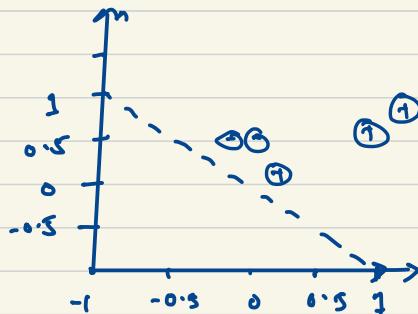


~~Prob 1~~Perception

x_1	x_2	Class
1	1	1
-1	-1	-1
0	0.5	-1
0.1	0.5	-1
0.2	0.2	1
0.9	0.5	1

- (i) In 3 Steps the learning algorithm will converge , 2 Step of w updates & last Step to check the correct prediction.
- (ii)

Initial Decision boundaryStep 1: $w = [1, 1]$

$$\text{Pred} = \mathbf{w}^T \mathbf{x}$$

Class +1 for pred > 0 Class -1 for pred < 0

Step 1

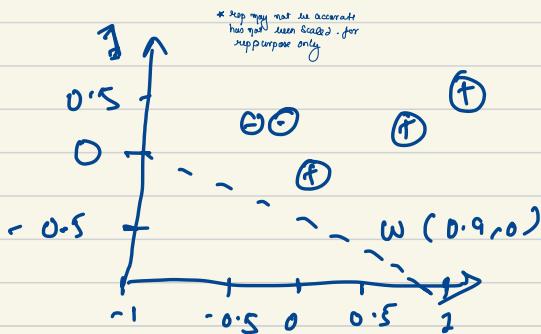
$$w = [1, 1]$$

$$\begin{aligned} p &= \text{pred} = w^T x \\ \text{Class} &\rightarrow +1 \quad y \geq 0 \\ \text{Class} &\rightarrow -1 \quad y < 0 \end{aligned}$$

		(W ¹)	
		pred	
x_1	x_2	Class	
1	1	+1	2 (Correct)
0.2	0.2	+1	0.4 (Correct)
0.9	0.5	+1	1.4 (Correct)
-1	-1	-1	-2 (Correct)
0.0	0.5	-1	0.5 (Wrong)
0.1	0.5	-1	0.35 (Wrong)

Correct Class pred = 4
Wrong $\leftrightarrow - = 2$

$$\begin{aligned} \rightarrow w &= [1, 1] - [0, 0.5] = [1, 0.5] \\ \rightarrow w &= [1, 0.5] - [0.1, 0.5] = (0.9, 0) \end{aligned}$$



C = Correct
WR = Wrong
(W¹)

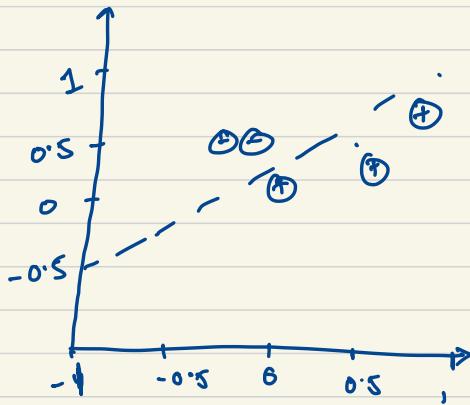
Step 2

$$w = [0.9, 0]$$

$$\begin{aligned} p &= \text{pred} = w^T x \\ \text{Class} &\rightarrow +1 \quad p \geq 0 \\ &\rightarrow -1 \quad p < 0 \\ \Rightarrow \text{Correct} &= 5 \\ \text{Wrong} &= 1 \end{aligned}$$

