How do we validate ML algorithms?

Typical setting
1. Data given (i.e. tables)
2. Choose & train a method
3. Test it!
4. if good enough, release it

How do we make sure of this?

do this by Cross Validation

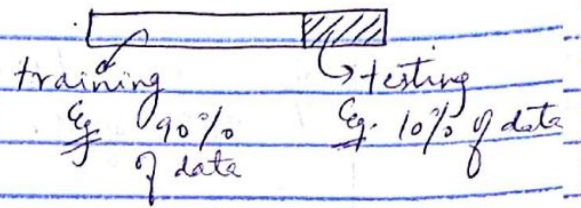Cross Validation
(hold out method)
(rotation estimation)
Given entire data ☐

input/features + targets
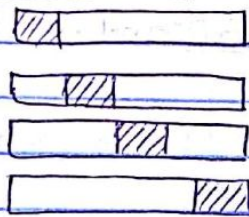(unsupervised)
(supervised)

① <u>hold out</u> : keep one part for testing

training → testing
Eg. 90% of data    Eg. 10% of data

we hold out a particular part of data for testing & use the rest of for training.

This is not reliable : the part you take may be biased (lucky / unfortunate)
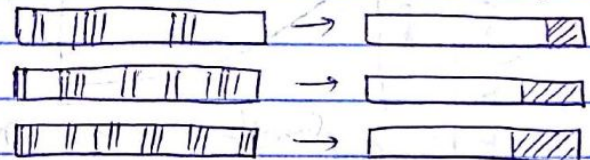
② <u>Cross Validation</u>.

(smaller part)
use the ▨ part for! testing and ☐ part for training

4 fold cross validation "k-fold" cross validation.    Usually between 3-5 fold.

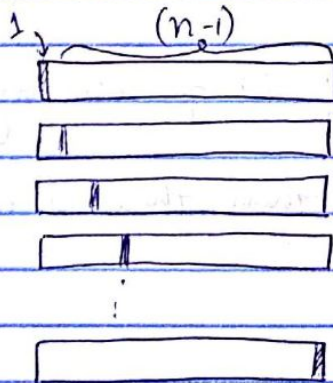③ <u>Random Sampling</u>.

We do not slice as in cross validation. (we do random sampling)

(same result as in cross validation)
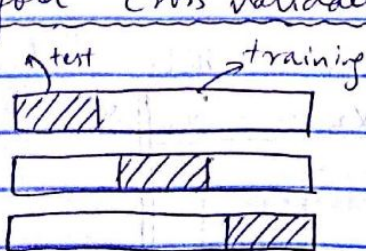
④ <u>Leave-one-out validation</u>

1,    (n-1)

❶ we use $n-1$ for training and only one for testing.... repeat this "$n$" number of times.

❷ very exhaustive.   (very)

❸ done when there is small data set

<u>k-fold cross validation</u>

↗test   ↗training    (training size is always > ▨ testing)

error (testing) = $e_1$
error (testing) = $e_2$        $E = \frac{1}{B} \sum e_i$
error (testing) = $e_3$

average error and standard deviation (very imp)

error of $10\% \pm 2\%$ is better than error of $4\% \pm 35\%$

eg.    $e = [10\%, 22\%]$ (confidence interval) with confidence $0.05 \Rightarrow 95\%$

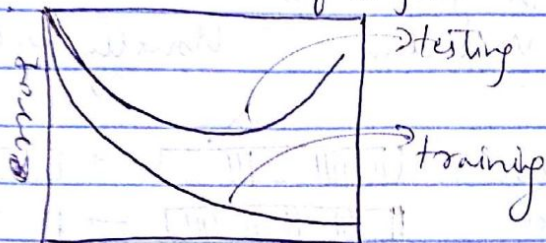95% of the times, the error generated is between 10% to 22%

## Usually:

$\mathcal{E} =$ acceptable error.



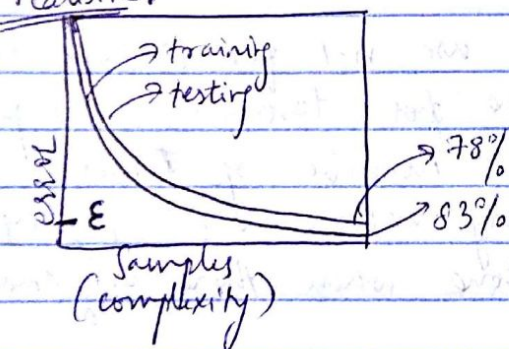if gap is small, the training is done properly, (it has learnt properly)

this shows overfitting:



(Complexity → neural n/w with 30 layers or 3 layers?)

memorized data, not learnt.

Subnetworks in human brain change their configuration

realistic:



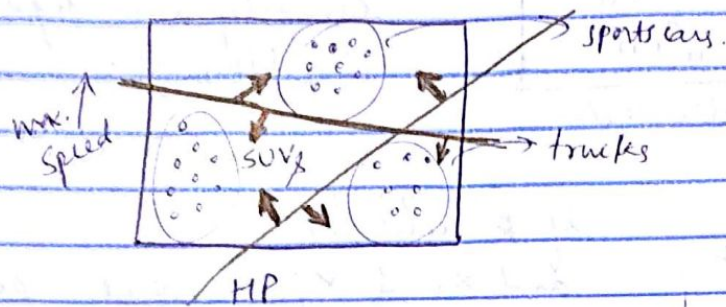if traing accuracy is < testing accuracy then there is a bug

## Clustering

Grouping given "unlabelled" data into "clusters" of similar objects.



no output.
→ (generally clustering)
unsupervised ≡ clustering

Eg horse power and maximum speed.



② Nicely separated
(space between clusters)

② linearly separated.
(we can draw lines, by just 2 lines, we have separated them. we need not know the actual contour

lines in feature space acting as decision boundaries

lines in hyper dimensional space → hyperplanes.

separation → distinction → understanding.

finding hyperplanes to separate clusters means we understand the meaning of features.

Clustering ⟨ ~ we need to know the # of clusters
don't ~~need~~ to ~~know~~
know → (hierarchical clustering)
→ (density based → finds dense regions in feature space
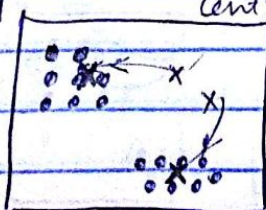
K-Means algorithm.
finds ⓚ centroids of data.
↳ has to be told how many.

① Randomly place k centroids.
② Assign each data point to a centroid.
[ Similarity measure = L2 (euclidean) ]
③ Update the ~~centroid~~ centroids
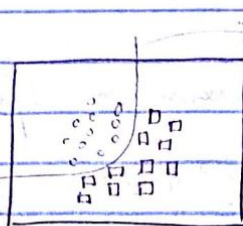


update the centroid

Centroid = prototype of the cluster.
= average of all data points that belong to that centroid.

→ curve, not a line!

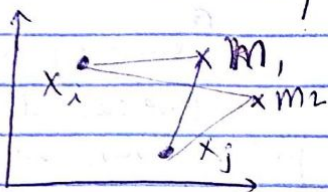if we have.

linearly not separable.
∴ difficult.

↳ this is complicated.

if $k$ = centroid $m_1 \& m_2$
and $x_i \& x_j$ are data points.

$$d(i,j) = \frac{\|x_i - x_j\|}{= \sqrt{|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + \cdots}}$$

Error:

$$E = \sum_{k=1}^{K} \sum_{x=X} \|x - m_k\|_{L_2}$$

(minimize this)   eucl. dist between all data points & centres.

trail of centroids: