

Linear Regression	vs	logistic Regression.
* when dependent Variable is continuous eg: 2.1, 3.6		when dependent Variable is binary. eg : 0, 1, 2.
↓		↓
* Straight line is used to separate the data.		some linear function of independent variables is used to predict
* Dependent and independent Variables is required to be linearly related		* no need of linear relation bw Dependent & independent Variables
* independent variable could be correlated		* independent should not be correlated.
↓		↓

$$Y = b_0 + x_1 b_1 + b_2 x_2 \dots$$

$$y\_model = tf.matmul(w, x) + b$$

$$cost = tf.square(Y - y\_model)$$

$$Y = b_0 + x_1 b_1 + x_2 b_2 \dots$$

$$P = \frac{1}{1+e^y} ; \ln\left(\frac{P}{1-P}\right)$$

$$y\_model = tf.sigmoid(tf.matmul(x, w) + b)$$

$$cost = tf.reduce_mean(-tf.reduce_sum(Y * tf.log(y\_model)))$$

### Soft max.

Same is the case with soft max, where as in logistic regression label will be 0 or 1 only & chances on the other hand in softmax we can have multiple based on index place of the output  
 eg  $Y = [ [1, 0, 0], [0, 1, 0], [0, 0, 1] ] \rightarrow 3 \text{ labels}$

$$y\_model = tf.nn.softmax(tf.matmul(w, x) + b)$$

$$cost = tf.reduce_mean(-tf.reduce_sum(Y * tf.log(y\_model)))$$

obs: of these '3' models soft max is best fit for classification of data which is linearly separable so that we can classify in numbers

Note: Data should be linearly separable otherwise weights will be NaN.

## K nearest neighbour (KNN)

It is based on similarity measure eg: distance functions

KNN has been used in statistical estimation and pattern recognition.  
distance functions:

$$\text{Euclidean } \sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

$$\text{Manhattan } \sum_{i=1}^k |x_i - y_i|$$

How to choose K value :

Generally K is more means more precise but no guarantee else go for cross-validation method generally  $K = (3-10)$  better  $K = \text{odd}$  to avoid tie

Drawback

distance measure from direct dataset might be a problem if the data is mixture of numerical and categorical values.

eg data is having annual income in dollars and age in years  
then income have higher influence on distance calculated.

Sol: Solution is to standardize the data set.

Standardization.

$$x_i = \frac{x - \text{min}}{\text{max} - \text{min}}$$

Note: In KNN there will be no weights because for every test case we are calculating distance b/w all train data and placed based on that distance.

distance = tf.reduce\_sum(tf.abs(tf.add(xtr, tf.negative(xte))), reduction\_indices=1)

pred = tf.arg\_min(distance, 0)

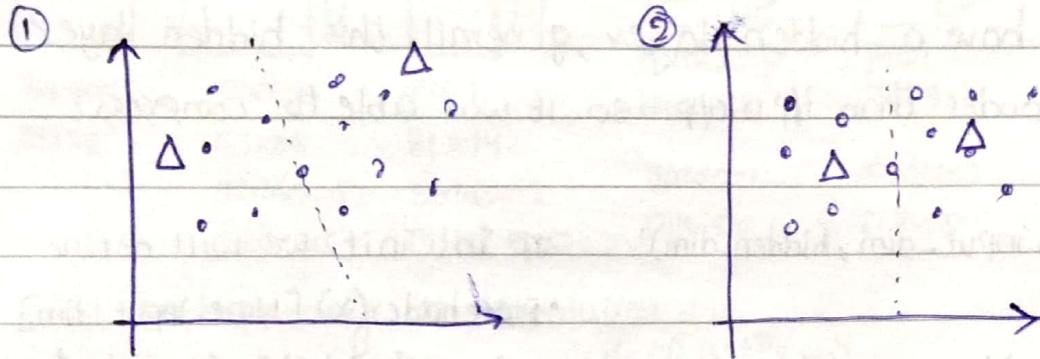
for i in range(len(xte)):

nn\_index = sess.run(pred, feed\_dict={xtr: xtr, xte: xte})

print('test', i, 'prediction:', np.argmax(Ytr[nn\_index]),  
'True class:', np.argmax(Yte[i]))

## Kmeans clustering

Clustering comes under unsupervised learning where data is mapped as clusters (groups), where this clustering (grouping) is done in a procedure.



- In the first stage random positions of clusters are taken and every point is compared in terms of distance to the clusters based on least distance it is "mapped" to it
- In second stage clusters are "recentered" in such a way that it will be in optimal distance from all points mapped to it. This process will repeat till no longer clusters are recentered.

`kmeans = KMeans(inputs=X, num_clusters=k, distance_metric='cosine', use_minibatch=True)`

`training_graph = kmeans.training_graph()`

`(all_scores, cluster_idx, scores, cluster_centers_initialized, cluster_centers_var, init_op, train_op) = training_graph()`

`cluster_idx_k = cluster_idx[0]`

`avg_distance = tf.reduce_mean(scores)`

- Start session run `init`, `init_op` to initialize variables and random clusters
- for  $i$  in range ( $1$ ,  $\text{num\_steps}+1$ ):

- , `d, idx = sess.run([train_op, avg_distance, cluster_idx], feed_dict={X: X_batch})`

Training completed then • Assigning label to ~~their closest~~ each centroid.

- Sample to their closest centroid given by `idx` • Assign most frequent label to centroid.
- Now test and predict accuracy if test data is available.

## Auto encoder :

Auto encoder is like a .mp3 file which is of less in size but still given similar sound o/p, while this o/p of encoded is not exactly same as the one which is before encoding.

Auto encoders will have a hidden layer, generally this hidden layer will be having fewer nodes than i/p or o/p so it was able to compress

class Autoencoder:

```
def __init__(self, input_dim, hidden_dim):    # In init we will define
    pass                                         • placeholder(x) [None, input_dim]
def train(self, data):                         • tf.name_scope('encode') → w[i, h], b[h], encoded → tanh
    pass                                         • tf.name_scope('decode') → w, b, decoded → tanh.
def test(self, data):                          • self.loss → square error   • bam.op → RMS optimizer
    pass
```

# In train the autoencoder:

for i in range(epoch): → len(data)

- start session & run variables → for j in range(num\_samples):

- Last save the session for reuse, !, - = sess.run([self.loss, self.bam.op]) → feed\_dict = {self.x: [data[j]]}

# In testing we restore train session.

- and we will run encode & decode feeding test data.

In this method we can compress data and reproduce it also

e.g. i/p → [8, 4, 6, 2] → hvalue

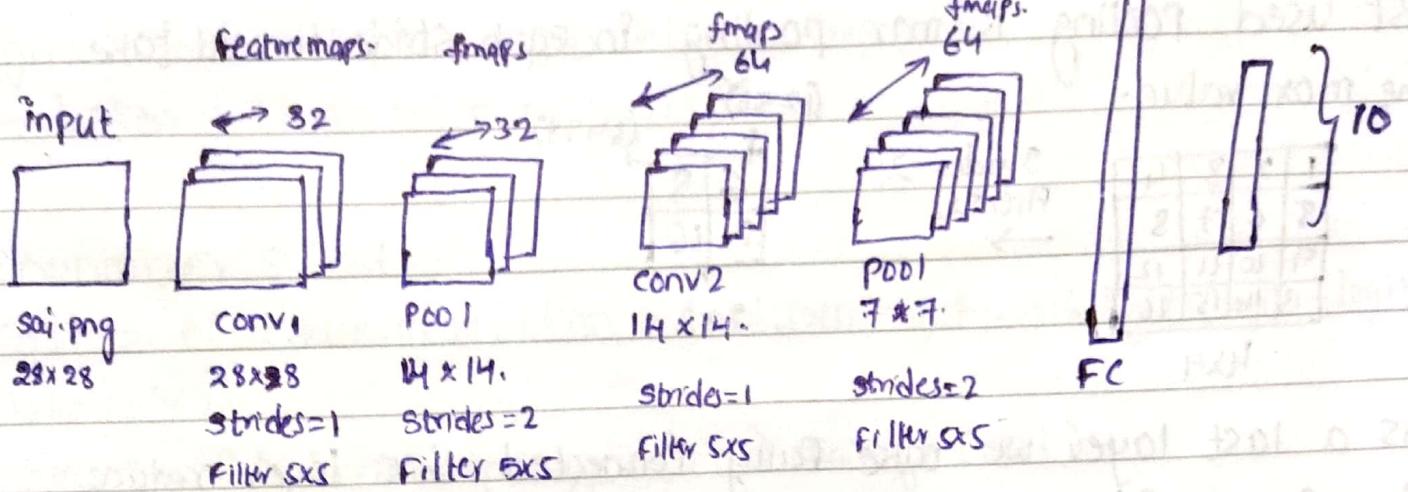
compressed → [0.78238308] → lvalue.

