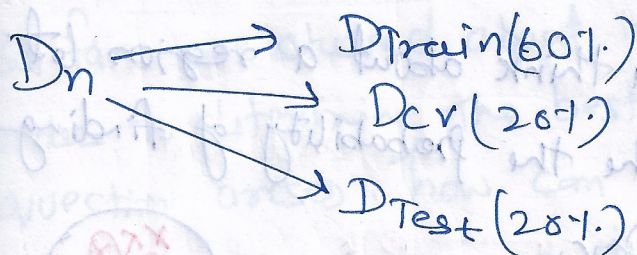


## How to Visualize Train, CV and test dataset

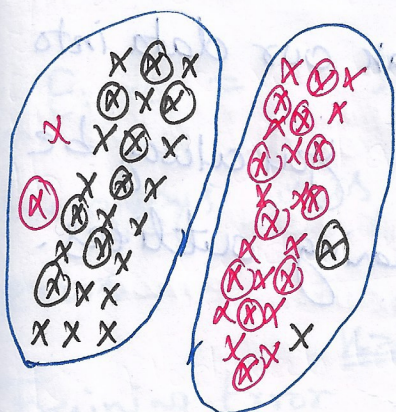


$\times$ : -ve datapoint in  $D_{train}$

$\otimes$ : +ve data point in  $D_{train}$

$\otimes$ : -ve datapoint in  $D_{cv}$

$\otimes$ : +ve data point in  $D_{cv}$



→ As we have randomly selected the points  
- we have both the test data and cross validation points close to each other.  
- we do also have few outliers.

- ①  $D_{train}$  and  $D_{cv}$  do not overlap perfectly.
- ② If there are many +ve/point from  $D_{train}$  in a region, then it is highly likely to find many +ve point from  $D_{cv}$  in the same region.
- ③ If there are very few +ve/-ve point in a region from  $D_{train}$ , it is very unlikely +ve/-ve points from  $D_{cv}$  in such region.



Such points are c/an Outliers / noise points

④ Intuitively, if you think about a region lot of point  $D_{train}$ , the the probability of finding lot of point from  $D_{cv}$ .

i.e. Density of true pts from  $D_{train}$   $\uparrow$  then the probability of true points from  $D_{cv}$  is higher.



⑤ Similarly, if density is less, then you may find the point  $D_{cv}$ .

So, in a nutshell, if we break our data into multiple parts  $D \rightarrow \begin{matrix} D_{train} \\ D_{cv} \\ D_{test} \end{matrix}$  you always be having outliers.

