

## Cat boost.

↓  
The library which handles  
Categorical data automatically.

Out-of-the-box support for the more descriptive data  
formats.

Catboost implements Symetric trees.  
(helps decreasing prediction time).

Default max-depth = 6

Dataset ordered in time

o (Catboost creates an artificial time for each datapoint)

time	datapoint	label
	$x_1$	10
	$x_2$	12
	$x_3$	9
	$x_4$	7
	$x_5$	52
	$x_6$	22
	$x_7$	33
	$x_8$	31
	$x_9$	32
	$x_{10}$	12

10 points  $\rightarrow$  9 models.

X Computationally expensive

$\log(\text{no. of datapoints})$  models

$$\Rightarrow \log(10) = 2.7$$

Model ~~order~~ trained on  $n$  data  
points is used to calculate the  
residuals

Ordered boosting

Catboost divides a data into Random permutations  
and applies ordered boosting.  
↓  
default = 4.

bagging-temperature.  $\rightarrow$  tuning parameter for randomness.

feature	label
$x_1$	10
$x_3$	15
$x_1$	20
$x_3$	15
$x_1$	30
$x_2$	20

response  
Coding.

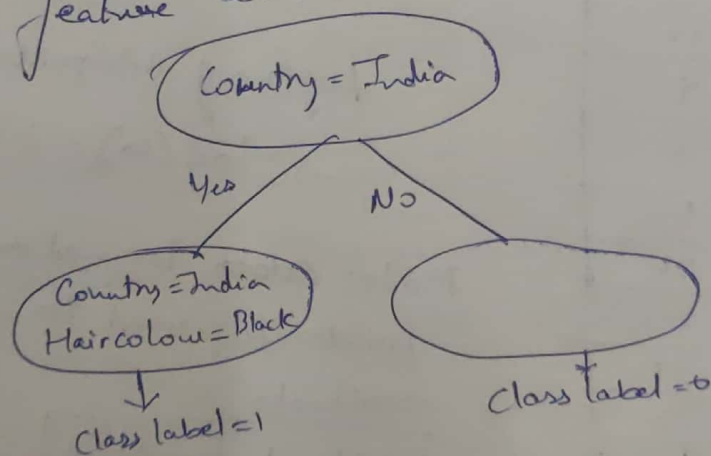
feature	label
20	10
15	15
20	20
15	15
20	30
20	20

This leads to target leakage

↓  
To avoid this, only datapoints that are past in time to a data point are considered.

Sometimes 0, → Use Laplace smoothing.

Categorical feature Combinations.



One-hot Encoding → onehot-max-size = N.

Handling Numerical features → Same as other Algos.  
(best split based on Information Gain).

### Limitations

- 1) Not suitable for Sparse matrices.
- 2) Not suitable if it has many numerical features  
(Catboost takes more time to train than Light GBM).

# T-SNE

T-distributed Stochastic Neighborhood Embedding.

Best for visualization of Data.

PCA:  $\rightarrow$  Basic }  $\rightarrow$  2-dim.  
MNIST } Christopher Olah blog.

X MDS, Sammon mapping, Graph-based techniques  $\leftarrow$  20 yrs.

✓ t-SNE  $\rightarrow$  2008, Geoffrey Hinton.

d-dim  $\rightarrow$  2d & 3d  
t-SNE  
Clean visualization.

PCA: Preserve global shape of data

t-SNE - local shape / struc. preserve.

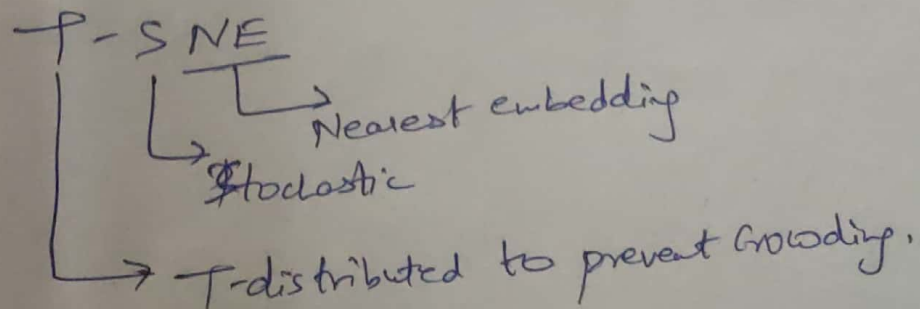
(+ global struc. by changing 1 parameter)

Neighborhood; Embedding  $D \rightarrow \boxed{x, y} \rightarrow$  same value

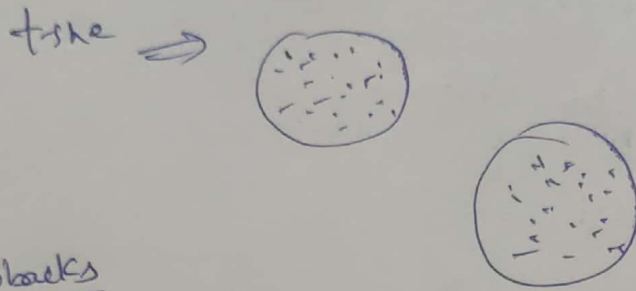
T-SNE is not a Deterministic algorithm

T-SNE is a Probabilistic algorithm

Every time you run T-SNE with a step, iteration,  
it might turn out different.



t-SNE  $\rightarrow$  expands dense clusters  
shrinks sparse clusters



### Drawbacks

- ① One of the things you cannot read from t-SNE is whether a cluster is dense or sparse.  
i.e. You cannot come to some density conclusion using t-SNE.
- ② t-SNE does not preserve distance between clusters.

### Thumb rule

- ① Run steps/iter till shapes stabilize
- ② perplexity  $2 \leq p < n$   
Never Run t-SNE just once and start reading it.
- ③ re-run t-SNE with  $(p, \text{step})$  multiple times.

trafjal knot

Epsilon  $\rightarrow$  How fast you should change from one iteration to another.