

II

History

2020 → Photo realism

2025 → Voice Cloning → Class demo.



extract patterns from data with neural networks.

Lab 1: Music Generation using neural networks

Lab 2: Facial detection systems

→ Pix

Lab 3: Fine-tune LLMs

Why deep learning

→ Learn features directly from data

Faces: ① Lines + Edges
② Nose, Ear etc.
③ Facial Structure } Bottom up

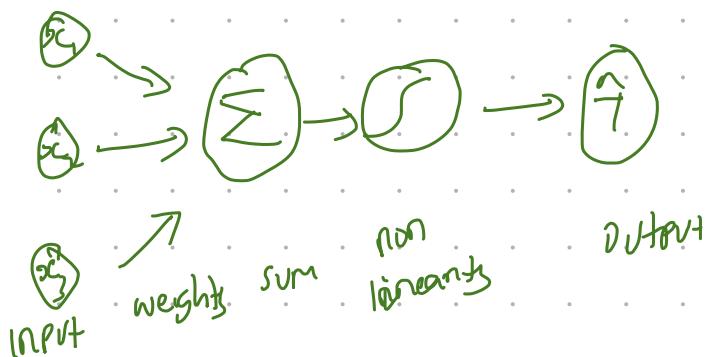
④ Can this be probabilistic

→ Discussion??

→ mix & match features.

III Perception

→ Single neuron. output



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

Linear Combination of output
Non linear activation function
Bias:

(IT)

why Transpose?

$$\hat{y} = g(w_0 + \vec{x}^T w) \quad \downarrow$$

where $\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $w = \begin{bmatrix} w_0 \\ \vdots \\ w_m \end{bmatrix}$

III Common Activation Functions

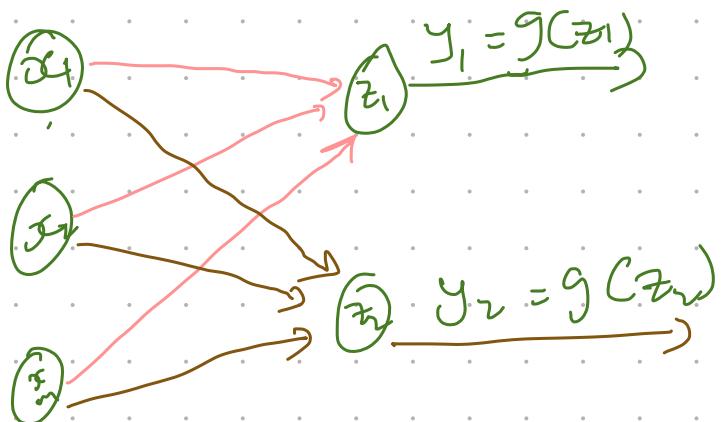
- tf.math.sigmoid(z)
torch.Sigmoid(z)
- ① Sigmoid $g(z) = \frac{1}{1+e^{-z}}$; $g'(z) = g(z)(1-g(z))$ → good for probability
- tf.m.relu(z)
torch.m.ReLU(z)
- ② ReLU $g(z) = \max(0, z)$, $g'(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}$

③ Hyperbolic Tangent

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad \left. \begin{array}{l} \text{tf.math.tanh(z)} \\ \text{torch.tanh(z)} \end{array} \right\}$$

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

III Multi-Output Perception



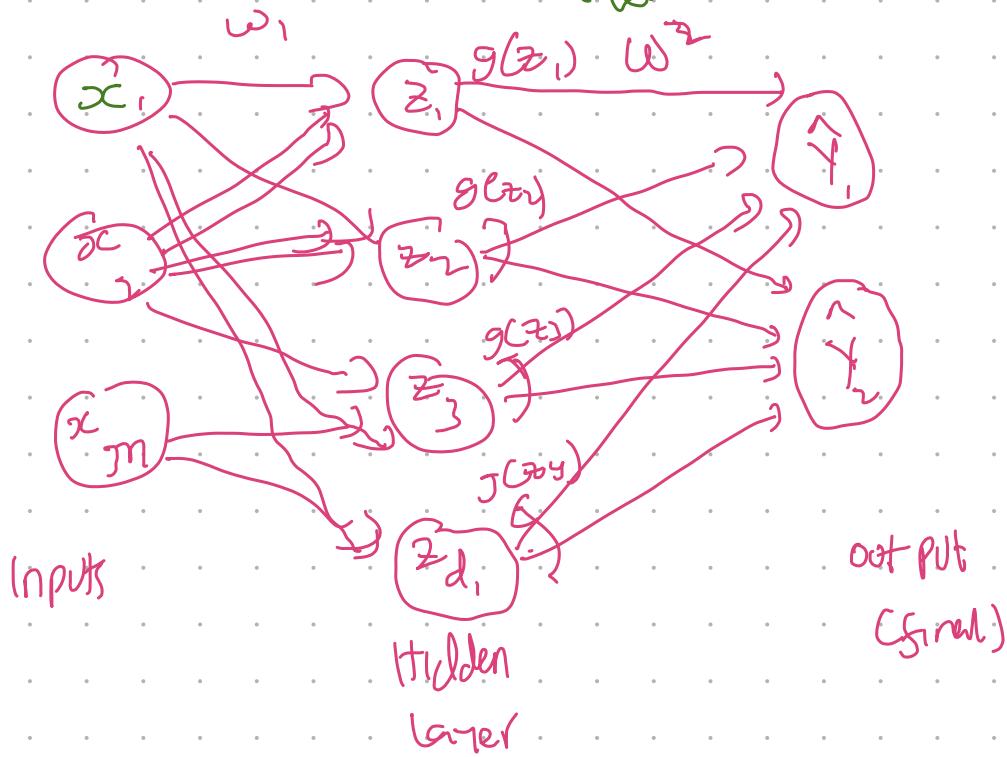
All inputs densely connected to all outputs never called Dense layers

$$z_i = w_{0i} + \sum_{j=1}^m w_j x_{j,i}$$

[Review Code: Min 3D]

