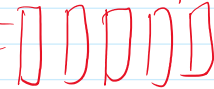


# Vanishing Gradient Problem

Thursday, April 7, 2022 7:45 AM

In machine learning, the **vanishing gradient problem** is encountered when training artificial neural networks with gradient-based learning methods and backpropagation. In such methods, during each iteration of training each of the neural network's weights receives an update proportional to the partial derivative of the error function with respect to the current weight. The problem is that in some cases, the gradient will be vanishingly small, effectively preventing the weight from changing its value. In the worst case, this may completely stop the neural network from further training. ~~As one example of the problem case, traditional addition~~

1)  $0.1 \times 0.1 \times 0.1 \times 0.1 = 0.0001 \rightarrow \text{VGP}$

2) Deep NN  $\rightarrow$  

3) Sigmoid / tanh  $\rightarrow$  Af  $\rightarrow$  deep

Backprop  $\rightarrow$   $\hat{y} - y = L$

$W_n = W_0 - \eta \left[ \frac{\partial L}{\partial W} \right]$  derivative of L wrt weight

$W_0 = 1$  small change

$W_n = 1 - 0.01 \times 0.0001$

$W_n = 0.99999$  vanishingly small

did not change

$\frac{\partial L}{\partial W_{11}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial W_{11}}$

$0 < x < 1$   $0 < x < 1$   $0.1$

$0.0001$

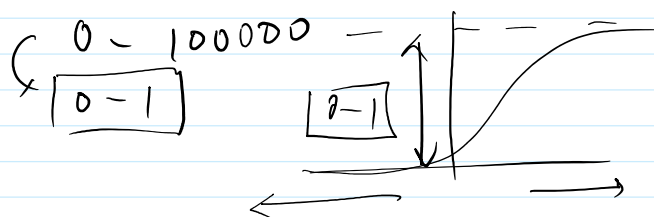
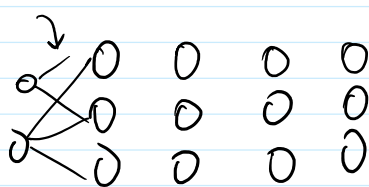
$\frac{\partial L}{\partial W_{11}} = 0.000001$

Backprop 2

model training  $\rightarrow$  x

80's, 90's 2 main reasons

DL failed



How to recognize?

- 1) Loss focus  $\rightarrow$  epoch  $\rightarrow$  no changes  $\rightarrow$  VGP
  - 2) weights  $\rightarrow$  graph value
- $W_{11}$
- Keras
- loss after epoch
- epoch
- Tensorboard
- callbacks

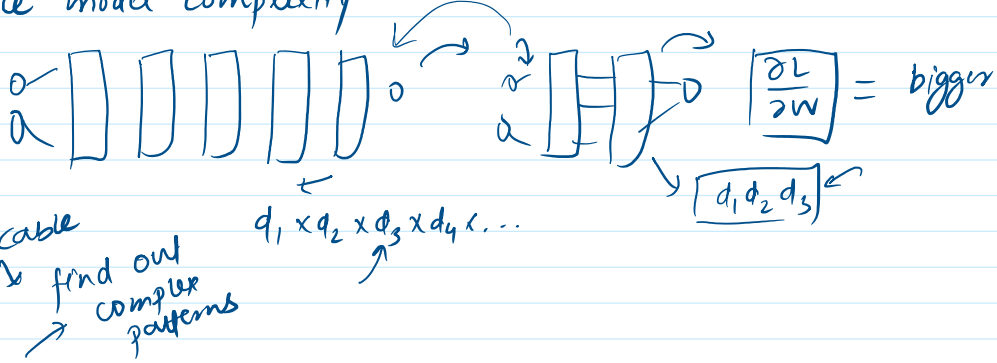
$W_n = W_0 - \eta \left[ \frac{\partial L}{\partial W} \right]$