reconstructed [6.87, 2.798, 6.25, 2.23]

It is better to make it batchwise training

Autoencoders are unsupervised Neural Networks that use ML to do compression for us. The aim of autoencoders is to learn a compressed, distributed representation for given data for purpose of dimensionality reduction.

# Convolution Neural Networks (CNN)

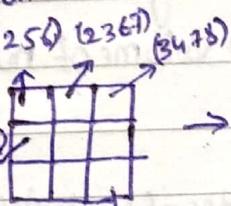


What is happening at convolution:

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Image of  $(4 \times 4)$

using  
filter  $(2 \times 2)$   
strides = 1  
without padding.



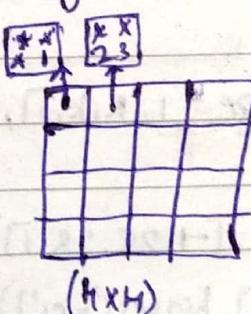
filter size is specified in weights  
and strides are specified in  
conv2d and max pool  
so in without padding size will  
reduce.

If padding = 'SAME' then it will automatically adjust the size in different ways

*	*	*	*
*	1	2	3
5	6	7	8
9	10	11	12
13	14	15	16

$(4 \times 4)$  with  
added due  
to same padding.

using  
filter  
 $(2 \times 2)$   
strides=1  
padding = SAME



If Padding is 'same'  
strides is  $\rightarrow 1 \rightarrow$  size will  
not change  
strides is  $\rightarrow 2 \rightarrow$  size will  
divide by 2.

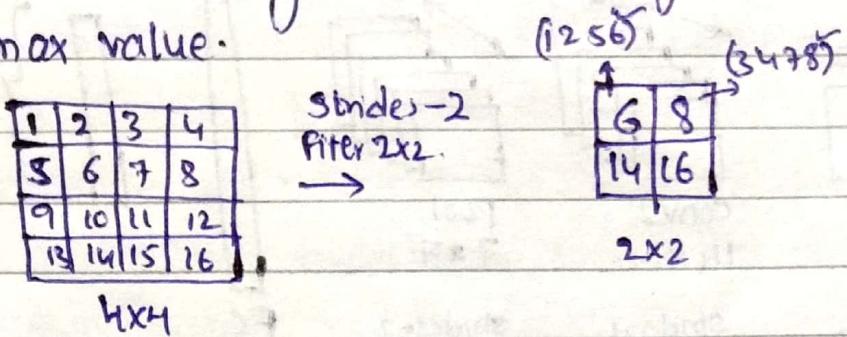
After convolution of image it will be called as feature maps.

size of filters and no of filters are specified in weights  $[5, 5, 1, 32]$

$\downarrow$   $\downarrow$   $\downarrow$   
filter  $\leftarrow (5 \times 5)$  1 image creates  
32 fmaps

After convolution we usually do Pooling to reduce dimensionality, this enables to reduce the dimensionality, no of parameters & overfitting & time taken for process.

most used pooling is max pooling in each stride it will take the max value.



As a last layer we take fully connected this layer reduce the fmaps of (84) into one big 1-D array called flattening  
flattening is arranging 3D volume of numbers into 1D vector

- conv2d :

```
x = tf.nn.conv2d(x, W, strides=[1, strides, strides, 1], padding='SAME')
```

```
x = tf.nn.bias_add(x, b)
```

```
return tf.nn.relu(x)
```

- maxpool2d :

```
tf.nn.max_pool(x, ksize=[1, K, K, 1], strides=[1, strides, strides, 1], padding='SAME')
```

- conv-net :

```
x = tf.reshape(x, shape=[-1, 28, 28, 1]) # input x is 1D now we convert it to 4D
```

```
conv1 = conv2d(x, weights['wci'], biases['bci']) (-1)*(width)*(height)*(colour-channel-number)
```

```
conv1 = maxpool(conv1, K=2)
```

```
conv2 = conv2d(conv1, weights['wc2'], biases['bc2']) '1' is standard feature from numpy
```

```
conv2 = conv2d(conv2, K=2)
```

same as (7\*7\*64)

```
fcl = tf.reshape(conv2, [-1, weights['wd1']].get_shape().as_list[0]])
```

```
fcl = tf.add(tf.matmul(fcl, weights['wd1']), biases['bd1'])
```

```
fcl = tf.nn.relu(fcl) # apply dropout if need fcl = tf.nn.dropout(fcl, dropout)
```

```
fcl = tf.add(tf.matmul(fcl, weights['out']), biases['out'])
```

return out

#model creation

logits = conv\_net(x, weights, biases, keep\_prob)

prediction = tf.nn.softmax(logits)

#optimizer & cost

cost = tf.reduce\_mean(tf.nn.softmax\_cross\_entropy\_with\_logits(logits=logits, labels=Y))

train\_op = tf.train.AdamOptimizer(0.01).minimize(cost)

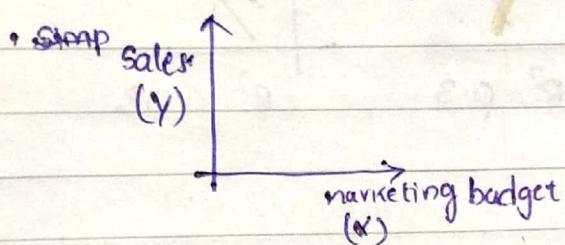
## Types of ML

- ① Regression  $\rightarrow$  continuous Variable
- ② Classification  $\rightarrow$  categorical Variable
- ③ clustering  $\rightarrow$  NO labels only on properties/feature

Supervised  $\xleftarrow{\text{Regression}}$  classification unsupervised  $\rightarrow$  clustering.

### • Regression:

- ① Simple linear regression.
- ② multiple linear regression.



$X \rightarrow$  independent Var (Predictor)  
 $Y \rightarrow$  Dependent Var (target var)

Optimal line  $\rightarrow$  best fit line

Residual  $\rightarrow$  error  $\rightarrow$  difference b/w actual and prediction.

$$e_i = y_i - y_{\text{pred}} \Rightarrow e_1^2 + e_2^2 + e_3^2 + \dots = \text{RSS} \text{ (Residual sum of squares)}$$

$$\text{RSS} = e_1^2 + e_2^2 + e_3^2 + \dots + e_n^2$$

$$e_i = y_i - y_{\text{pred}}$$

$$y_{\text{pred}} = \beta_0 + \beta_1 x_i \quad \therefore \text{RSS} = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

• Our linear model  $\Rightarrow y = mx + c$ .

