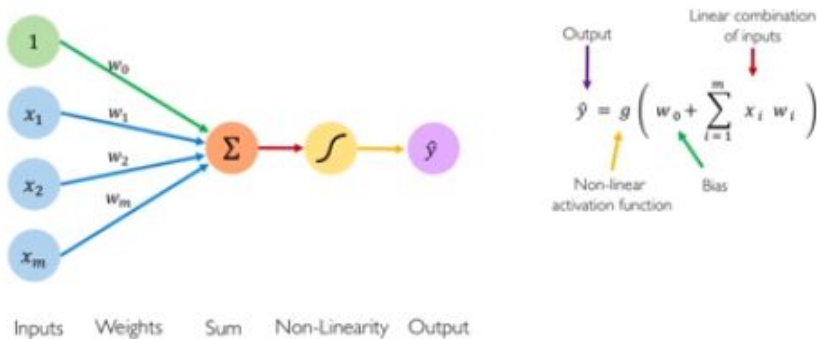Learning from the things are working later ~ learning features from underlying features.

Big Data -> Hardware -> Software

Perceptron :

## The Perceptron: Forward Propagation
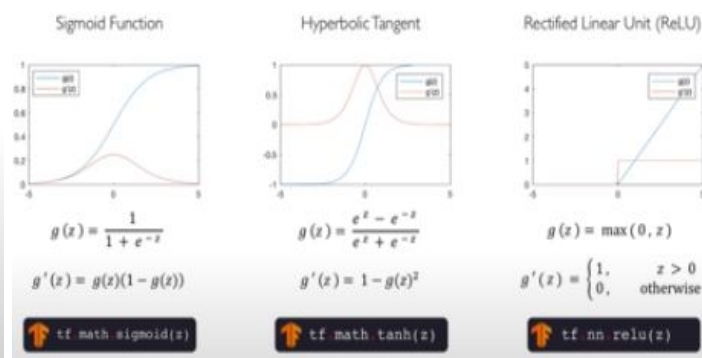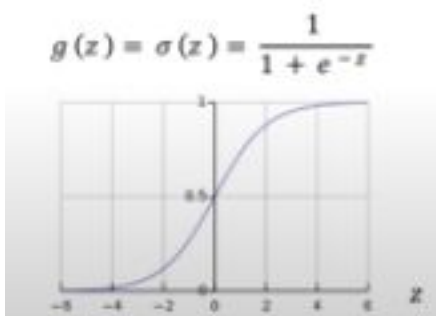


$$\hat{y} = g\left(w_0 + \sum_{i=1}^{m} x_i\, w_i\right)$$

Output — Linear combination of inputs

Non-linear activation function — Bias

Inputs    Weights    Sum    Non-Linearity    Output

Activation Function :

Takes a real number in x axis and converts it into a number in 0 and 1.

$$\hat{y} = g\left(w_0 + X^T W\right)$$

## Common Activation Functions

- Example: sigmoid function

$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



**Sigmoid Function**

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

tf.math.sigmoid(z)

**Hyperbolic Tangent**

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

tf.math.tanh(z)

**Rectified Linear Unit (ReLU)**

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

tf.nn.relu(z)

Purpose of AF is - to introduce non-linearities in functions.

Hidden layers , dense layer

Quantifying Loss : measures the cost incurred from incorrect predictions.
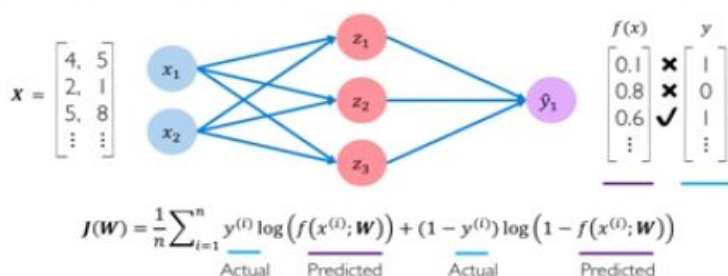
$$\mathcal{L}\left(\underbrace{f\left(x^{(i)};W\right)}_{\text{Predicted}},\underbrace{y^{(i)}}_{\text{Actual}}\right)$$

Empirical Loss : total loss over the entire dataset.