		(C/WR)	
		Pred	
x_1	x_2	Class	
1	1	1	0.9 (C)
0.2	0.2	1	0.18 (C)
0.9	0.5	1	0.21 (C)
-1	-1	-1	-0.9 (C)
0.0	0.5	-1	0.0 (W)
0.1	0.5	-1	0.15 (C)

$$\rightarrow w = [0.9, 0] - [0, 0.5] = [0.9, 0.5]$$



* If w may not be accurate
has not been scaled - for
suppose only

Step 3

$$w = [0.9, 0.5]$$

Correct prediction = 6

$$p = \text{prod} = w^T x$$

Class $\rightarrow +1 \quad p \geq 0$

$\rightarrow -1 \quad p < 0$

	x_1	x_2	Class	pred
	1	1	1	0.4 C
	0.2	0.2	1	0.08 C
	0.9	0.5	1	0.56 C
	-1	-1	-1	-0.4 C
	0.0	0.5	-1	-0.25 C
	0.1	0.5	1	~0.13 C

Problem 2:

GitHub Repo Link:
Code provided in Repo

OBSERVATIONS:

The Gurumukhi digit classifier is a simple Neural Network trained to recognize handwritten Gurumukhi digits (0-9). These are my observations.

(i) Data Loading & preprocessing:

- Data has been loaded from two folders 'train' and 'val' with 'train' folder having subfolders for each class.
- The train is used for training, 'Val' is used to test/Val & tune the folders
- There is a separate function to test images placed in a separate folder / given a path.
- The images are read in grey scale flattened and normalized.

(ii) Model Architecture:

- The classifier is designed in pytorch, it is a feed forward neural network with 4 fully connected layers
- The dropout rate is 0.21, as this prevents overfitting by not a proportion of neurons during training. This parameter has been tuned from different values.

- The hidden layers have size of [128, 64, 32] i.e. the first layer has 128 neurons, 2nd layer has 64 and 3rd layer has 32 neurons.
- Since there are 10 classes, the output size of 10 is mentioned in last layer. Each image is flattened to a size of 32*32 which becomes input size.

(3) Training & Validation:

- The model is trained with ADAM optimizer with a learning rate of 3e-3 and L2 regularization (weight decay) of 9e-4. The L2 regularization is used to prevent overfitting.
- Loss function used is cross-entropy loss, which is suitable for multi-class classification.
- Model is trained for 100 epochs and training & validation loss is mentioned at each epoch.

(4) Final performance & further improvement.

- The training & validation loss suggested that model did not overfit and generalizes well.
- Despite with less training examples, the model utilizes L2 regularization & dropout to prevent overfitting.
- However for even better results we can utilize a CNN along with a better GPU to generate a better model.

Problem 3

Git hub repo link:

(1) Model Summary

(A) Chart image classification using CNN

Architecture:

- 2 convolution layers each followed by batch Normalization, ReLU activation & max pooling.
 - first Conv layer: 4 input, 4 output channels, Kernel size: 3, Stride: 1, padding: 1
 - Second Conv layer: 4 input, 4 output, kernel: 3, stride: 1, padding: 1
- CNN layer are followed by a linear layer for classification
 - Linear Layer: 4 + 32 * 32, output features (5)
- Maxpooling is used twice right after each ReLU activation function.

Training:

- Trained for 25 epochs, ADAM optimizer, learning rate: 0.001
- Loss funcⁿ: Cross Entropy loss

task 3: finetuning pretrained 'AlexNet'.

- Imported necessary libraries including Pytorch, torchvision, numpy & pandas.
- Defined a Custom Dataset Class 'ChordDataset' to handle the chord data
- Loading, preprocessing / transforming the data using torchvision transform. Compose().
- Creating datasets and data loader by instantiating the custom dataset class.
- Load pre-trained 'AlexNet' & modify the first layer to accept 4 input channels and replace the last layer to match the no. of classes in dataset. (5 classes)
- prepare model for training: Move to GPU if there. loss: CrossEntropy
optimizer : SGD, lr: 0.001, momentum = 0.9
- Trained for 25 epochs & validated the model.

Problem 4

(a) Stochastic Gradient Descent (SGD) and Mini Gradient Descent are optimization algorithms used to minimize loss function. They differ in how they process the training data to update models params.

- SGD (Stochastic Gradient Descent):

- ⇒ Models parameters are updated after processing each individual training example.
- ⇒ This means gradient is computed and the weights are updated for every single data point
- ⇒ Can be expensive but can make model escape local min.