To determine how good is this model by  $\frac{R^2}{RSE}$

coefficient of determination  
 $R^2 \rightarrow$  R-square (var)

$RSE \rightarrow$  Residual Standard Error

Problem with RSS:

RSS value changes if our data is changed its units/dimensions.

X	Y
1kg	100rs
2K	150rs

$\rightarrow$

X	Y
1000 grams	12 \$
2000 grams	23 \$

$\rightarrow$  (grams & \$)

Change in units bring change in RSS value to overcome that we go for TSS (Total sum of square) method

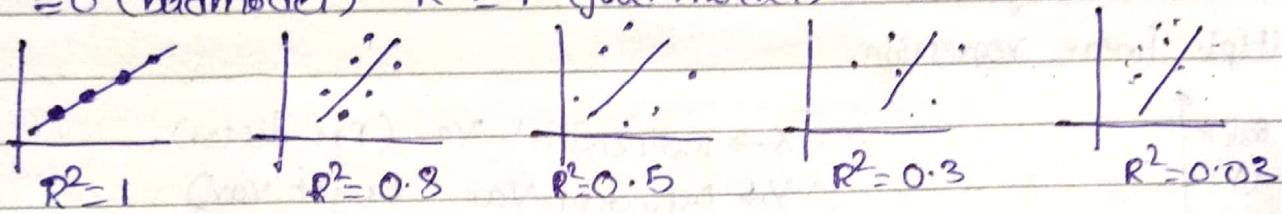
$TSS \rightarrow$  sum of the error for your data points from the mean of your response :  $TSS = \sum_{i=1}^n (y_{pred} - \bar{y}_i)^2$  (worst model)

$R^2$  (coefficient of determination)  $\rightarrow R^2 = 1 - \frac{RSS}{TSS}$

RSS is already least value , now we want TSS to be least so  $R^2$  vary from 0 to 1.

if  $RSS \approx TSS$  that means RSS equal to worst model  $\therefore R^2 = 0$  (model bad)

$R^2 \leq 0$  (bad model)  $R^2 \geq 1$  (good model)



### case study :

- ① took a dataset for a T.V selling company
- ② separated  $x$  &  $y$ , split and using Sklearn Train & Predict.
- ③ plotted residual error graph for index vs  $(y - y_{pred})$   $\rightarrow$  residual error
- ④ calculated from sklearn.metrics import mean\_squared\_error, r2\_score.

$mse = \text{mean\_square\_error}(Y_{\text{test}}, Y_{\text{Pred}})$

$r\text{-squared} = r2\text{-score}(Y_{\text{test}}, Y_{\text{Pred}})$

- mse (mean square error) means  $(1 - 0.957) \rightarrow$  our model ~~model was~~ was unable to match 95.7% values only that's nearly 8% values it can't match (good)
- $r\text{-square}(R^2\text{-score}) \rightarrow 1 - 0.59 \Rightarrow 59\%$  (actual 60%  $\Rightarrow R^2$  is good model)

### Residual square error (RSE)

$$RSE = \sqrt{\frac{RSS}{df}} \rightarrow \text{degrees of freedom}$$

$$df = n - 2$$

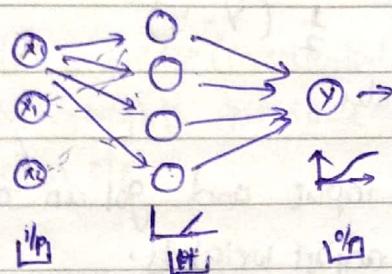
$$n \Rightarrow \text{no. of data points}$$

$RSE \rightarrow$  has some advantages & disadvantages like RSS so  $\uparrow$  good to go with  $R^2$

H Most used activation functions:

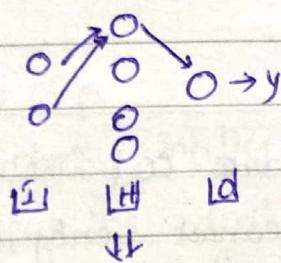
- ① Threshold function.  $y = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$
- ② Sigmoid function.  $y = \frac{1}{1+e^{-x}}$  (mostly in o/p layer)
- ③ Rectifier function  $y = \max(0, x)$
- ④ Hyperbolic tangent (tanh)  $y = \frac{1-e^{-2x}}{1+e^{-2x}}$

exercise : ①



Most cases we use Rectifier in the "hidden" layer and sigmoid at "o/p"

Why do we need a hidden layer?



$$x_1 \rightarrow 0 \rightarrow y \quad \Rightarrow \quad y = w_1 x_1 + w_2 x_2.$$

$x_2 \rightarrow 0 \rightarrow y$  has only one combination.

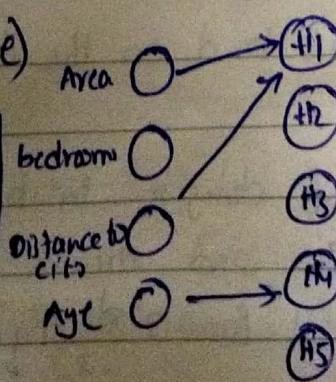
Have plenty of combinations (Y)

so each hidden layer can be specified to particular case.

\* If there is "no hidden" layer then obviously o/p will be a linear combination of inputs  $[Area, Bedrooms, Distance\ to\ city, Age] \rightarrow [Price]$

\* whereas if it has hidden layers it can have other specific combinations. (i.e)

In "H1" case it only considers "Area" & "distance" as major parameters which is general case



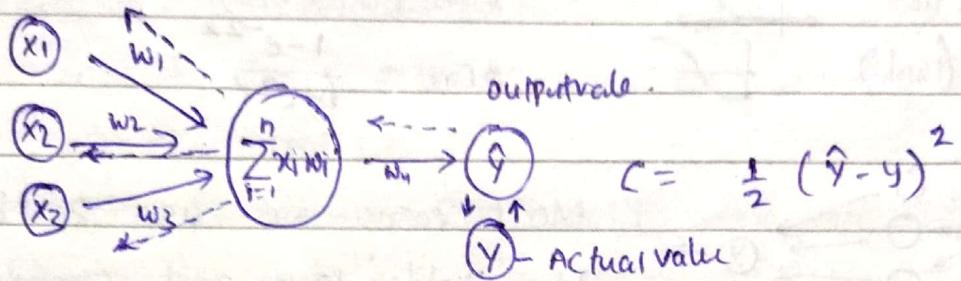
Y

In "H4" case it only considers Age as a major parameter in this case for "vintage" building it happens

$\therefore H_1 \& H_4$  the params are different though both are for house pricing: hidden layer are useful

## How do Neural networks learn?

While teaching neural network we don't teach it with hardcoded rules instead we let it to figure out the rules.

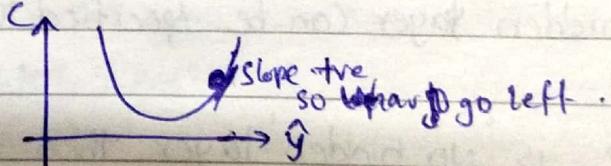
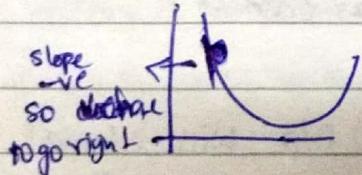


Once all the dataset is given as an input and got an o/p now you compare with actual value and adjust weights.

"epoch"  $\rightarrow$  inputting all the input's is called as one epoch. (usually)  
but for batch training it is different.

How weights are adjusted?

The primary reason for adjusting weights is to reduce cost function.  
[best optimal value for the cost function]. Simply we cannot verify all the possible values of 'w' to know optimal values. So we opted Gradient Descent method.