$$J(W) = \frac{1}{n}\sum_{i=1}^{n}\mathcal{L}\left(\underbrace{f\left(x^{(i)};W\right)}_{\text{Predicted}},\underbrace{y^{(i)}}_{\text{Actual}}\right)$$
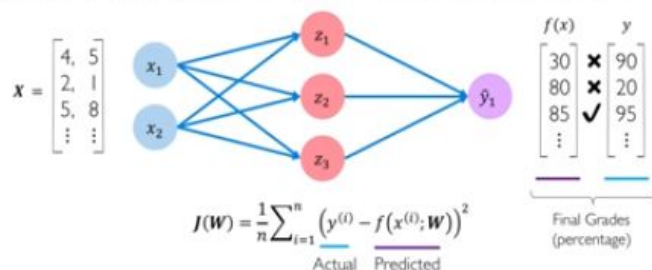
## Binary Cross Entropy Loss

*Cross entropy loss* can be used with models that output a probability between 0 and 1



$$J(W) = \frac{1}{n}\sum_{i=1}^{n}\underbrace{y^{(i)}}_{\text{Actual}}\log\left(\underbrace{f(x^{(i)};W)}_{\text{Predicted}}\right) + \underbrace{(1-y^{(i)})}_{\text{Actual}}\log\left(\underbrace{1-f(x^{(i)};W)}_{\text{Predicted}}\right)$$

```
loss = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_logits(y, predicted) )
```

## Mean Squared Error Loss

*Mean squared error loss* can be used with regression models that output continuous real numbers



$$J(W) = \frac{1}{n}\sum_{i=1}^{n}\left(\underbrace{y^{(i)}}_{\text{Actual}} - \underbrace{f(x^{(i)};W)}_{\text{Predicted}}\right)^2$$

Final Grades (percentage)

```
loss = tf.reduce_mean( tf.square(tf.subtract(y, predicted)) )
```

Loss Optimisation : finding the network weights that achieve the lowest loss through maybe gradient descent.

## Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3.     Compute gradient, $\frac{\partial J(W)}{\partial W}$

4.     Update weights, $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$

5. Return weights

```python
import tensorflow as tf

weights = tf.Variable([tf.random.normal()])

while True:      # loop forever
    with tf.GradientTape() as g:
        loss = compute_loss(weights)
        gradient = g.gradient(loss, weights)

    weights = weights - lr * gradient
```

Backpropagation : chain rule ~

## Optimization through gradient descent

$$W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$$

How can we set the
learning rate?

Adaptive Learning Rates ~

- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
    - how large gradient is
    - how fast learning is happening
    - size of particular weights
    - etc...

# Gradient Descent Algorithms

| Algorithm | TF Implementation | Refe |
|---|---|---|
| • SGD | tf.keras.optimizers.SGD | Kiefer & Wolfowitz. "Sto Maximum of a Regressi |
| • Adam | tf.keras.optimizers.Adam | Kingma et al. "Adam: A I Optimization." 2014. |
| • Adadelta | tf.keras.optimizers.Adadelta | Zeiler et al. "ADADELT; Method." 2012. |
| • Adagrad | tf.keras.optimizers.Adagrad | Duchi et al. "Adaptive S Learning and Stochastic |
| • RMSProp | tf.keras.optimizers.RMSProp | |

```
import tensorflow as tf
model = tf.keras.Sequential([...])

# pick your favorite optimizer
optimizer = tf.keras.optimizer.SGD()

while True: # loop forever

    # forward pass through the network
    prediction = model(x)

    with tf.GradientTape() as tape:
        # compute the loss
        loss = compute_loss(y, prediction)

    # update the weights using the gradient
    grads = tape.gradient(loss, model.trainable_variables)
    optimizer.apply_gradients(zip(grads, model.trainable_variables)))
```

Stochastic gradient descent ~ means using a single point.

Using a single point for gradient descent in very noisy and using all the data points is very computationally expensive , therefore the solution is ~

Using mini batches of points ~ and true is calculated by taking the average , therefore smoother convergence and allows greater learning rates , enables parallelism therefore faster.
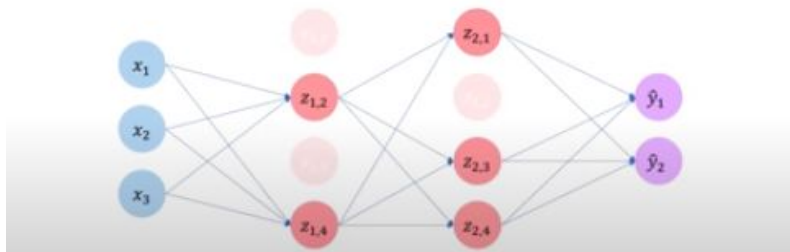
Overfitting exists so to deal with it ~

Regularisation ~ technique that constraints our optimisation problem to discourage complex models ~ improve generalisation.

