

DL - A2 : Team 101

Dharaneshwar - dshrisai | Suman - sumanman

Part 1-1 :

Autoencoder Architecture:

- Input Layer: 1000 neurons
- Hidden Layer 1: FC, 512 neurons, ReLU
- Bottleneck Layer: FC, 32 neurons
- Hidden Layer 2: FC, 512 neurons, ReLU
- Output Layer: FC, 1000 neurons, Sigmoid
- Loss Function: MSE
- Regularization: L2 Regularization on all weight matrices with a hyperparameter λ

Task 1) Trainable Parameters :

Layer	Parameters
Input	0
Hidden Layer 1 (FC 512)	$1000 \times 512 = 512,000$ + 512 (Bias)
ReLU	0
Bottleneck Layer	$512 \times 32 + 32 = 16,384$
Hidden Layer 2 (FC 512)	$512 \times 32 + 512 = 16,896$
ReLU	0
Output (FC 1000)	$512 \times 1000 + 1000 = 513,000$
Total	= 1,058,824

Task 2) Loss with L₂ Regularization:

$$\text{Loss} = \text{mse loss} + \lambda \sum \frac{\|w\|^2}{L_2 \text{ Term}}$$

regularization \rightarrow
hyperparameter

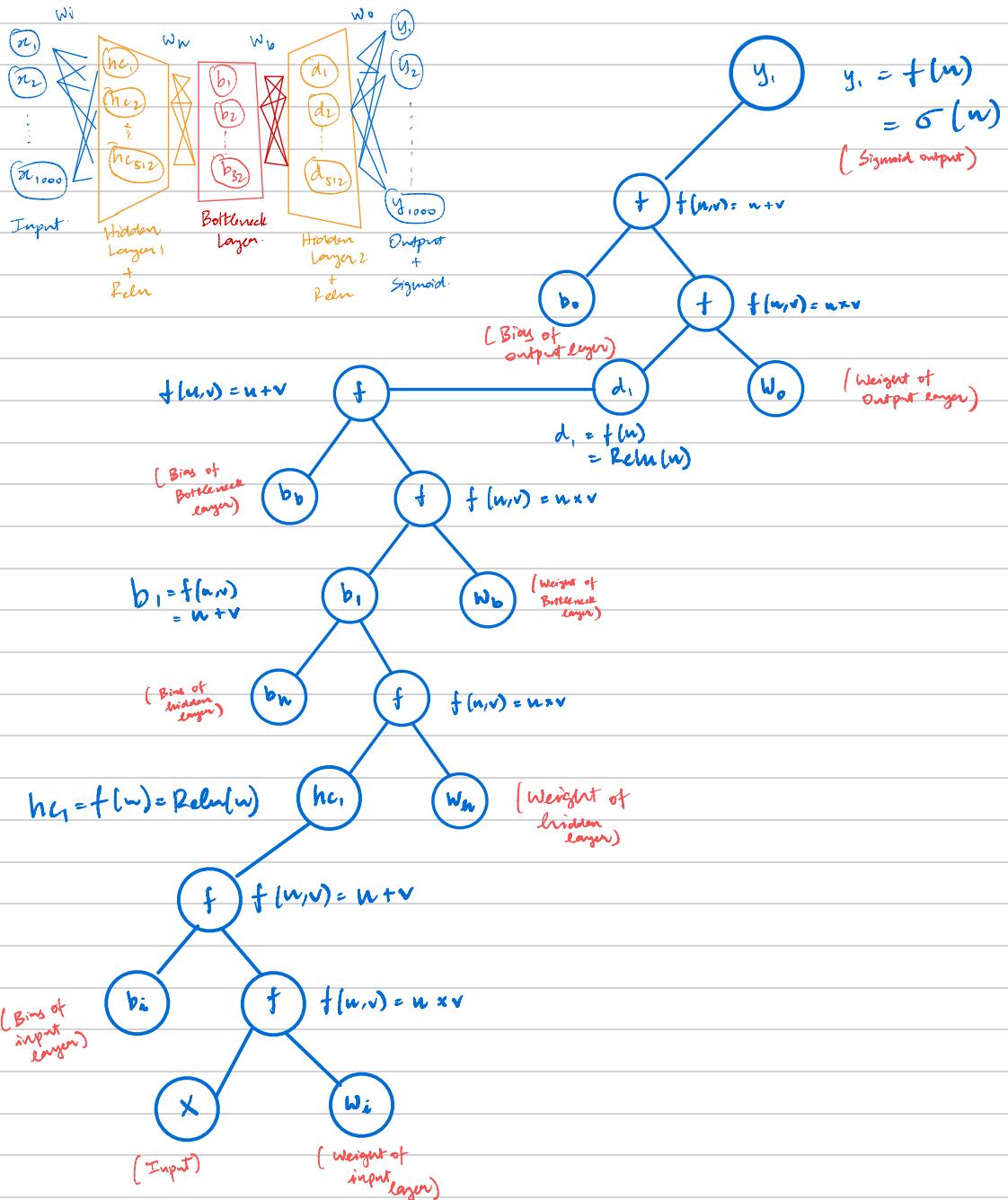
When $w \rightarrow 0$, the regularization term tends to zero which results in low loss. This can help prevent overfitting as smaller weights reduce model's complexity and its ability to fit noise.