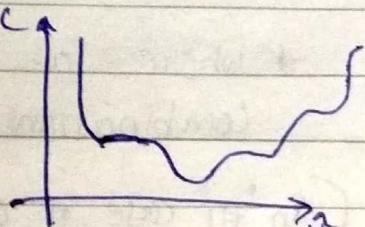


Stochastic Gradient descent?

This topic comes when cost fun is not proper convex

\* there is every change it might land on the local maxima or minima

In stochastic GD method instead of doing it in batches (all the inputs) at a time we take only inputs once and calculate  $\hat{y}$  and find cost fun and adjust weights and repeat it for 2nd inputs.



## Batch G.D

Whole batch at once  
↓  
(may lead to local minima)  
↓  
it has to load all inputs at once  
so it's slow.  
↓  
it is a deterministic (consistent)  
same starting weights give same  
accuracy in op

## Stochastic G.D

only one ip at a time.  
↓  
(only leads to global minima)

↓  
it is faster & lighter.  
(cause it don't load all the inputs)

↓  
it is stochastic (inconsistent)  
(it has random choose rows &  
weights so not possible to  
get same result every time).

Note: There is also mini batch gradient descent method.

### Training Steps

- ① Randomly initialise the weights to small numbers close to 0 but not '0'.
- ② Input the first observation of your dataset in the input layer - each feature in one input node.
- ③ Forward Propagation from Left to right neurons are activated in a way that the impact of each neuron's activation is limited by weights propagate the activations until getting the predicted results.
- ④ compare predicted result to actual results and generate error.
- ⑤ Back propagation Right to Left - update the weight according to how much they are responsible for the error the learning rate decides how much we update weight.
- ⑥ Repeat steps 1 to 5 and update the weight after each observation (Reinforcement learning) (or) Repeat 1 to 5 steps and update weight only after batch of observation (Batch learning).
- ⑦ When all the dataset get passed through it it is called an Epoch we do more & more epochs.

pd.read\_csv('...') → To read a csv file it gives a DF as o/p  
 ie = LabelEncoder() → (Categorical 10 → 1, 2, 3)  
 Oh = OneHotEncoder() → (1, 2, 3 → (0 0 0) (0 0 1))  
 sc = StandardScaler() → Standardize/normalize

In data frame we use → data.iloc[:, 1:3] → to access rows & cols  
 data frame cannot be processed directly so we convert to array from DF  
 using data.values → this converts it to array format

SKLearn::module::deletion → train-test-split (X, Y, test\_size=0.2, random\_state=1)  
 even for onehotencoder we should use Oh.fit\_transform(X).toarray()

```

import keras; from keras.models import Sequential; from keras.layers import Dense
classifier = Sequential()
classifier.add(Dense(64))

```

input\_dim = (input + output), init = "uniform", activation='relu'  
 from second layer we don't need to specify the input\_dim.  
 for final layer we use sigmoid for 1 o/p softmax for multiple o/p.

```

classifier.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

```

```

import keras
from keras.models import Sequential
from keras.layers import Dense
from sklearn.metrics import accuracy, confusion_matrix
classifier = Sequential()
classifier.add(Dense(units=6, input_dim=11, activation='relu', kernel_initializer='uniform'))
classifier.add(Dense(units=1, activation='sigmoid'))
classifier.compile(optimizer='adam', loss='binary_crossentropy', metrics=[accuracy])
y_pred = classifier.predict(x_test)
print accuracy_score(y_test, y_pred)

```

## Linear Regression.

- \* Blw Linearly Separable Variables only
- \* To Predict continuous Variable
- \* Straight line.
- \* forecasting Sales etc
- \* cost = least square error
- \*  $Y = b_0 + b_1 x + e$

## Logistic Regression.

- \* Blw linearly separable & non-separable.
- \* To Predict Categorical Variable
- \* S-curve
- \* classification etc
- \* cost
- \*  $\log \left( \frac{Y}{1-Y} \right) = C + B_1 X_1 + B_2 X_2$ .

Bias-variance tradeoff: (it is an error due to bias and variance).

generally we try to achieve "low bias" and "low variance" for a model.

(to get consistent predictions on different datasets) (to avoid underfit and overfit)

Machine learning errors → Bias error → expected error in the prediction of model.  
→ Variance error →

Expectation of error in prediction of a model)

Bias errors: are due to our wrong assumption in a model. for example

we apply linear regression to a dataset assuming that both variables have a linear relationship blw them, but they might not have linear relationship necessarily, this leads to high bias (ie. model is underfitting).

Variance error: measure of variability in the results given by model

when the dataset is changed.

BULL'S eye diagram.



innermost circle is our target

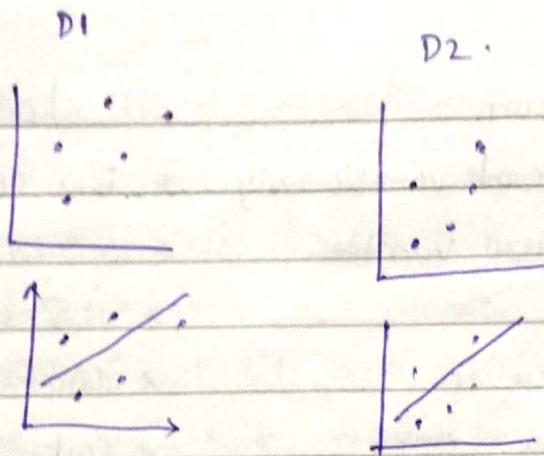
① model with low-bias, low-variance have predictions near to target

② high-bias high-variance says the model predicts vary rapidly. ✗ (model is not consistent)

③ low-variance high-bias says model is consistent but far from target ✗

④ high-bias high-variance -✗

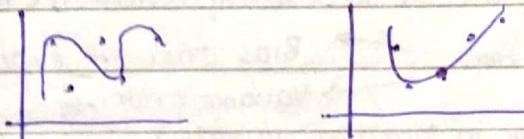
example take 2 data sets



Apply LR model on both

in this case we can say the model is having bias error because for both the data sets its <sup>Prediction</sup> ~~outp~~ is almost similar but actual o/p's of both are different therefore model is consistent but having bias error (high bias low variance)

Apply Polynomial LR model



in this case the model is having variance error because for both datasets the prediction is very different ∴ we can say that this model is inconsistent (high-variance)

### How to reduce Bias-variance tradeoff?

- Dimensionality Reduction
- Regularization in Linear models / ANN
- Using mixture model and ensemble learning
- Optimal value of K in KNN.

To check Bias-variance tradeoff we use "K-fold cross validation" (Split data set into more parts and selects random sets from them to test and train).

\* Transfer learning is a technique used to reduce the development of model from scratch to use the model & its weights by adding our desired layers to it.

Dropout is the solution for overfitting.

classifier = classifier.Dropout(0.1) →  $P = \frac{1}{10}$  → add after every layer in which you want to fire some neurons

Parameter tuning to improve best model:

Parameters such as "no of hyperparameters", "batch size", "Learnrate"

classifier = kerasClassifier(build\_fn=build\_classifier) ← in keras.wrappers.scikit\_learn  
from the definition of classifier (same function)

Parameters = {  
 'batch\_size': [25, 32],  
 'epochs': [100, 500],  
 'optimizer': ['adam', 'rmsprop']}

grid\_search = GridSearch.fit(X\_train, Y\_train)  
grid\_search.best\_params\_ ←  
best\_accuracy, best\_optimizer ←

Cross validation to check the model with different test sets:

classifier = kerasClassifier(build\_fn=build\_classifier) ← keras.wrappers.scikit\_learn  
batch\_size = 10, epochs = 100

accuracy = cross\_val\_score(estimator=classifier, X=X\_train, Y=Y\_train, cv=10, n\_jobs=-1)

accuracy.mean(), accuracy.std().