## Regularization 1: Dropout

- During training, randomly set some activations to 0
  - Typically 'drop' 50% of activations in layer
  - Forces network to not rely on any 1 node

## Regularization 2: Early Stopping

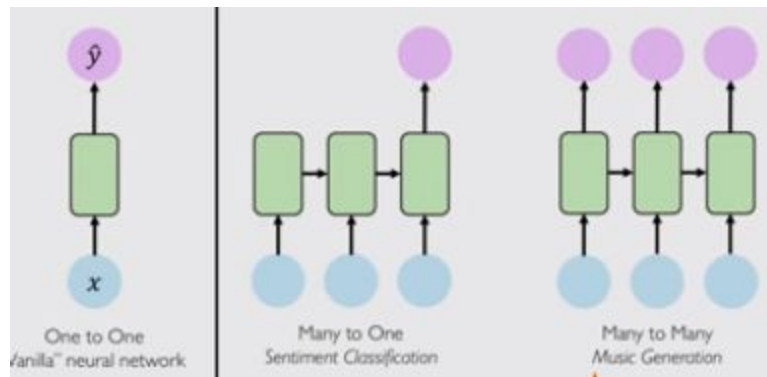- Stop training before we have a chance to overfit



---

# RECURRENT NEURAL NETWORKS

Sequential Modelling : predict the next thing that will happen after a previous states are given
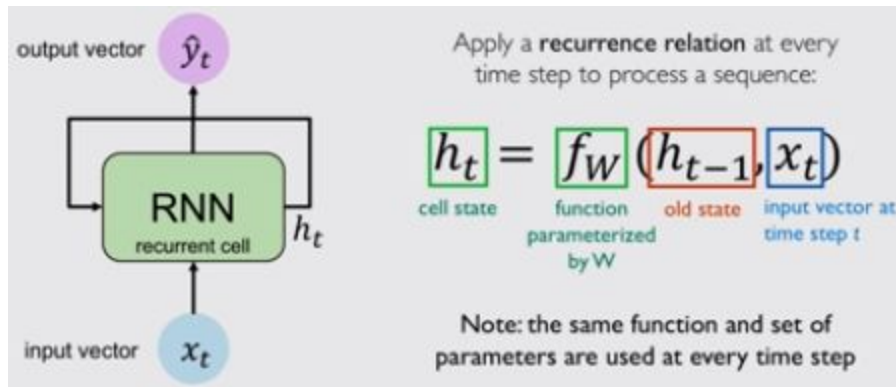
If word prediction then if using BagOfWords then the order is not reserved therefore Separate parameters should be used to encode the sentence.

In standard Feed Forward -> it goes one to one passing of data

RNN for Sequence Modeling :



RNNs have a loop in them that enables them to maintain an internal state , not like a vanilla NN where the final output is only present.
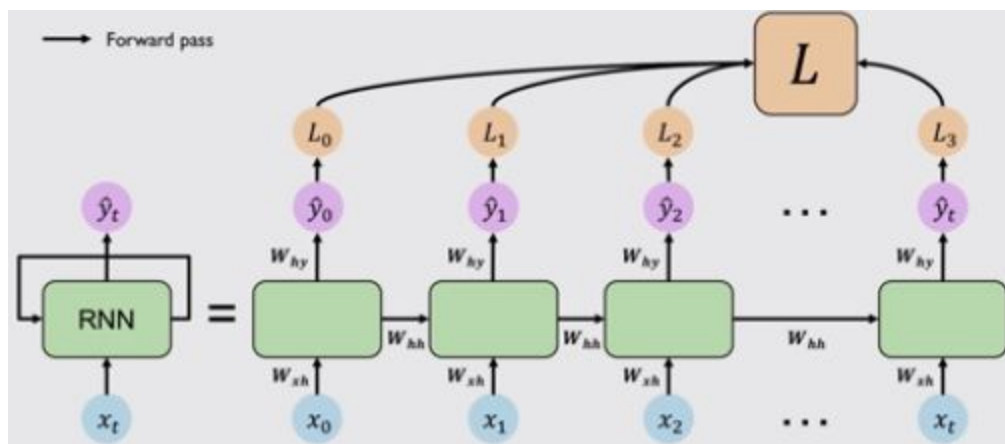
output vector $\hat{y}_t$