- Mini-Batch Gradient Descent:

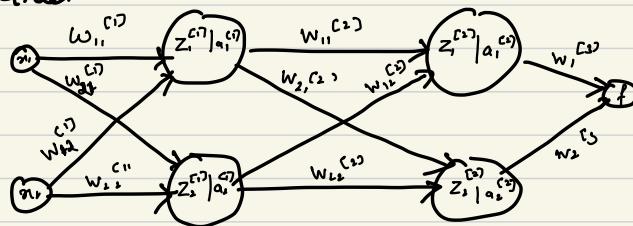
- ⇒ This can be considered as in between Batch Gradient Descent (where gradient computation & weights are updated after the entire dataset) and SGD.
- ⇒ In MBGD, the models parameters are updated after processing a small subset of training data called a mini batch.
- ⇒ This has computationally efficiency of Batch GD and escapes local minima like SGD.

(2)

Given $f = w_1^{(2)} a_1^{(2)} + w_2^{(2)} a_2^{(2)}$

compute $\frac{\partial f}{\partial z_1^{(2)}}, \frac{\partial f}{\partial z^{(2)}}, \frac{\partial f}{\partial Z^{(2)}}, \frac{\partial f}{\partial w_{11}}$

Given



$$A^{(2)} = \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \end{bmatrix} = \begin{bmatrix} \sigma(z_1^{(2)}) \\ \sigma(z_2^{(2)}) \end{bmatrix}$$

$$A^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \end{bmatrix} = \begin{bmatrix} \sigma(z_1^{(1)}) \\ \sigma(z_2^{(1)}) \end{bmatrix}$$

①

$$\frac{\partial Z}{\partial z_1^{(2)}} = \frac{\partial Z}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}}$$

$$f = w_1^{(2)} a_1^{(2)} + w_2^{(2)} a_2^{(2)}$$

$$a_1 = \sigma(z_1^{(2)})$$

$$= \sigma(w_{11}^{(2)} \cdot a_1^{(1)} + w_{21}^{(2)} \cdot a_2^{(1)})$$

\hookrightarrow Activation function

$$\Rightarrow \frac{\partial z}{\partial a_1^{(2)}} = w_1^{(2)}$$

$$\frac{\partial z_1^{(2)}}{\partial z_1^{(2)}} = \sigma'(z_1^{(2)}) = \sigma(z_1^{(2)}) * (1 - \sigma(z_1^{(2)}))$$

Hence: $\frac{\partial z}{\partial z_1^{(2)}} = w_1^{(2)} * \sigma(z_1^{(2)}) * (1 - \sigma(z_1^{(2)}))$

$$(3) \quad \frac{\partial f}{\partial z_1^{(1)}} = \left(\frac{\partial f}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial z_1^{(1)}} \right)$$

$$= \frac{\partial z}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial z_1^{(1)}}$$

$$+ \frac{\partial f}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial a_2^{(1)}} \cdot \frac{\partial a_2^{(1)}}{\partial z_1^{(1)}}$$

$$= w_1^{(2)} \cdot (w_1^{(2)} \cdot \sigma(z_1^{(2)}) * (1 - \sigma(z_1^{(2)})) \cdot w_1^{(2)}$$

$$+ \sigma(z_1^{(2)}) * (1 - \sigma(z_1^{(2)}))$$

$$= w_1^{(2)} \cdot (w_1^{(2)} \cdot \sigma(z_1^{(2)}) * (1 - \sigma(z_1^{(2)}))$$

$$(w_1^{(2)} \cdot \sigma(z_1^{(1)}) * (1 - \sigma(z_1^{(1)})) +$$

$$w_1^{(2)} \cdot \sigma(z_2^{(1)}) * (1 - \sigma(z_2^{(1)}))$$

$$(2) \frac{\partial f}{\partial z^{(2)}} = \left(\frac{\partial e}{\partial z_1^{(2)}}, \frac{\partial f}{\partial z_2^{(2)}} \right)$$

$$= \frac{\partial f}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} + \frac{\partial f}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}}$$