Standardisation

$$x_{\text{stand}} = \frac{x - \text{mean}(x)}{\text{standard deviation}(x)}$$

Normalisation

$$x_{\text{norm}} = \frac{x - \text{min}(x)}{\text{max}(x) - \text{min}(x)}$$

```

import numpy as np
np.random.seed(123)
from keras.models import Sequential
from keras.layers import Dense, Flatten, Accuracy, Dropout, Convolution2D, MaxPooling2D
from keras.utils import np_utils
from keras.datasets import mnist
(x_train, y_train), (x_test, y_test) = mnist.load_data()
# preprocessing input data
x_train = x_train.reshape(x_train.shape[0], 1, 28, 28)
x_test = x_test.reshape(x_test.shape[0], 1, 28, 28)
x_train = x_train.astype('float32')
x_test = x_test.astype('float32')
x_train /= 255
x_test /= 255
y_train = np_utils.to_categorical(y_train, 10)
y_test = np_utils.to_categorical(y_test, 10)
# model architecture.
model = Sequential()
model.add(Convolution2D(32, 3, 3, activation='relu', input_shape=(1, 28, 28)))
model.add(Convolution2D(32, 3, 3, activation='relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
model.add(Flatten())
model.add(Dense(128, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(10, activation='softmax'))
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
model.fit(x_train, y_train, batch_size=32, nb_epoch=10)
model.predict(x_test) == model.eval(x_test, y_test, verbose=20)

```

binary-crossentropy  $\leq^0$   
 categorical-crossentropy  $\leq^1$   
 mse  $\rightarrow$  for regression problem "contains"

```

from keras.models import Sequential
from keras.layers import Dense, Flatten, Dropout, Conv2D, Maxpool2D
model = Sequential()
model.add(Conv2D(64, input_shape=(120, 120, 3), kernel_size=(3, 3), strides=(1, 1),
                activation='relu'))
model.add(Conv2D(32, input_shape, kernel_size=(3, 3), strides=(1, 1), activation='relu'))
model.add(Maxpool2D(strides=2, pool_size=(2, 2)))
model.add(Flatten())
model.add(Dense(units=16, activation='relu', kernel_initializer='uniform'))
model.add(Dense(units=8, activation='softmax'))
model.add(Dense(units=2, activation='softmax'))
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
model.fit(x_train, y_train, batch_size=10, epochs=10)
model.batch_size(x_test, y_test)
(convolution) → (maxpooling) → (flattening) → (full connection)

```

Max pooling: (spatial invariance) means features what it is looking for need not to be in the same position or size (ratio) they can vary a bit. There are maxpooling, mean pooling, sumpooling, average pooling - -

ReLU layer: (rectilinear operation) ReLU increases non linearity. because basically images are non linear (at edges/borders) during convolution it might turn into linear so that ~~ReLU~~ <sup>is needed</sup> will tag

Convolution: (Yann Lecun)

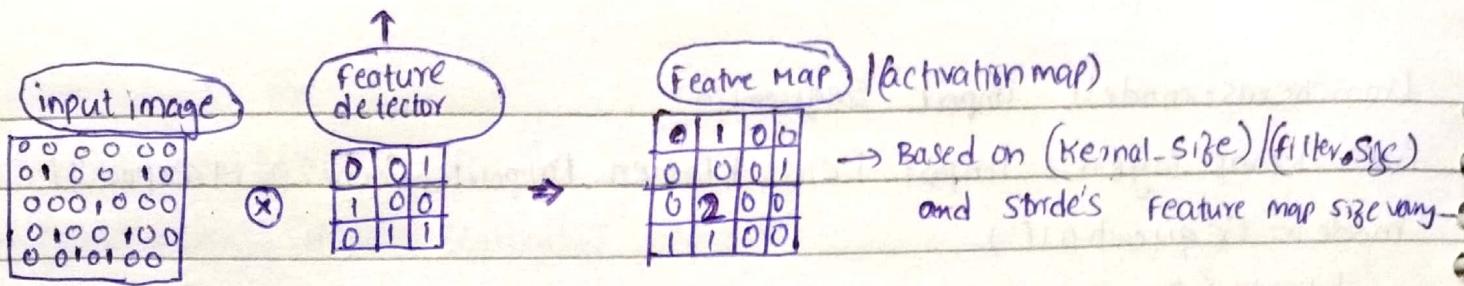
Flattening: is rearranging the pooled featuremaps in single array that goes into N:N

Input layer: - flattened pooled featuremaps are fed to input layer.

Fully connected: - (hidden layers) (must be fully connected) to input layer

Output layer: each for each class → connected to fully connected layer.

Machine will be trying to learn these (feature detectors) / (Filters) / (Kernel)



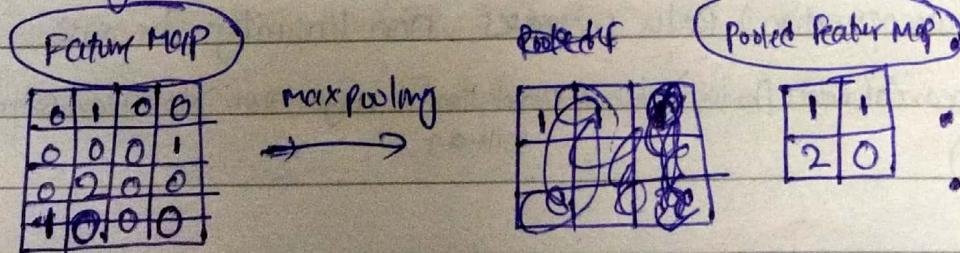
We took an image and we took one random feature detector / filter and we convoluted both then we got convolved features called "feature map"

\* In this process based on stride size & kernel size we can reduce feature map size

\* What exactly is a feature detector?

- Feature detector is a pattern / shape which we are looking for in a image.
- Feature map will have highest values if the specific portion of the image and feature detector match.
- Along the training we are making it to learn these feature detectors in terms of weights. i.e  $(64, 32 \rightarrow \text{in conv2D}(64))$  This can look for 64 different patterns (feature detectors) in a image.  
Kernel size = (3,3)  $\rightarrow$   $\therefore$   $64 \times 3 \times 3$  filters  
 $\therefore (64 \times 3 \times 3)$  weights learned.

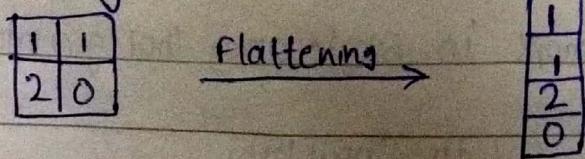
Max Pooling



size got reduced

- Improving spatial invariance
- Reducing no. of features to look for.