RNN recurrent cell $h_t$

input vector $x_t$

Apply a **recurrence relation** at every time step to process a sequence:

$$h_t = f_W(h_{t-1}, x_t)$$

cell state    function parameterized by W    old state    input vector at time step t

Note: the same function and set of parameters are used at every time step

**Output Vector**
$$\hat{y}_t = W_{hy}^T h_t$$

**Update Hidden State**
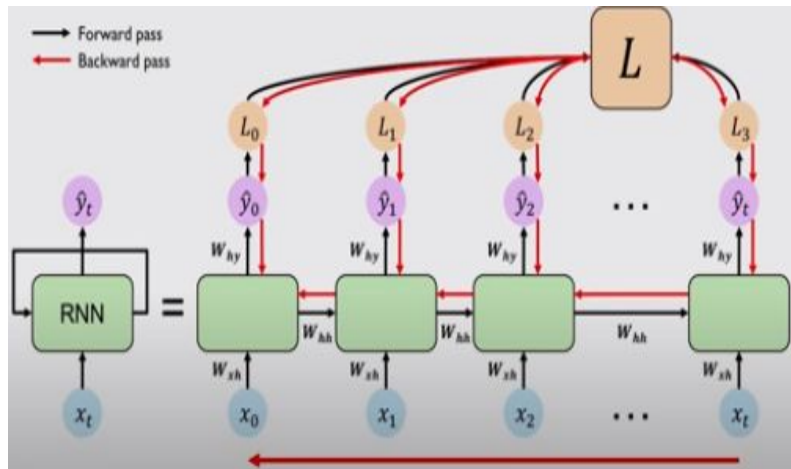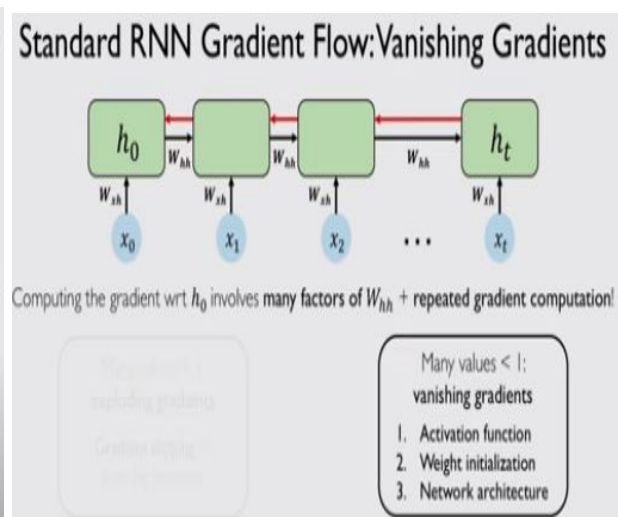$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$
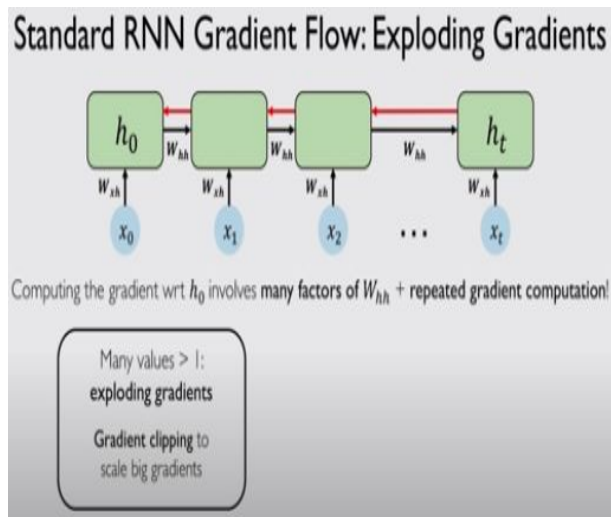
**Input Vector**
$$x_t$$



```
tf.keras.layers.SimpleRNN(rnn_units)
```

Backpropagation through time :

