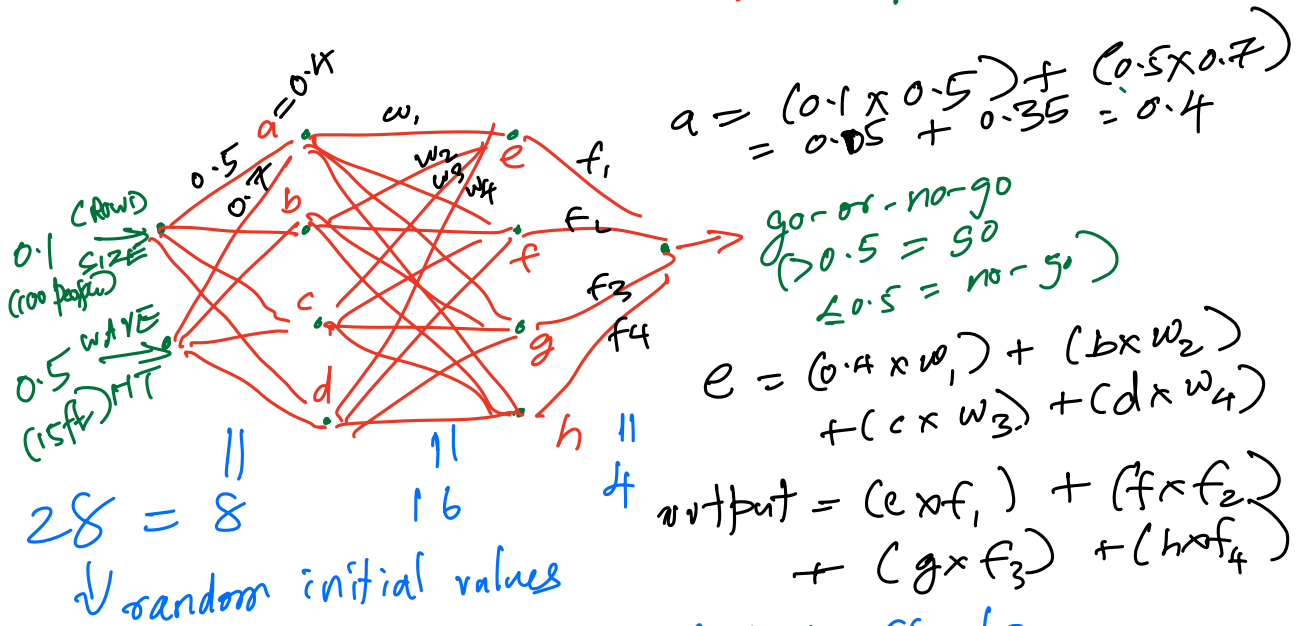


INTRO TO NNs

decide whether to go to Beach

inputs: crowd size (0-1000) $\xrightarrow{\text{SCALE (0-1)}} \text{e.g. 500 people } 0.5$
 wave size (1-30ft) $\text{e.g. 10ft } 0.3333$



$$28 = 8 \times 4$$

↓ random initial values

weights typically 2 or 4 byte floats

-1 to +1 e.g. 0.239872346
 -2 to +2 -2.4981789325167
 -3 to +3

of connections = 28 = # of weights

these 28 weights are called the
"parameters" of this NN

So, a NN is comprised of:

- its parameters (in a text file)
- the runnable executable to do a prediction flow using the parameters

Building a NN

- decide on architecture (# layers, # nodes)
- start with random parameters (weights)
- VERY COMPUTE HEAVY THE MODEL "LEARNS" ITS WEIGHTS
- TRAIN the model with training data set to modify the weights (W_{FINAL})
- TEST model for reliability, quality
- deploy !!

Use the model (NN)

- give it input (unseen data)
- COMPUTE HEAVY let NN compute the results using the deployed W_{FINAL} and generate output
- continuously monitor model's quality of results with user feedback

also called
a "prediction"

OPEN AI (Large Lang Models)

GPT 3.5 - 175B params

GPT 4 - 1.7T params

Some Chinese LLMs - 4 or 5T params

Google's Gemini Pro 1.5 - \approx 2T params

- NN architecture (layers, nodes, _{conns})
- weights being "learned" via training data
- deploy & predict

175 B weights \rightarrow apt 3.5 (open AI)

$$175 \text{ B} \times 2 \text{ bytes} = 350 \text{ B bytes} \\ = 350 \text{ GB RAM}$$

\therefore the weight file $\approx 350 \text{ GB bytes}$
may not be possible in laptop or
ordinary servers

Llama 2 $\Rightarrow 70 \text{ B params}$
 $\Rightarrow 31 \text{ B}$
 $\Rightarrow 7 \text{ B}$
 \swarrow $\rightarrow 140 \text{ GB RAM}$
deep cloud

$$7 \text{ B} \times 2 \text{ bytes} = 14 \text{ GB RAM}$$

Gemini Pro 1.5 = 2T Paramm

4000 GB RAM

MISTRAL

$8 \times 7 \text{ B}$