Flattening



## Softmax & Cross entropy:

softmax takes the values and normalize them so that their sum is equal to '1' (softmax is a logistic function  $f_j(z) = \frac{e^{z_j}}{\sum_k e^{z_k}}$ )

$$\text{Cross entropy} \rightarrow L = -\log \left( \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right) \rightarrow H(p, q) = -\sum_x p(x) \log q(x)$$

linear  $\rightarrow$  mse  $\rightarrow$  3.14

softmax  $\rightarrow$  categorical-crossentropy  $\rightarrow [0 \ 0 \ 1 \ 0]$

relu/sigmoid  $\rightarrow$  binary-crossentropy  $\rightarrow 0$

image augmentation   
`keras.preprocessing.image import ImageDataGenerator`  
generates images in different scales from the given

Supervised	Artificial Neural Networks Convolution " " Recurrent "	used for regression & classification used for computer vision used for time series Analysis
unsupervised	Self organizing maps Deep Boltzmann Machines Auto encoders	used for feature detection used for recommendation systems used for recommendation systems

"Cerebrum , cerebellum , Brainstem

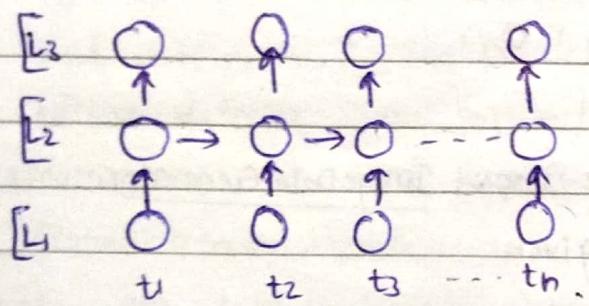
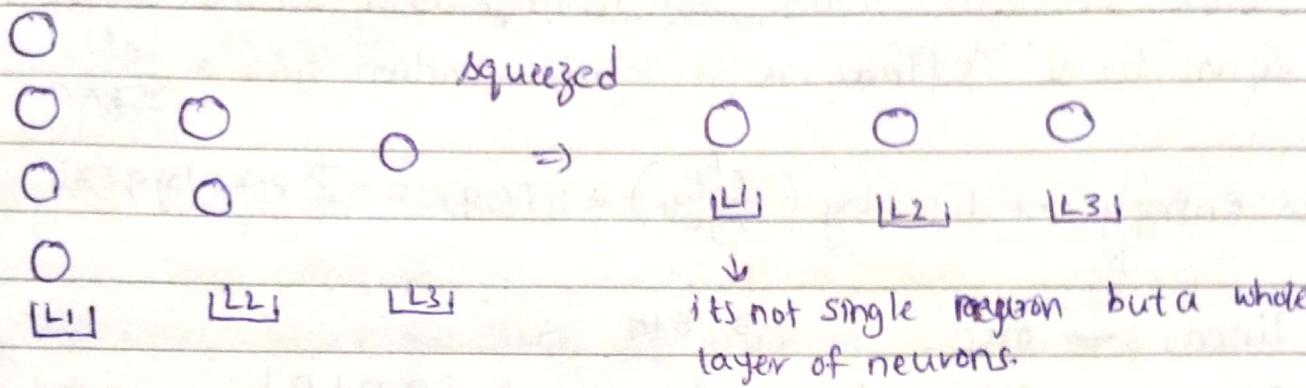
Frontal lobe  $\leftarrow$  (Temporary memory) RNN (Personality, behaviour, working memory)

Parietal lobe  $\leftarrow$  (sensation & perception to create a 3D world in mind)

Temporal lobe  $\leftarrow$  weights of ANN (responsible for long term memory)

Occipital lobe  $\leftarrow$  CNN (vision)

## Recurring neural network (RNN):

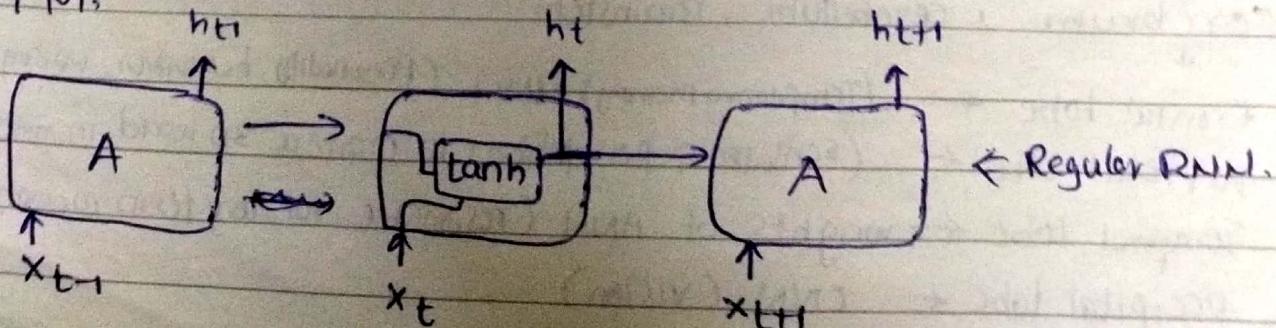


every neuron represents a layer of neurons, and  $t_1, t_2, t_3$  are different time's.

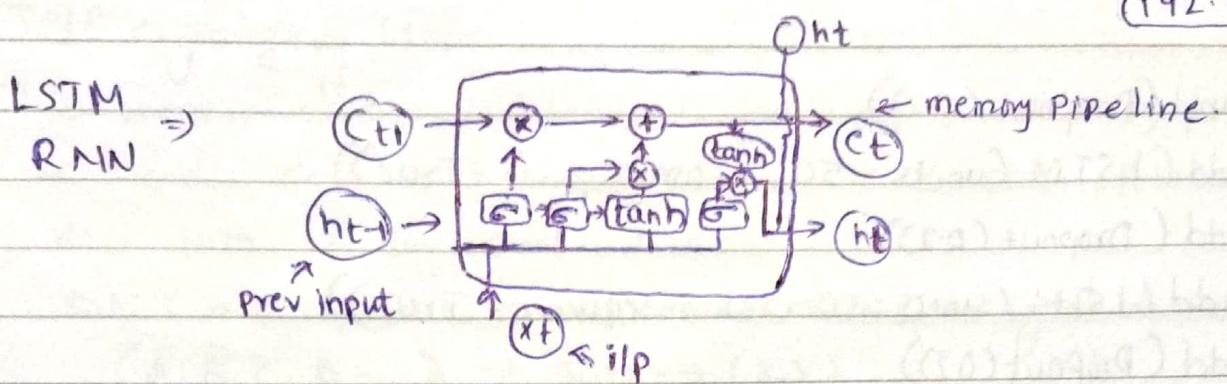
In RNN while calculating weights (Adjusting weights) we will face a problem called vanishing / exploding gradient solutions:

- Truncated Back propagation,
  - Penalties
  - gradient clipping
  - weight initialization
  - Echo state Networks
  - long short term memory networks
- } gradient exploding ( $w_{rec} > 1$ )
- } vanishing gradient ( $w_{rec} < 1$ )

LSTM:



bag  
hefi  
•pcap file  
(192.168.115.140)



every line in this architecture is a vector  $\rightarrow$   
concatenate  $\Rightarrow$  copy  $\Rightarrow$  O pointwise operation. (like valve)

implementation of RNN:

```
from sklearn.preprocessing import MinMaxScaler
```

```
sc = MinMaxScaler ( feature_range = (0,1) )
```

```
training_set_scaled = sc.fit_transform(training_set)
```

} normalized

#creating a data structure with 60 timesteps and 1 output.  
(it will consider previous 60 days stock price and predict the next).

```
x_train = []
```

```
y_train = []
```

```
for i in range(60, len(training_set)):
```

```
    x_train.append ( training_set_scaled[ i-60 : i , 0 ] )
```

```
    y_train.append ( training_set_scaled[ i , 0 ] )
```

```
x_train, y_train = np.array (x_train), np.array (y_train)
```

```
x_train = np.reshape(x_train, ( x_train.shape[0] , 60 , 1 ))
```

# Building the RNN

↑  
no of ilp rows

↑  
ilp cols

← shape expected by RNN

```
from keras.models import Sequential
```

```
from keras.layers import Dense, LSTM, Dropout
```

```
regressor = Sequential() ← we call it regression as we are predicting.
```

```
regressor.add(LSTM (units = 50, return_sequences = True, input_shape = (60,1)))
```

↑  
no of units of LSTM cells

↑  
Return sequen LSTMs T/F

↑  
input shape

```
regressor.add(Dropout(0.2))
```

```
regressor.add(LSTM(units=50, returnSequence=True))
```

```
regressor.add(Dropout(0.2))
```

```
regressor.add(LSTM(units=50, returnSequence=True))
```

```
regressor.add(Dropout(0.2))
```

```
regressor.add(LSTM(units=50, returnSequence=False))
```

```
regressor.add(Dropout(0.2))
```

```
regressor.add(Dense(units=1))
```

```
regressor.compile(optimizer='adam', loss='mse')
```

```
regressor.fit(x_train, y_train, batch_size=32, epochs=100)
```