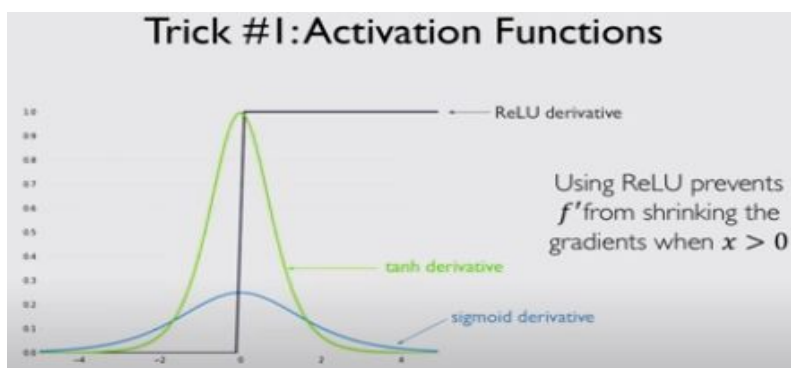In RNNs at each time backprop is done also in total ~

Standard RNN Gradient Flow can have 2 problems :



The problem of LongTermDependencies (vanishing gradient)~

Multiply many small numbers together -> errors which are in further back time steps keep having smaller and smaller gradients -> this bias es parameters to capture dependencies in models , thus standard RNNs becoming less capable.

To solve the above problem :

Both tanh and sigmoid have derivatives less than 1.
Therefore use RELU but x>0 only then.

## Trick #2: Parameter Initialization

Initialize **weights** to identity matrix

Initialize **biases** to zero

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

This helps prevent the weights from shrinking to zero.

## Solution #3: Gated Cells

Idea: use a more **complex recurrent unit with gates** to control what information is passed through
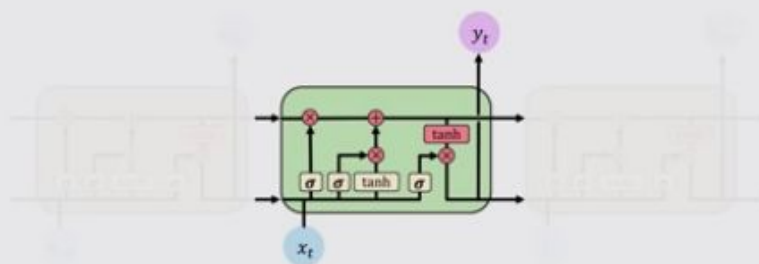
gated cell
LSTM, GRU, etc.

**Long Short Term Memory (LSTMs)** networks rely on a gated cell to track information throughout many time steps.

The above is well suited for learning tasks.

Long Short Term Memory (LSTM) Networks :

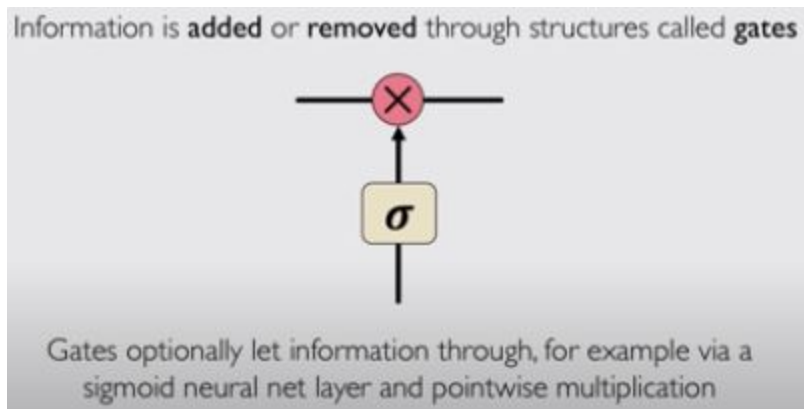LSTM modules contain **computational blocks** that **control information flow**



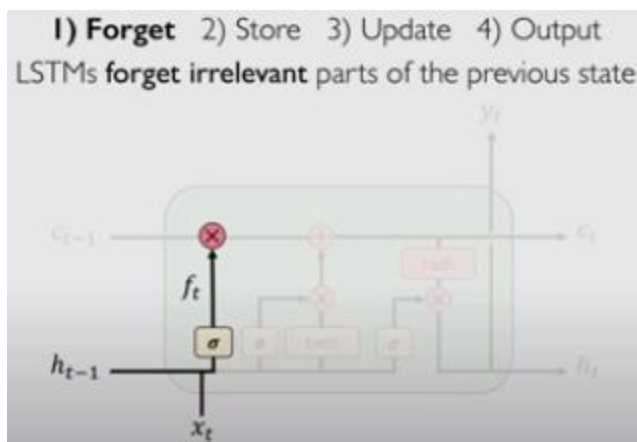LSTM cells are able to track information throughout many timesteps
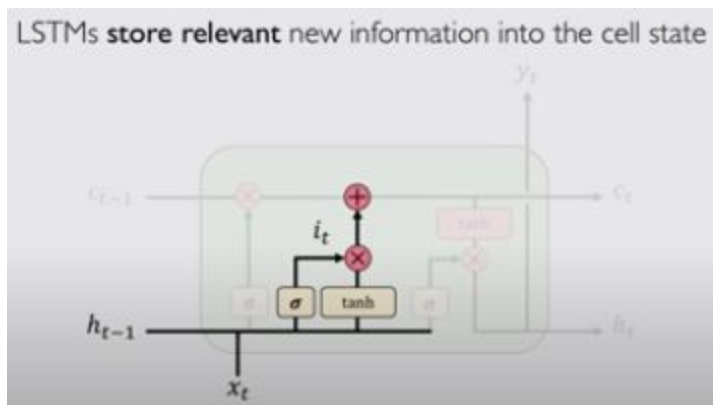tf.keras.layers.LSTM(num_units)

Gates are an important part in LSTMs

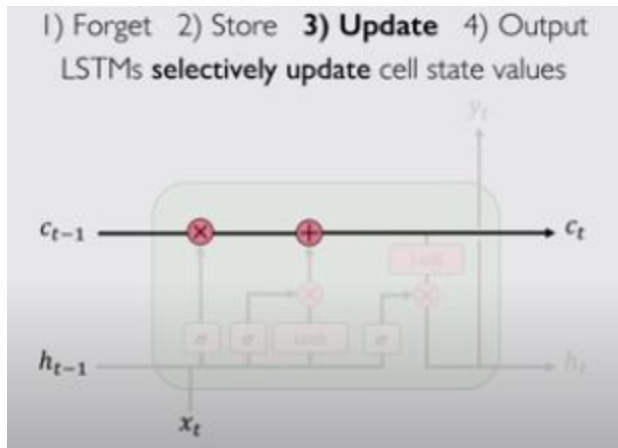Information is **added** or **removed** through structures called **gates**



Gates optionally let information through, for example via a sigmoid neural net layer and pointwise multiplication

LSTMs work in 4 steps :

**1) Forget**  2) Store  3) Update  4) Output
LSTMs **forget irrelevant** parts of the previous state



1.

LSTMs **store relevant** new information into the cell state



2.

1) Forget  2) Store  **3) Update**  4) Output
LSTMs **selectively update** cell state values



3.

1) Forget  2) Store  3) Update  **4) Output**
The **output gate** controls what information is sent to the next time step
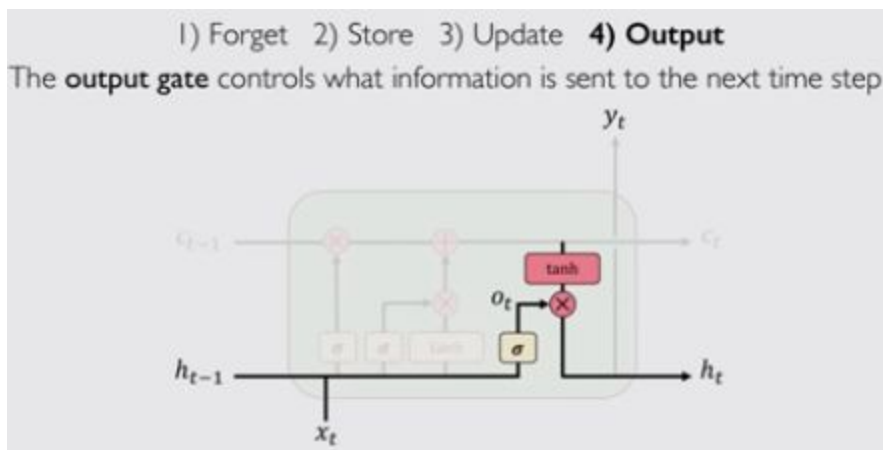


4.

All these allow LSTMs to have an uninterrupted gradient flow.

RNN Applications :

Example :

Music Generation ~ at each level the next tune is predicted.
Sentiment Classification ~
Machine Translation ~ encoding bottleneck is a problem so attention mechanism is used ~ all states of time steps are accessed and trained upon individually.
Trajectory Prediction ~ Self Driving Car
Environmental Modelling