [0.53155137, 0.53182759, 0.63440096],
 [0.84943179, 0.72445532, 0.61102351],
 [0.72244338, 0.32295891, 0.36178866],
 [0.22826323, 0.29371405, 0.63097612],
 [0.09210494, 0.43370117, 0.43086276],
 [0.4936851 , 0.42583029, 0.31226122],
 [0.42635131, 0.89338916, 0.94416002],
 [0.50183668, 0.62395295, 0.1156184 ],
 [0.31728548, 0.41482621, 0.86630916],
 [0.25045537, 0.48303426, 0.98555979],
 [0.51948512, 0.61289453, 0.12062867],
 [0.8263408 , 0.60306013, 0.54506801],
 [0.34276383, 0.30412079, 0.41702221],
 [0.68130077, 0.87545684, 0.51042234],
 [0.66931378, 0.58593655, 0.6249035 ],
 [0.67468905, 0.84234244, 0.08319499],
 [0.76368284, 0.24366637, 0.19422296],
 [0.57245696, 0.09571252, 0.88532683],
 ...])
```

2. Generate target data:

Size = 51, dimensions = 1

count(0) = 26, count(1) = 25

3. Build 10 intervals (sections):

Group data in intervals using
0.1 window

4. Build "naive classifier":

default_forecast_value = 0

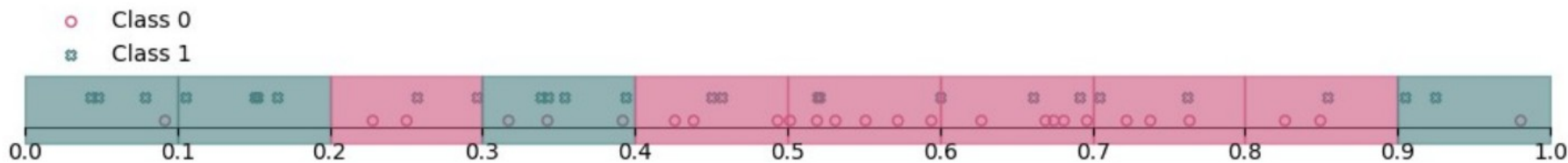
Logic: the most number of
points will set the class for
the interval. If equal number
of 0/1 values: class is set
to default_forecast_value

How many Dimension is a lot?

1 Dimension:

Misclassified points: 17

Empty sections: 0



How many Dimension is a lot?

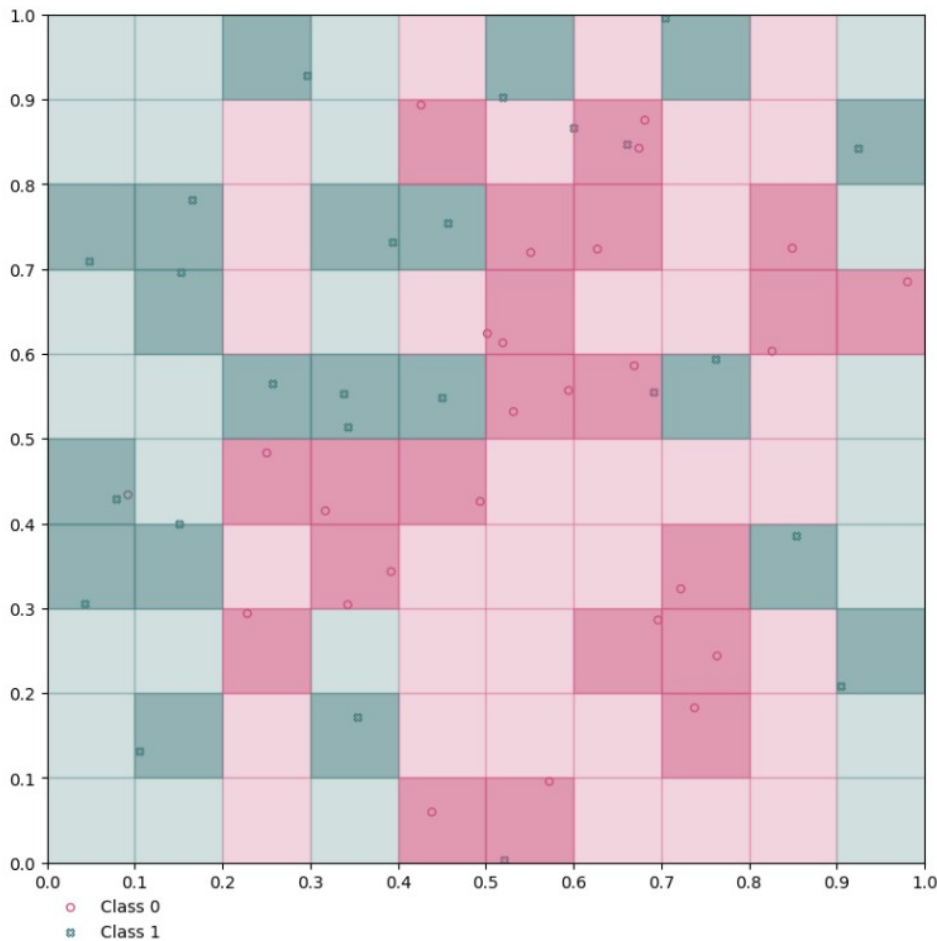
2 Dimensions:

Misclassified points: 5

Empty sections: 59

Is our classifier doing better? *No!!*

The data is already too sparse.