# Testing to predict 2017 Jan 1 data. (we have 20 days of testing data from Jan 1 to Jan 20) To predict a day we need previous 60 days data so that, along with test data 2017 Jan(1-20) we also take data in 2016 Dec & Sep. To do this we will add all the original data from train & test and then we will normalize from that

```
dataset_total = pd.concat([dataset_train['open'], dataset_test['open']], axis=0)
```

```
inputs = dataset_total[-len(dataset_total) - len(dataset_test) - 60:-1].values
```

```
inputs = inputs.reshape(-1, 1)
```

```
inputs = sc.transform(inputs) → must not fit (if we fit prev values might change)
```

```
x_test = []
```

```
for i in range(60, len(inputs)):
```

```
x_test.append(inputs[i-60:i, 0])
```

```
x_test = np.array(x_test)
```

```
x_test = np.reshape(x_test, (x_test.shape[0], 60, 1))
```

```
pred = regressor.predict(x_test)
```

```
Pred = sc.inverse_transform(pred) → To convert into actual values
```

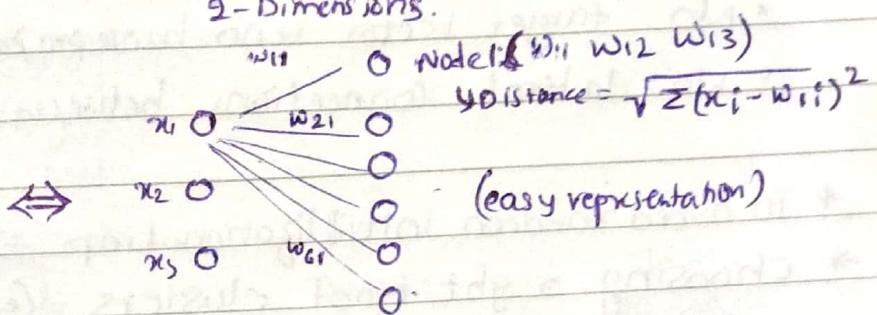
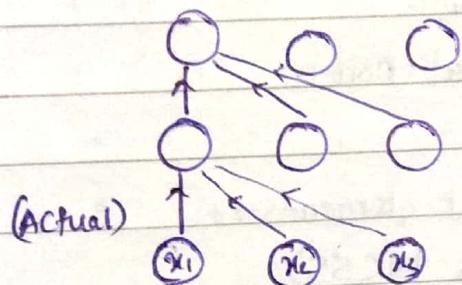
# Now plot the actual and predicted values to see the predictions accuracy

## Self organizing Maps:

This comes under unsupervised learning.

we provide data with 'n' nof columns it constrain that entire data into '2' dimensional data.

SOM's are useful to reduce the dimensionality of dataset



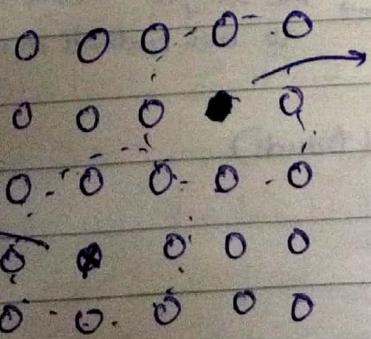
SOM are bit different for ANN.

SOM are bit different.  
In ANN weights are multiplied by inputs and activation function decides whether to fire the node or not; whereas in SOM there is no activation fun. Weights are characteristics of node itself.  
• Here weights acts as inputs to coordinates in input space.

(First random weights are given to entire network)

- we calculate euclidian distance =  $\sqrt{\sum (x_i - w_{(B\mu)})^2}$ , and we will choose the least distance node as Best matching unit (BMU)
  - Now weights of (BMU) node will change  $\rightarrow$  (in a way to decrease euclidian dist)
  - The a whole range of node's around (BMU) node's weights got updated.

This is BMU for  
(A<sub>1</sub>, A<sub>2</sub>, A<sub>3</sub><sup>BMU</sup>, A<sub>5</sub>)  
inputs  
so its weights got  
updated.



BMU for (B<sub>1</sub> B<sub>2</sub> B<sub>3</sub> B<sub>4</sub> B<sub>5</sub>) input.  
BMU weights are updated and  
other nodes around it weights got  
updated based on the distance  
from BMU.

and with in certain range all the node weight got updated.

- Based on more iterations the region around BMW getting updated will shrink.
- SOM retain topology of the input set
- SOM reveal correlations that are not easily identified.
- SOM classify data without supervision.
- No target vector  $\rightarrow$  no backpropagation.
- No lateral connections between output nodes.

\* To avoid random initialization trap we use kmeans++  
 \* choosing right noof clusters  $\Rightarrow$  (metrics WCSS)

$$WCSS = \sum_{P_i \in \text{cluster}(i)} \text{distance}(P_i, c_i)^2$$

The method we use is called elbow method.

```
from sklearn.preprocessing import MinMaxScaler; from minisom import MiniSom,
```

```
Som = MiniSom(x=5, y=5, input_len = data.shape[1])
```

```
Som.random_weights_init(data)
```

```
Som.train_random(data, num_iteration = 1000)
```

```
from pylab import bone, pcolor, colorbar, plot, show
```

`bone()`  $\rightarrow$  white chart

`pcolor(Som.distance_map().T)`  $\rightarrow$  Colour map of  $(5 \times 5)$  based on distance of each unit

`colorbar()`  $\rightarrow$  color bar

`mappings = Som.win_map(data)`  $\rightarrow$  This will give members in data that are mapped to each coordinate in map  $(0,0) (0,1) \dots (4,4)$

`frauds = mappings[(2,2)]`

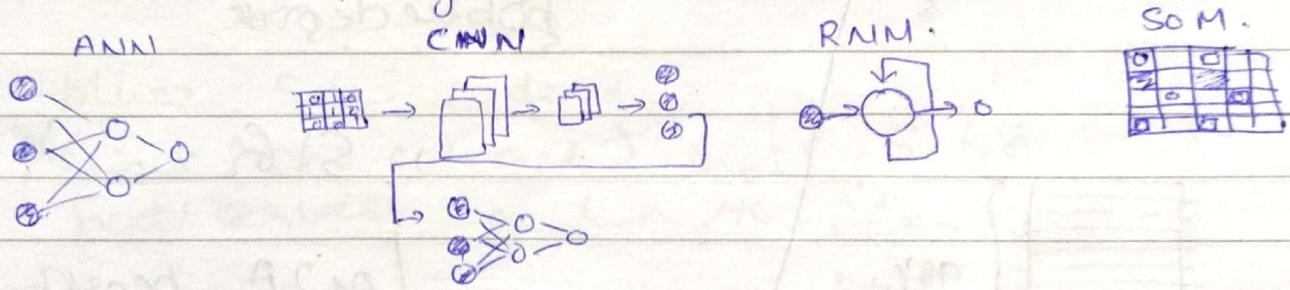
`frauds = m[m.inverse_transform(frauds)]`

## Boltzmann Machines.

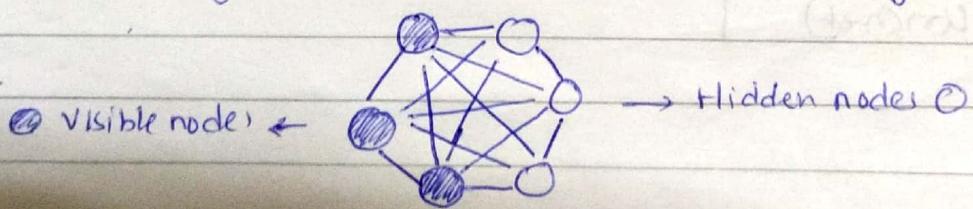
Boltzmann Machines can be seen in two different point of views

1. An energy based model (EBM)
2. A probabilistic Graphic Model. (

- (EBM) energy based model.
- (RBM) Restricted boltzmann machines
- (CD) contrastive divergence.
- (DBN) deep Belief Networks
- (DBM) deep boltzmann machines



All the above (ANN, CNN, RNN, SOM) have a direction in which they grow whereas boltzmann machines doesn't have any directions



- Boltzmann machines have only input/ <sup>visible</sup> hidden layers but no output layer.
- everything is connected to everything.
- There is no direction of flow

[in boltzmann machine even input/visible nodes are connected to each other because boltzmann machines are fundamentally different from other machines. They don't expect data instead they generate data therefore boltzmann machines are generative machines].