During Back propagation,

$$w \leftarrow w - \alpha \left(\frac{\delta \text{loss}}{\delta w} + \lambda w \right)$$

\rightarrow This additional term (λw) during backpropagation, encourages smaller weights and smoothens the model to fit better.

Task 3) Computational Graph:



Task 4)

1. Autoencoders are trained on unlabelled data, which can lead to choosing the wrong threshold for loss.
2. New materials need additional training period causing degrade.
3. Imbalanced Dataset can cause generalization failing to flag subtle defects
4. Not all defects affect construction significantly which is an assumption.
5. High Dimensional data may cause inference lag in real-time usage.

Task 5)

1. We can use Variational Encoder to introduce probabilistic reasoning.
2. Incorporation of Temporal context in LSTM Autoencoder to capture time dependency.
3. Instead of static threshold we can introduce dynamic threshold.
4. Data augmentation to generate synthetic anomalies using simulation or GRANS
5. Deploying lightweight models for faster eval time quality control on site.

PART I. II

①

a) Linear Transformations

Each input vector x_i is linearly transformed into three different vectors

$$\text{Query } q_i = W_q x_i + b_q$$

$$\text{Key } k_i = W_k x_i + b_k$$

$$\text{Value } v_i = W_v x_i + b_v$$

where,

$\rightarrow W_q, W_k, W_v \in \mathbb{R}^{d_{\text{model}} \times d_K}$ are learned weight matrices

$\rightarrow b_q, b_k, b_v \in \mathbb{R}^{d_K}$ are bias vectors

$\rightarrow d_{\text{model}}$ is the input embedding size; d_K is the dimension of query, key & values

$q_i \rightarrow$ Represents what the token is 'looking for'

$k_j \rightarrow$ Represents what each token "offers"

$v_j \rightarrow$ Represents the actual content to use in the weighted sum.

b) Scaled Dot-Product Attention:

for each query vector q_i , attention scores are computed against all key vector k_j using the dot product:

$$\text{Score}_{ij} = \frac{q_i \cdot k_j}{\sqrt{d_K}}$$

one dot products product measures similarity between query q_i & key k_j

role of $\sqrt{d_K}$: to prevent large values when d_K is large, which could cause softmax to have extremely small gradients.

softmax over all j to get normalized attention weights:

Apply softmax over all j to get normalized attention weights:

$$a_{ij} = \frac{\exp\left(\frac{q_i \cdot k_j^T}{\sqrt{d_K}}\right)}{\sum_{j=1}^n \exp\left(\frac{q_i \cdot k_j^T}{\sqrt{d_K}}\right)}$$

c) The attention scores α_{ij} are used to weight the value vectors v_j using the following formula.

$$z_i = \sum_{j=1}^N \alpha_{ij} \cdot v_j$$

where,

- α_{ij} is the normalized attention weight between query q_i & k_j
- v_j is the value vector for token x_j
- z_i is the output for token x_i , a weighted sum of all value vectors based on their relevance to x_i

d) Optional Output transformation:

After computing z_i , we can often apply one more linear transformation:

$$z_i^{\text{final}} = w_o z_i + b_o$$

where, $w_o \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$

$b_o \in \mathbb{R}^{d_{\text{model}}}$

Impact of the final representation:

- 1) projects the attention output back to the original embedding space
- 2) enables stacking of transformer layers
- 3) learns a better latent representation by mixing attended information.

2) Computational Graph