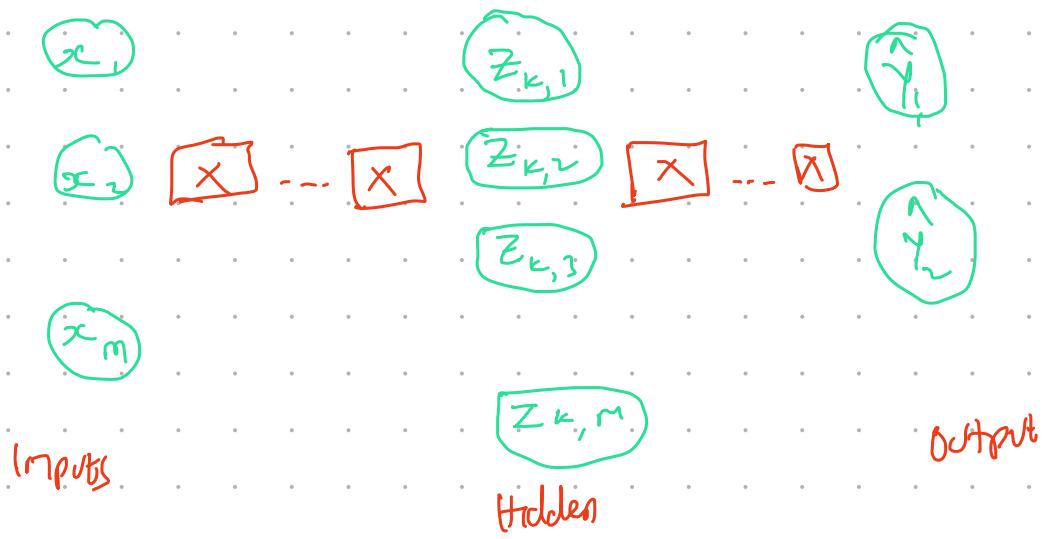
IV Single Layer Neural Network

→ contains hidden layer



Deep Neural Networks

Output of a hierarchical linear, non-linear combination



$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=0}^{n_{k-1}} g(z_{k-1,j}) w_{j,i} \quad (le)$$

Empirical Loss

The loss of our network measures the cost incurred from incorrect predictions.

$$L(f(x^i; w), y^i)$$

Predicted
Actual

Empirical loss \Rightarrow total loss on entire dataset.
 (cost function)

objective
function

$$J(w) = \frac{1}{n} \sum_{i=1}^n L(f(x^i; w), y^i)$$

predicted
actual

Average of every loss in the dataset. for every point

Binary Cross Entropy loss

maps to 0, 1, if 0.6 and 1, then no loss

$$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad [code]$$

✓ correct

Mean Squared Error Loss

$$= \frac{1}{n} \sum_{i=1}^n (\text{Actual} - \text{Predicted})^2$$

[code]

Loss Optimization.

TRAINING NEURAL NETWORKS

- Improve the network
- Build models that minimize the loss on a dataset
- Find weights that minimize loss over dataset (w) (\hookrightarrow)
- Loss function is a function of weights.

How to

- Start at random point.
- Compute gradient g
 - says which way is up. (locally)
 - Take a small step towards opposite direction
 - find local minimum

$$g = \frac{\partial J(w)}{\partial w}$$

Gradient Descent Pseudocode

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

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.

Gradient Computation is back propagation $(\frac{\partial J(w)}{\partial w})$



$$\frac{\partial J(w)}{\partial w_i} = \frac{\partial J(w)}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_1} \times \frac{\partial z_1}{\partial w_i}$$

Repeat for every weight in network using gradients from later layers.

* Tells us whether if we increase our weight, the loss will decrease.

(Just application of chain rule \rightarrow in theory)

In practice:

Training Neural Network is noisy

Loss Function Optimization

① difficult.

② Uses gradient descent.

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

↳ learning rate: } how quickly we move

How to set learning rates

① Try random learning rates and find one that just works

small \rightarrow get stuck

high \rightarrow unstable; diverge.

② Design adaptive learning rate that adapts as a function of gradient.

Multiple options: SGD, Adam etc

Gradient Descent vs Stochastic Gradient Descent

Gradient Descent: Computing gradient at every point is very EXPENSIVE

Stochastic Gradient Descent

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

2. Loop until convergence:

3. Pick a single data point (

if compute gradient, $\frac{\partial J_i(W)}{\partial W}$

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

5. return weights.

extra
J. N_i, S_i

Above but with mini batches

→ Pick a batch of B points

→ Smoother convergence

→ Larger learning rates

↳ possible to trust the gradients

↳ possible to parallelize

Overfitting

Underfitting → model does
not fully learn
the data

Training data is a
loose fit,
we want it
to generalize

Overfitting \Rightarrow over-complex, not general, value
more parameters

↑
Just memorizes all the training
set.

Regularization

→ D.S. converges complex model from
over-learning / memory.

How

① Dropout

- Set activations of some new rows to 0.
- forces network to not rely on any one node.
- forces to rely on multiple parameters

⑪ Early Stopping

- Stop training before it overfits