

7) ANN

Perceptron

- An independent building block (node) of ANN

Ex)

$$\text{output} \downarrow \quad \text{input} \downarrow \quad \text{weights} \downarrow \quad \text{bias} \downarrow \\ \hat{y} = \text{S}(x^T w + w_0) \quad \text{nonlinearity} \quad \text{Dot product} \quad (\text{intercept}) \\ \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \\ \text{activation fn} \quad (\text{linear weighted sum})$$

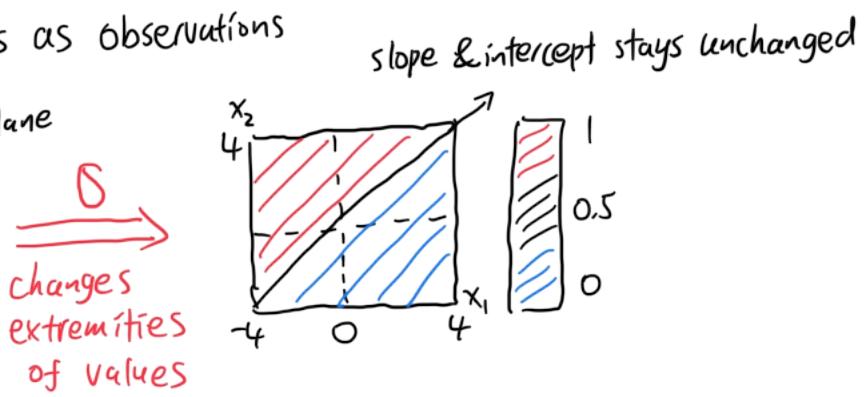
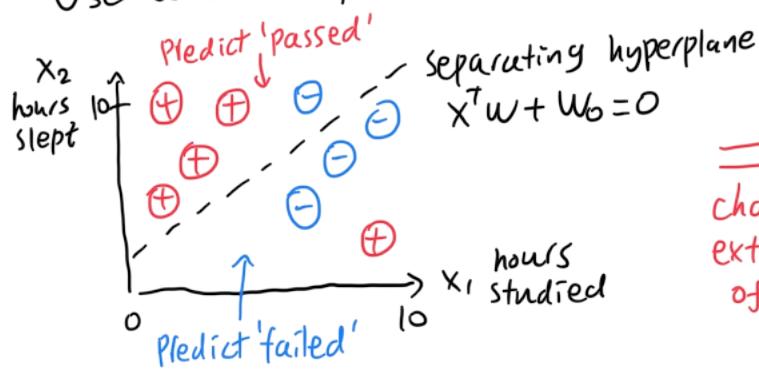
$$\hat{y} = \text{S}(x^T w + w_0)$$

nonlinearity
(activation fn) Dot product
 (linear weighted sum)

Forward Propagation

- It iteratively transforms **feature Space** (Geometric representation of data inputs)

- Use axes as features, coordinates as observations



Backward Propagation

- 2 common Loss functions

① **MSE (Mean-Squared Error) Loss**: $L = \frac{1}{2}(\hat{y} - y)^2$ for numerical predictions

② **Cross-Entropy / Logistic Loss**: $L = -\sum_{i=1}^n y_i \log(\hat{y}_i)$ for categorical predictions

- Train on individual loss function is **not ideal** bc:

i) Time consuming ii) can cause overfitting

- So instead, optimize **Cost function**: $J = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$, $m = \# \text{ batch of samples}$

• Goal: $W = \underset{W}{\operatorname{argmin}} J(W)$

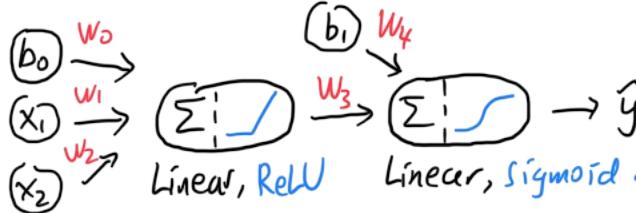
• Solution: **backpropagation** (gradient descent on W) $\rightarrow W := W - \alpha \frac{\partial J}{\partial W}$

ANN Single Layer

\downarrow
*i*th error(residual)

- Simple Regression: $y_i = \beta_0 + \sum \beta_i x_i + \epsilon_i$, the linear part of perceptron ($x^T w$)

\uparrow
ith data value \uparrow
predicted value

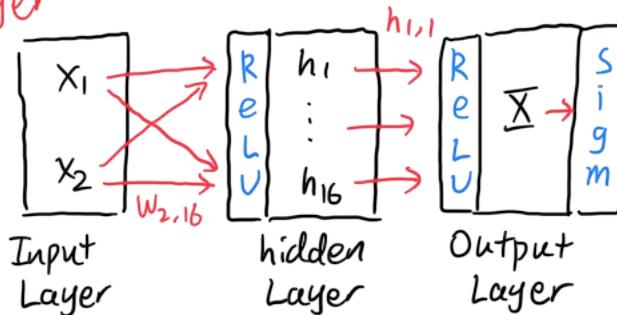
- model: 

Linear, ReLU Linear, Sigmoid

include this if ur doing classification
- Training workflow: for $i \in \text{epoch}$ { forward prop \rightarrow Compute & store loss \rightarrow backprop }
Make final forward prop \rightarrow Calculate accuracy:
- Testing workflow: Report final accuracy \rightarrow Show Losses vs. Epochs graph

ANN Multilayer

- architecture:

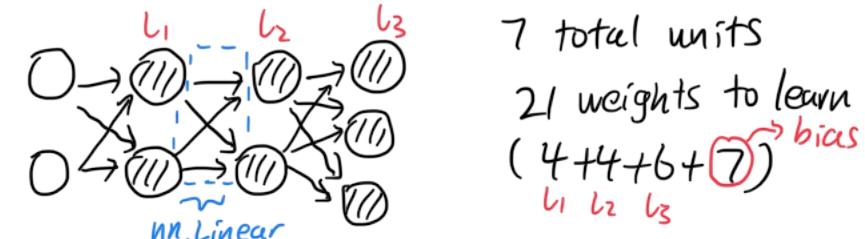
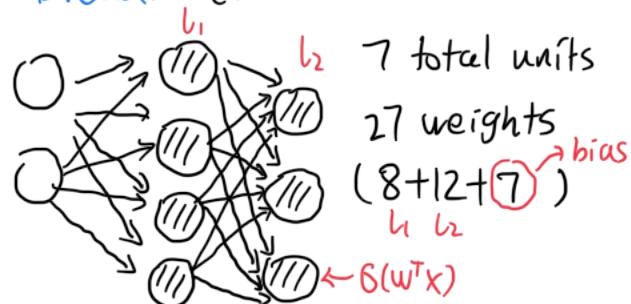


- b_i is bias in each layer that produces output to next layer (they don't take inputs!)
- W_{ij} (nn.Linear) comes from the transformation matrix into next layer

- U need **nonlinearities** in-b/t layers (else all layers will **collapse** into 1 layer)
- Use **softmax** over **sigmoid** for multi-classification : output is $\{P(y_1), \dots, P(y_n)\}$

Dimension of DL model

- **Breadth** (# nodes in hidden layer) vs. **Depth** (# hidden layers) are the dimension



Slower train time to learn more complex mappings