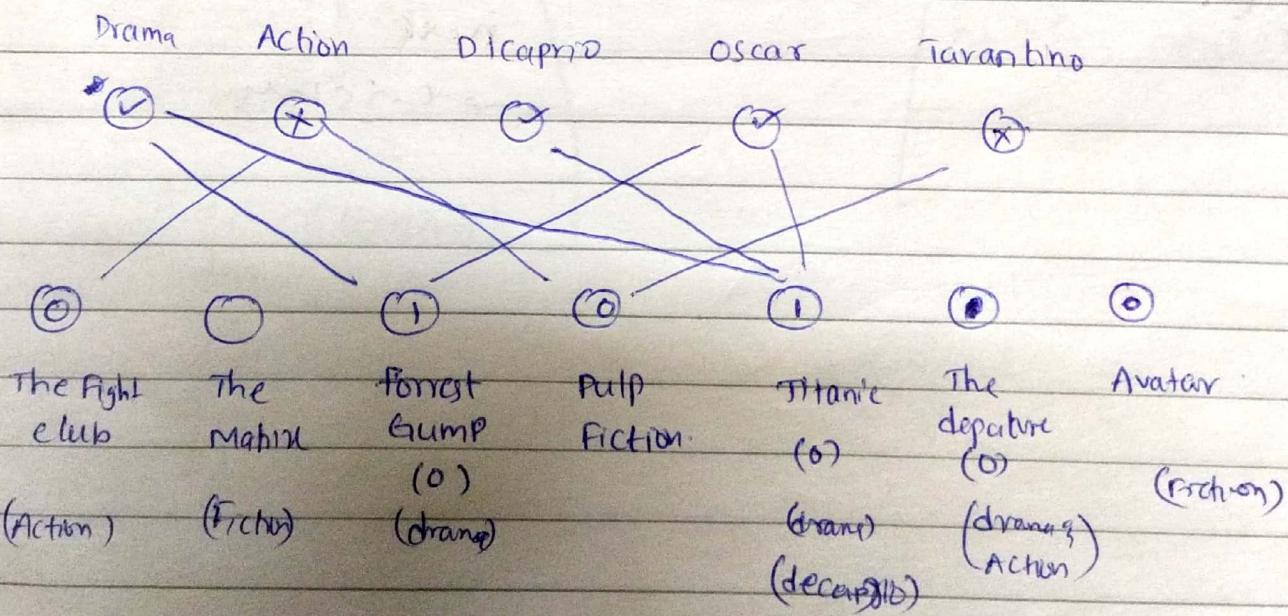
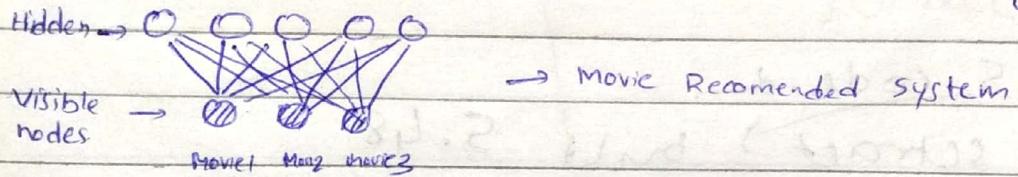
$$P_j = \frac{e^{-E_j/kT}}{\sum_{i=1}^M e^{-E_i/kT}} \rightarrow \text{Boltzmann equation.}$$

$P \rightarrow$  probability is inversely proportional to energy  $E$ .

$$E(v, h) = -\sum_i a_i v_i - \sum_j b_j h_j - \sum_i \sum_j w_{ij} v_i h_j \rightarrow \text{energy state.}$$

$$P(v, h) = \frac{1}{Z} e^{-E(v, h)} \rightarrow \text{probability of each energy state}$$

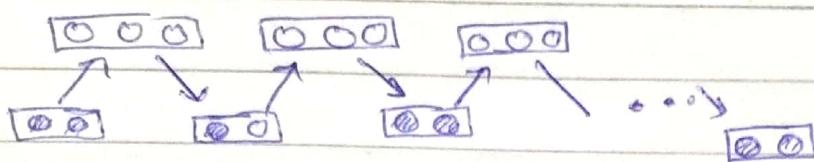
Due to implementation difficulties in fully connected Boltzmann machines they decided to implement (RBM) restricted boltzmann machines



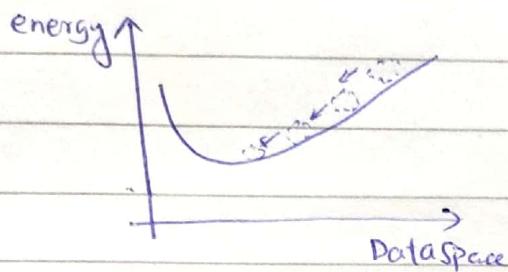
Contrastive Divergence

feeding input values to network ie we will generate input values again (which are slightly vary from actual) again those generated values are feeded into network and regenerated and again. it is repeated.

## Contrast Divergence 1



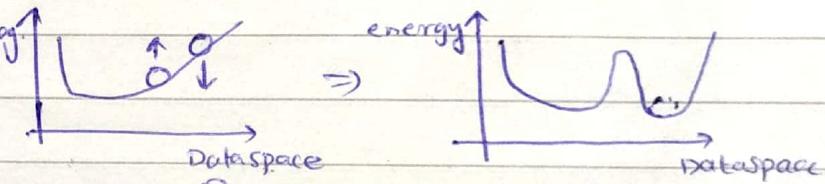
(unit inputs modify no further)



$$\frac{\partial \log P(v^0)}{\partial p_{0j}} = \langle v_i^0 h_j^0 \rangle - \langle v_i^\infty h_j^\infty \rangle$$

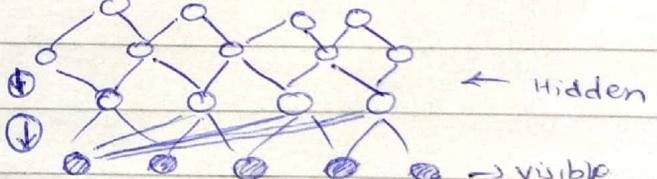
a system always tries to end up in lowest energy state  
weights are actually deciding the shape of the curve

Basically this method of generating inputs and feeding back to system take long time so geofry hinton found a shortcut which is contrastive divergence.



Deep Belief Network (DBN)

Stack of (RBM) on one another



implementation of Boltzmann Machines: (RBM)

```

import numpy as np
import pandas as pd
import torch
import torch.nn as nn
import torch.nn.parallel
import torch.optim as optim
import torch.utils.data
from torch.autograd import Variable
    
```

# Torch have tensors like tensorflow