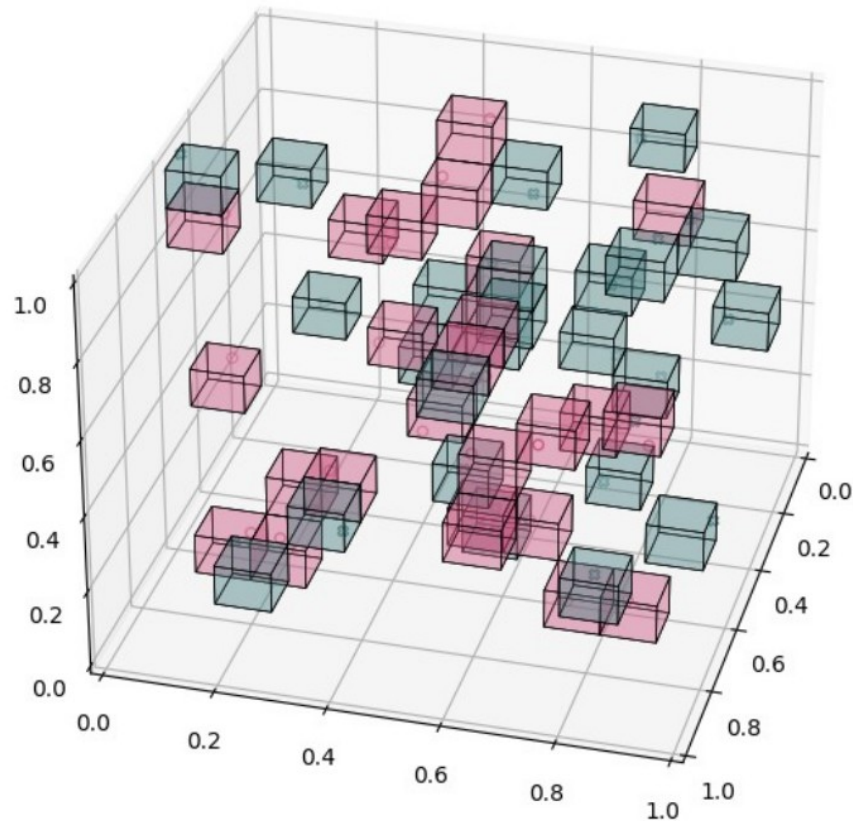
How many Dimension is a lot?

3 Dimensions:

Misclassified points: 1

Empty sections: 951

Though our naive classifier can correctly set 0/1 classes to 50 out of 51 points, it's pretty useless.



How many Dimension is a lot?

51 data points:

1 feature \rightarrow the density is 5.1 points per "box".

2 features \rightarrow 0.51 points per section.

3 features result in a density of 0.051 points per interval.

With more data, it becomes easier to separate it. We've almost perfectly separated 51 points using just 3 dimensions.

The results will be different if we use smaller interval ranges, but no matter what, it's always possible to separate $N+1$ points using N -dimensions.

In our case, it seems 2 dimension is already too much.