$$= w_1^{(3)} \cdot \sigma(z_1^{(2)}) + w_2^{(3)} \cdot \sigma(z_2^{(2)}) (1 - \cancel{\sigma(z_1^{(2)})} \cancel{(1 - z_1^{(2)})})$$

$$\textcircled{1} \quad \frac{\partial f}{\partial w_{11}} = \frac{\partial f}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial z_1^{(1)}} \cdot \frac{\partial z_1^{(1)}}{\partial w_{11}}$$

$$= w_1^{(3)} \cdot (\sigma(z_1^{(2)}) (1 - \sigma(z_1^{(2)})) \cdot w_{11}) \\ (\sigma(z_1^{(1)}) (1 - \sigma(z_1^{(1)})) \sigma_1)$$

$$z_1^{(1)} = w_{11}^{(1)} \cdot \sigma_1 + w_{12}^{(1)} \sigma_2$$

Problems

(1) Tuning hyperparams on test dataset pose a risk of overfitting the model which would fail to generalize, as we would be essentially optimizing the values on a specific dataset. Thus as a result the model may capture noise / patterns unique to test dataset.

⇒ The usual best practice is to have Train, Validate & test dataset where train set is used to train, validate is used to tune & test is the unseen dataset used to test. This would allow to obtain an unbiased model performance.

(2)

The two strategies to reduce overfitting are:

(1) Regularization: • Reg. is added to reduce overfitting by adding a penalty term to loss function.
• The most common reg methods are L_1 & L_2 regularization
• L_1 adds the absolute value of the weights while L_2 adds the squared value of the weights
• These penalty discourage model to learn overly complex representations

(2) Dropouts: • This technique is specific to NNs, During training dropout randomly "drops" or deactivates a certain percentage of neurons.
• This prevents the model from relying too heavily on

On any single neuron or a feature

(3)

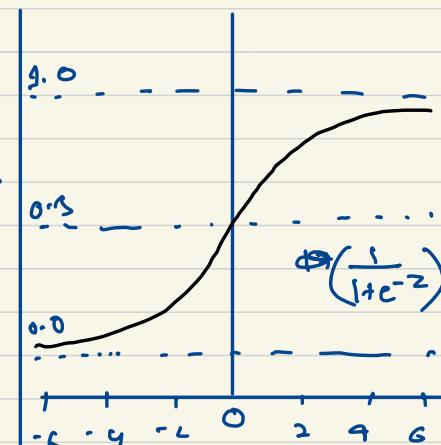
- Input layer size:
- Depends on the no of features in the dataset used for the specific problem
 - for ex., if we are dealing with image of size 28×28 pixels, the input layer would have $784(28 \times 28)$ nodes
 - for ex. for tabular data with 10 features input layer would have 10 nodes

- Output layer size:
- Depends on the type of problem being solved & no of classes or target variable.
 - for binary classification, we can have a single o/p neuron with sigmoid activation
 - for multiclass, o/p layer should have same number of neurons as classes with softmax activation
 - for reg., o/p with one neuron with single activation funcn.

(4)

The Sigmoid activation function: $\frac{1}{1+e^{-x}}$

- takes any real valued input and maps it to a value b/w 0 & 1
- It is S shaped smooth curve, with $f(n)$ approaching 0 as $n \rightarrow -\infty$ and 1 as $n \rightarrow +\infty$ and centered around 0.5 when $n=0$. It is differentiable.
- Used in binary classification.



(5)

- The learning rate is a hyper param that determines the step size during optimization process. It controls how much of a change is applied in updating the parameters of the model
- A smaller learning rate leads to smaller updates while a larger learning rate leads to larger updates
- Choosing the right lr is appropriate for proper convergence of the model (decreasing cost w.r.t parameters)
- If lr is too small, it will take a long time to converge can possibly get stuck in local min
- On the other hand a large lr may overshoot or miss the minimal/optimal condition and may never converge.
- A learning rate given in ques q.s is very very large and hence will cause instability in training process the network may not converge. It is recommended to start with smaller values of lr (0.001, 0.0001 etc)