# The Curse of Dimensionality (Material taken from Pattern Recognition and ML by Bishop)

In most practical applications of pattern recognition, we have to deal with spaces of high dimensionality comprising many input variables. As we now discuss, this poses some serious challenges and is an important factor influencing the design of pattern recognition techniques.

In order to illustrate the problem we consider a synthetically generated data set representing measurements taken from a pipeline containing a mixture of oil, water, and gas (Bishop and James, 1993). These three materials can be present in one of three different geometrical configurations known as 'homogenous', 'annular', and 'laminar', and the fractions of the three materials can also vary. Each data point comprises a 12-dimensional input vector consisting of measurements taken with gamma ray densitometers that measure the attenuation of gamma rays passing along narrow beams through the pipe.

Figure 1 below shows 100 points from this data set on a plot showing two of the measurements $x_6$ and $x_7$ (the remaining ten input values are ignored for the purposes of this illustration). Each data point is labelled according to which of the three geometrical classes it belongs to, and our goal is to use this data as a training set in order to be able to classify a new observation ($x_6$, $x_7$), such as the one denoted by the cross in Figure below. We observe that the cross is surrounded by numerous red points, and so we might suppose that it belongs to the red class. However, there are also plenty of green points nearby, so we might think that it could instead belong to the green class. It seems unlikely that it belongs to the blue class. The intuition here is that the identity of the cross should be determined more strongly by nearby points from the training set and less strongly by more distant points. In fact, this intuition turns out to be reasonable and will be discussed more fully in later chapters.

How can we turn this intuition into a learning algorithm? One very simple approach would be to divide the input space into regular cells, as indicated in Figure 2. When we are given a test point and we wish to predict its class, we first decide which cell it belongs to, and we then find all of the training data points that
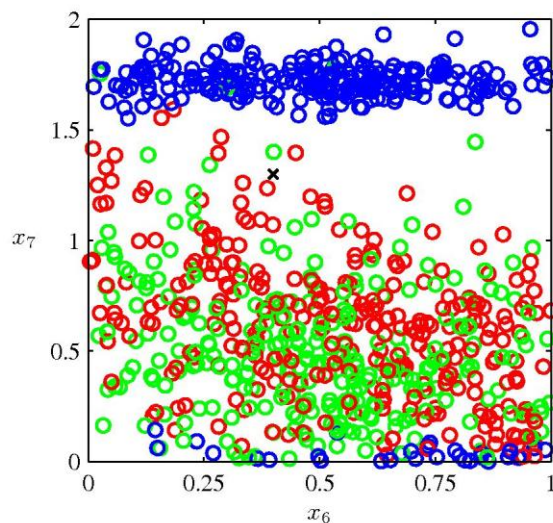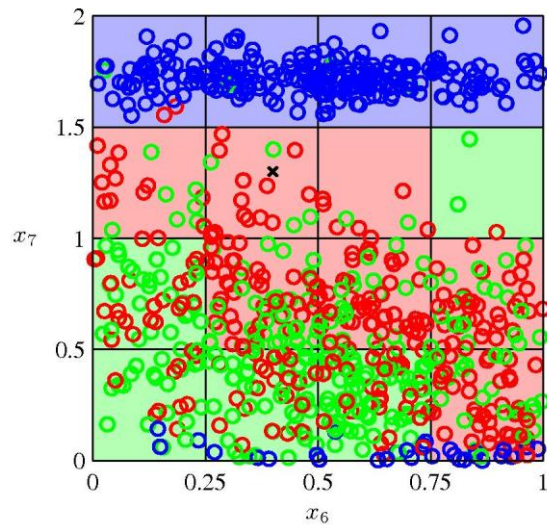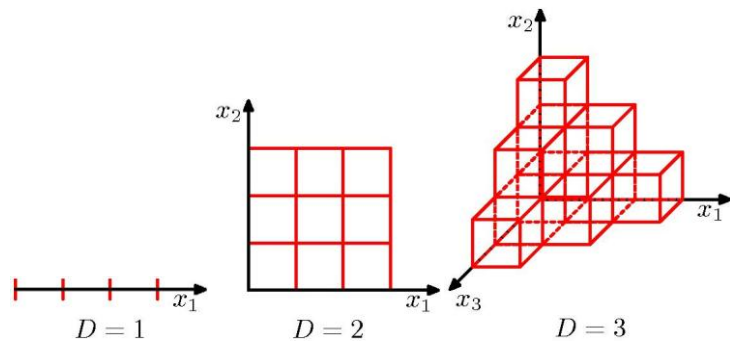


Figure 1
Scatter plot of the oil flow data for input variables $x_6$ and $x_7$, in which red denotes the 'homogenous' class, green denotes the 'annular' class, and blue denotes the 'laminar' class. Our goal is to classify the new test point denoted by '$\times$'.

**Figure 2** Illustration of a simple approach
to the solution of a classification
problem in which the input space
is divided into cells and any new
test point is assigned to the class
that has a majority number of representatives
in the same cell as
the test point. As we shall see
shortly, this simplistic approach
has some severe shortcomings.

fall in the same cell. The identity of the test point is predicted as being the same
as the class having the largest number of training points in the same cell as the test
point (with ties being broken at random).

There are numerous problems with this naive approach, but one of the most severe
becomes apparent when we consider its extension to problems having larger
numbers of input variables, corresponding to input spaces of higher dimensionality.
The origin of the problem is illustrated in Figure 1.21, which shows that, if we divide
a region of a space into regular cells, then the number of such cells grows exponentially
with the dimensionality of the space. The problem with an exponentially large
number of cells is that we would need an exponentially large quantity of training data
in order to ensure that the cells are not empty. Clearly, we have no hope of applying
such a technique in a space of more than a few variables, and so we need to find a
more sophisticated approach.

**Figure 3** Illustration of the curse of dimensionality, showing how the number of regions of a regular grid grows exponentially with the dimensionality $D$ of the space. For clarity, only a subset of the cubical regions are shown for D=3

CoD in context of DM - when the computation time of DM algorithm increases <u>exponentially</u> with the number of dimensions.