$\frac{W_0 - W_n}{\eta} =$

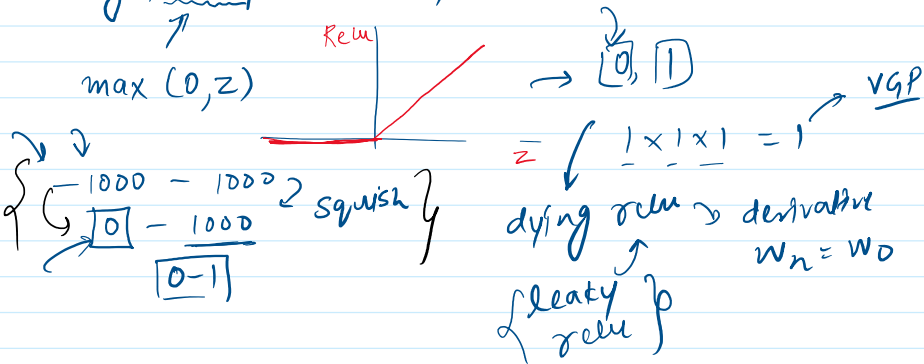
$$\left[ \frac{w_0 - w_n}{n} \right] =$$

## How to handle Vanishing Gradient Problem →

1) Reduce model complexity



2) Using ReLU Activation functions



3) Proper weight init  $\begin{cases} \rightarrow \text{Glorot} \\ \rightarrow \text{xavier} \end{cases}$

4) Batch norm  $\rightarrow$  layer  $\rightarrow$

5) Residual Network  $\leftarrow$  CNN  $\rightarrow$  ResNET  
 $\hookrightarrow$  building block  $\rightarrow$  ANN

Exploding Gradient Problem

10, 10, 10, 10  $\rightarrow$  10000

random

Loss ↓ x

1 - 100

{ Gradient clipping }

Diagram illustrating the calculation of the derivative of the loss function  $\mathcal{L}$  with respect to the weight  $w'_{11}$  in an RNN:

- The RNN is represented by a box labeled "RNN".
- The derivative  $\frac{\partial \mathcal{L}}{\partial w'_{11}}$  is shown in a circle.
- The derivative is calculated as the product of the derivative of the loss function with respect to the input  $d_1$  and the derivative of the input  $x_1$  with respect to the weight  $w'_{11}$ .
- The input  $x_1$  is shown as a vector of values  $d_1, x_1, x_2, x_3$ .
- The derivative of the input  $x_1$  with respect to the weight  $w'_{11}$  is shown as a vector of values  $0, 1, 0, 0$ .

big number

$$d_1, x d_1, x^2 d_1, x^3 d_1$$

01

1000

$$\hookrightarrow W_n = W_0 - \eta \frac{\partial L}{\partial w}$$

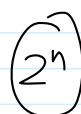
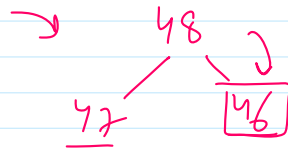
1000x

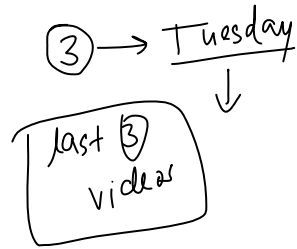
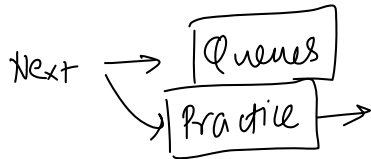
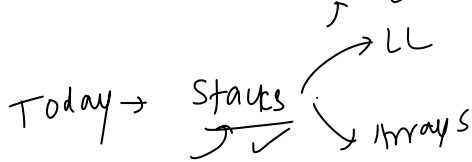
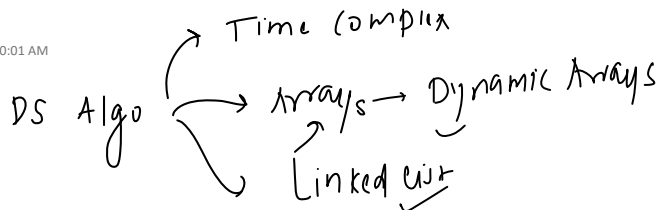
 $2 \times 10$ 

✓  
15)

18.

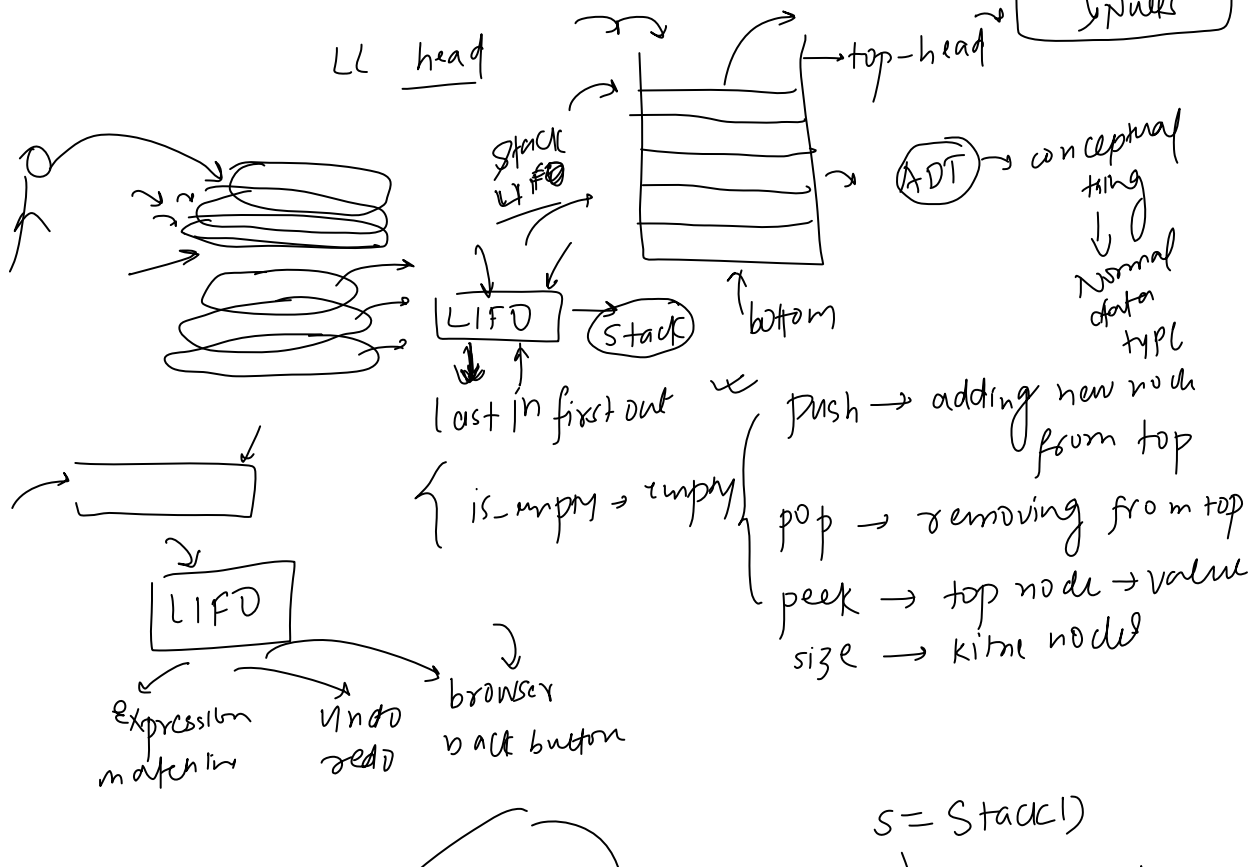
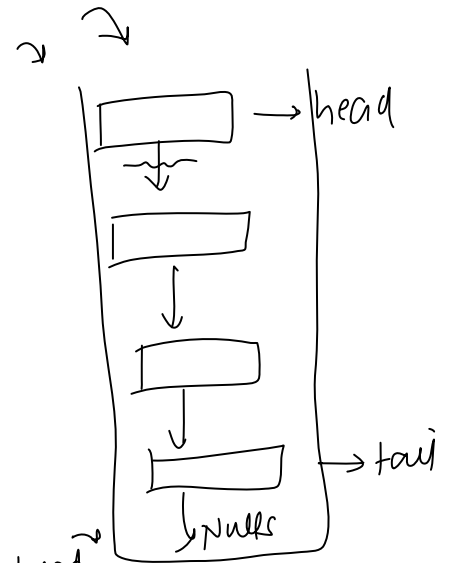
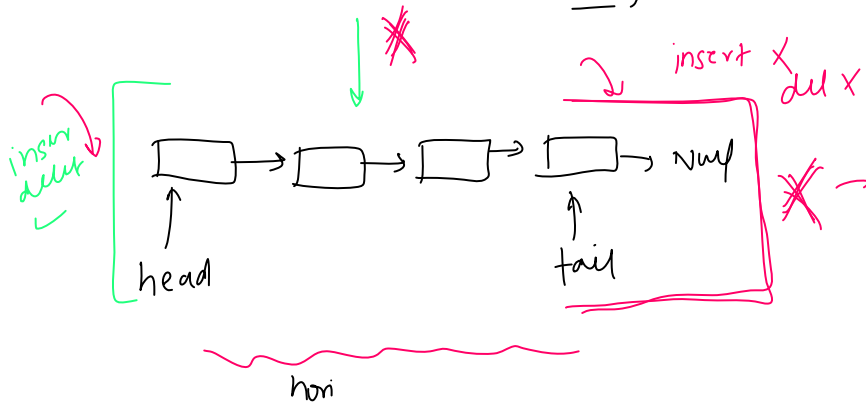
۱. ۵۲۱۰۰

[illegible]

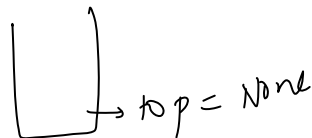
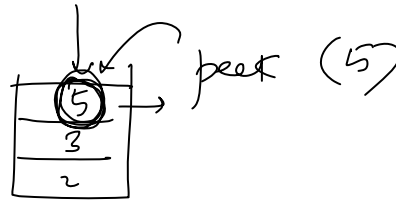
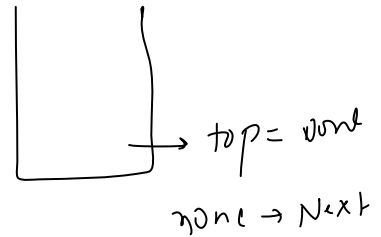
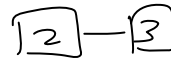
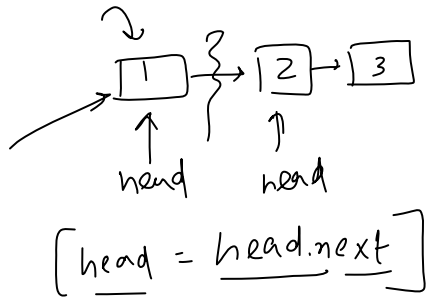
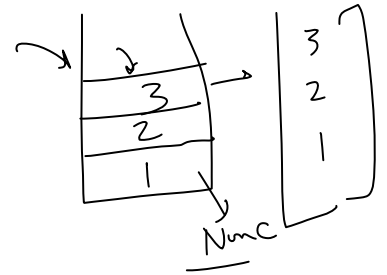
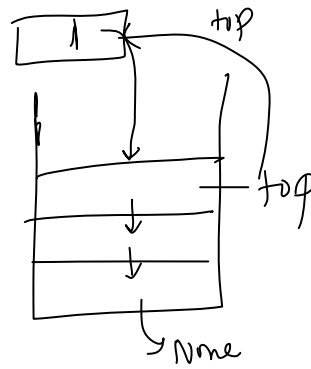
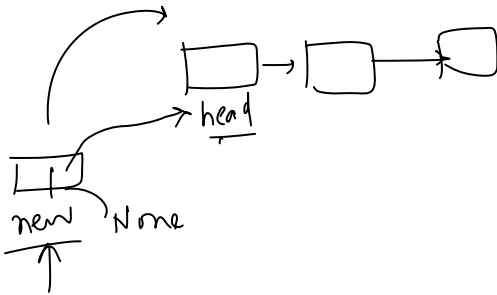
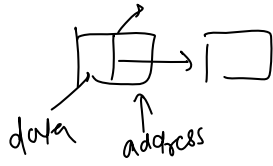
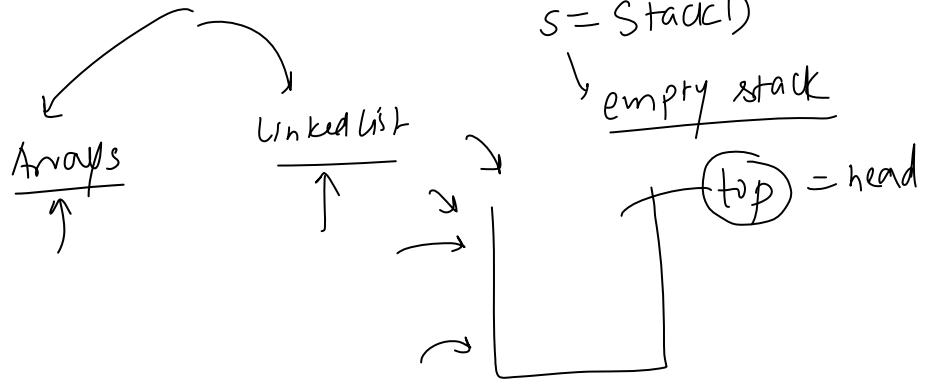


Stacks → DS

useful



memory



LIFO  
Advantage of Stack

Write

1) Insert → O(?) → O(1)  
push

Static → int a[50]

LL → Arrays

Dynamic



1)  $\frac{1}{\text{push}}$

2) delete  $\rightarrow O(?) \rightarrow O(1)$

int  $a[5]$

size = 4

top = -1  $\rightarrow 0 \rightarrow 1 \rightarrow 2 \rightarrow 3$

1, 2, 3, 4, 5

-1 = 4-1

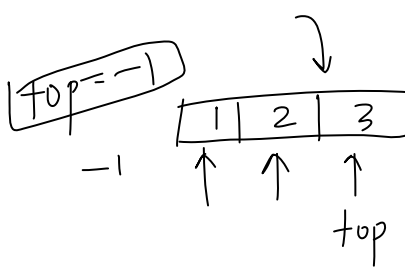
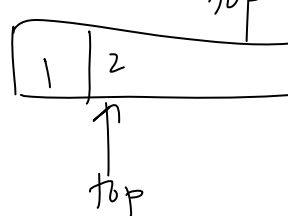
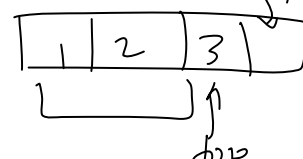
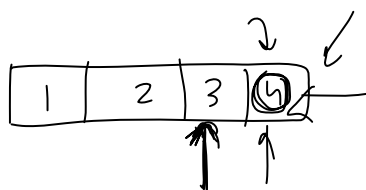
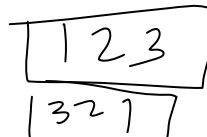
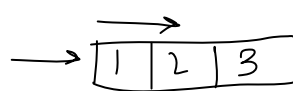
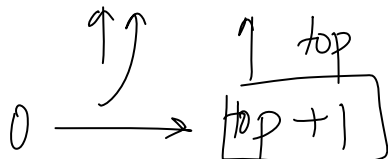
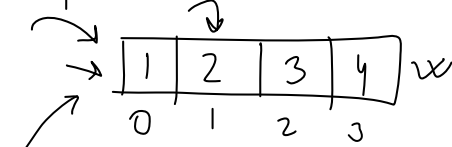
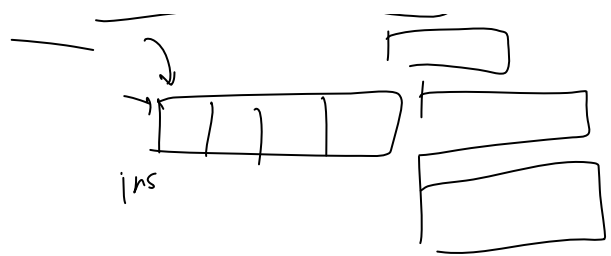
0 = 4-1

1 = 4-1

2 = 4-1

3 = 4-1

ans-

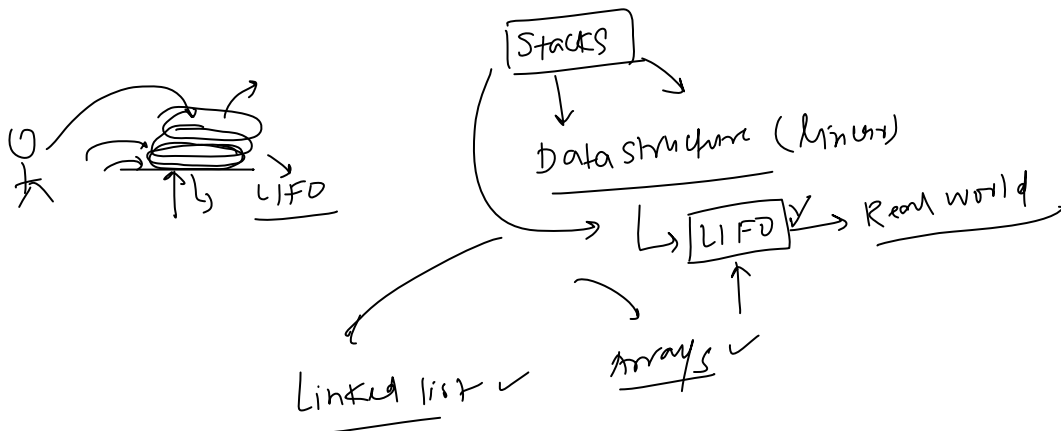
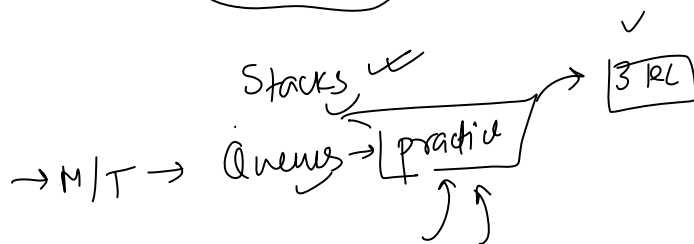
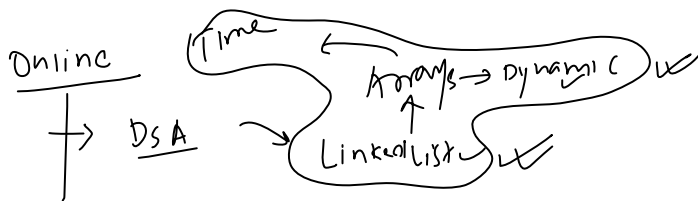
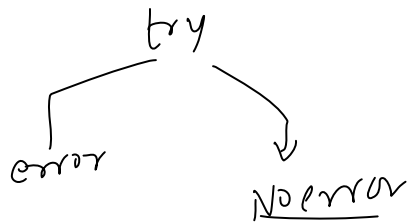
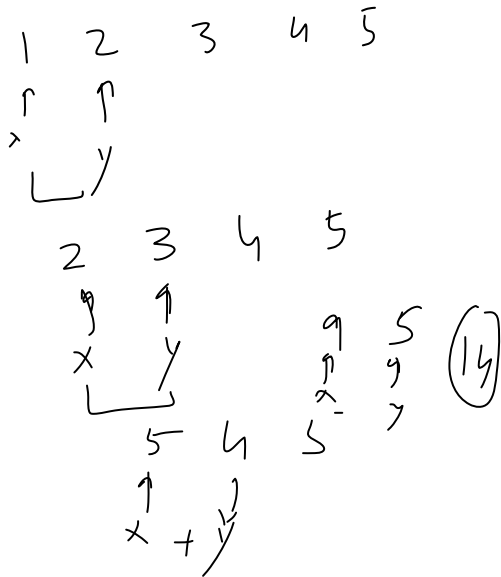


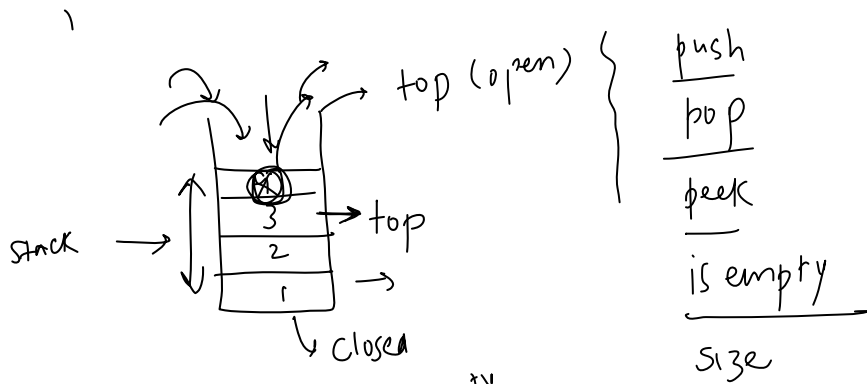
pop - 3

pop - 2

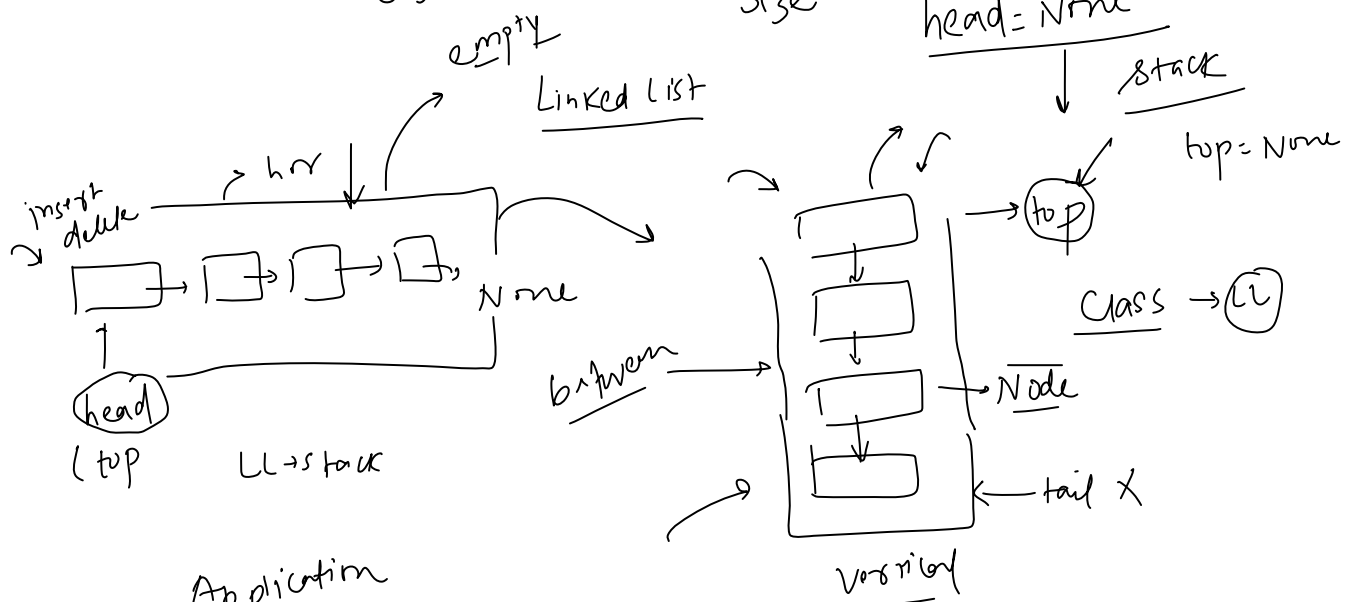
pop - 1

$x, y: x+y$   
 $\text{stack}[-1]$   $\xrightarrow{\text{pop}}$  1





Simple

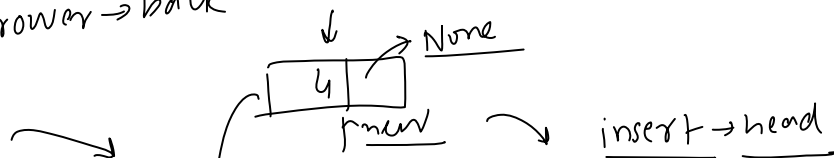


### Application

- expression
- undo/redo
- recursion →
- browser → back

LIFO

LL

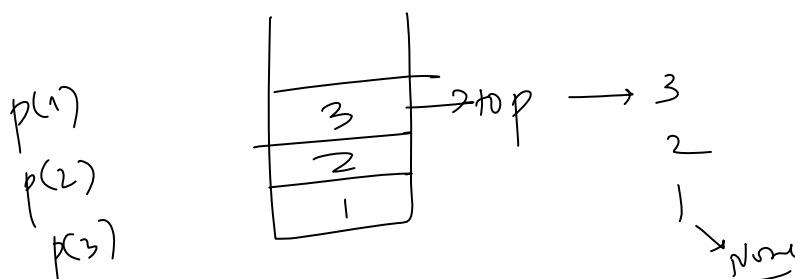


new.next = None

new.next = top

top = new

1 2 3



1. top = None

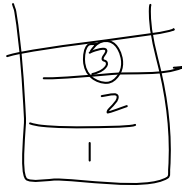


1  
p(3)

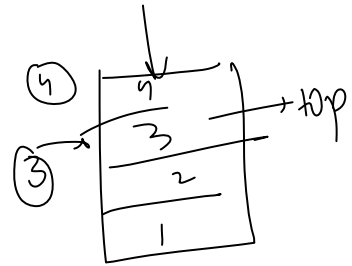
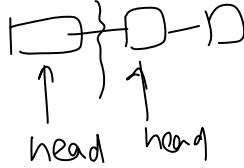
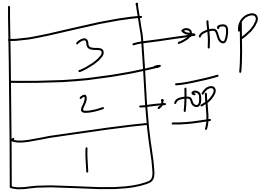
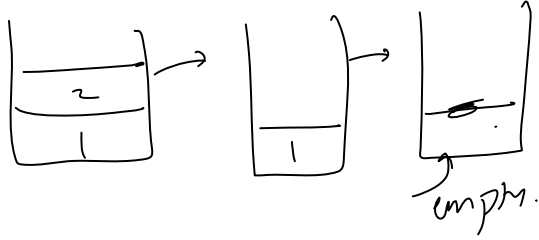


1  
None

top = none



pop



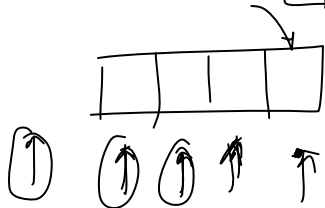
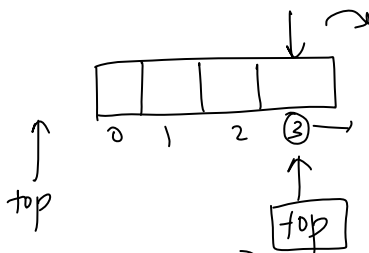
Array  $\rightarrow$  stack

Static (different)

Dynamic easy

HW

fixed  
size = 4 - 1  
value = 2



check  
 $top == size - 1$

