

Table of Contents

Headers

Machine Learning Basics

Page

3

Underfitting/Overfitting/Regularization

4

Cross-Validation

5-6

Approach 1 and 2 to Data Augmentation

7-8

The Perceptron Algorithm

9

Parameter Initialization

10

Forward Propagation of Multilayer Perception (MLP)

11-13

Backpropagation and Formula Derivations

14-20

Ways to Optimize Loss Function

21-25

Common Functions

26

Convolutional Neural Networks

27-29

Confusion Matrix and Example statistic Analysis

30-33

Understanding Colab/MPS/Multiprocessing

34-35

Efficiency of Different Compilers

36

Machine Learning Basics

• Feature - attribute/characteristic of what you are studying

◦ Often select, transform, or create new features to improve performance of the model (called feature engineering)

◦ Inputs to the model that are useful to "train" it

◦ Often encoded in a feature/design matrix

• Classification models aim to separate classes in the feature space

• Clustering algorithms groups "close" data points

• Regression looks at relationship between predictors and response

• Supervised machine learning - each training example is paired with output label (included in dataset)

◦ Difficult for large datasets to all be labeled

• Parameters like weights/biases are adjusted automatically by algorithm during training, but hyperparameters are determined by programmer manually

◦ Different NN's have different parameters; 6NN parameters include edge weights and knowing which nodes to connect in graph

• Epochs - number of times training dataset is passed through NN during training

◦ Hyperparameters include number of hidden layers and neurons in them, learning rate, etc.

• Inference - another name for testing (what model infers about)

• Even though we understand some mathematics of how NN works, it's difficult to interpret what each neuron is doing

◦ E.g. for image recognition all it does is minimize a cost function, doesn't actually observe patterns

• Tensor - multidimensional array

Artificial intelligence - theory/development of computers to do tasks normally requiring human intelligence

↓

machine learning - gives computer ability to "learn" (adjust parameters) without being explicitly programmed

↓

deep learning - machine learning algorithms using neural networks

Underfitting, Overfitting, and Regularization

Generalization - the ability to perform well on previously unseen inputs

In machine learning we not only optimize by reducing training error but also test/generalization error

Underfitting - high training and validation error

High bias, low variance
Some error remains of training set used

Model is too simple and can't recognize complex features

Solution: add complexity to model architecture

Overfitting - training error is low but validation error is high, indicating memorization/overfitting of training

Also known as capturing noise, random variations/errors that do not represent true underlying patterns (e.g. drawing a 1 that looks like a 7)

Low bias, high variance

Features in training recognized very quickly (problem for small datasets)

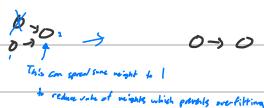
Solution: Use regularization or increase dataset size from data set using data augmentation

Capacity determines how well model can capture patterns in the training data

Regularization - nullifying effect of certain neurons (prevents overfitting)

↑ Complexity in NN \Rightarrow complex nonlinearity fit to training data \Rightarrow overfitting

Dropout regularization - For each batch, set half the neuron activations in each layer to 0; prevents a neuron from being too reliant on output of others



Say $\hat{a}^c = \begin{bmatrix} 0.3 \\ 0.4 \\ 0.5 \end{bmatrix}$. A binary dropout mask of same size as \hat{a}^c is created by first creating tensor filled with random values between 0 and 1.

Then, due to Law of Large Numbers, we can create the binary mask by seeing if the result is $\geq \text{keep_rate}$. Finally do $\hat{a}^c \odot \text{mask}$ to get final activation of layer C.

Dropout only used in training; the activations are also scaled in training by $\frac{1}{\text{keep_rate}} = \frac{1}{1-\text{drop_rate}}$ so that are similar values in training and inference

The less neurons I keep the more I scale their activations

Dropout is more useful when training over many epochs

No dropout \Rightarrow each neuron relies on precise patterns from all other neurons \Rightarrow learns very specific features/mode from training data \Rightarrow overfitting

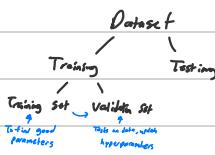
Dropout \Rightarrow each forward pass uses different subset of neurons \Rightarrow neuron doesn't rely on very specific neurons/features \Rightarrow less overfitting

BUT too much dropout \Rightarrow too few neurons in ff \Rightarrow can't capture patterns in data \Rightarrow underfitting \Rightarrow does worse on test set

Additional computations can slow down training

L2/ridge regularization - add term to cost function to penalize larger weights

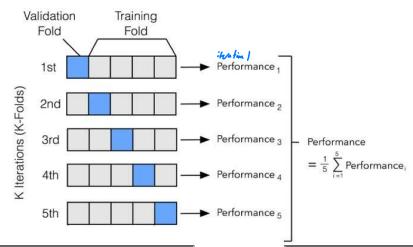
Cross-Validation



- Validation data estimates generalization error (does same calculation as test set) and updates hyperparameters (usually 80/20 split)
- Validation set is good because otherwise you adjust hyperparameters to make test accuracy better, which may overfit to the test data and not generalize
- Test data being completely independent of how the model is trained makes it an accurate measure of model performance on unseen data (\star)
 - If you tune it to test data with same distribution as training in theory that gives similar hyperparameters but then you have nothing to evaluate \rightarrow
 - If you already did hyperparameter tuning with validation accuracy, it is okay to train model multiple times and choose the one that generalizes best (high test accuracy)

Cross-validation involves splitting training data into K folds/subsets, and using $K-1$ folds for training and the last one for validating

- This process is repeated K times with each time having a different fold as validation set (each point in validation one and in training set K-1 times)
 - K-folds cross-validation follows this, where on the i-th trial, the i-th fold is the validation set
 - Each data point uses for training and validation eventually
 - Each iteration has a specific (usually same as rest) # of epochs before going on to next iteration with different validation set (hyperparameters didn't change much here)
 - After going through K-folds and seeing errors, adjust hyperparameters
- \star Every point is eventually in training data, so K-folds better captures distribution of all training data so hyperparameters won't be super different



- K-folds used to evaluate performance of deep learning model; averaging performance over multiple folds reduces variance of performance caused by randomness of single train-test split
 \uparrow
 (lower as good number)

- K folds \rightarrow ↑% of train data across all folds \Rightarrow model trained on higher proportion of data \Rightarrow model learns better patterns and validation accuracy/model performance

averaged over more folds \Rightarrow accuracy/performance metrics less influenced by any single split/fold class variance in val accuracy)

\Rightarrow hyperparameters are more independent of split used \Rightarrow better hyperparameter tuning

\uparrow
 be more training data and
 more folds to avg

- K folds \rightarrow smaller validation subset that may not be representative of all the data \Rightarrow more bias
- K folds \rightarrow less trials to average \Rightarrow more variability in performance

After finding best model with cross-validation, train your new model using all the training data

- Distribution - refers to spread of features across dataset (i.e. class distribution, brightness, RGB values, etc.)
- Theoretically (although training/validation split perfectly capturing all training data distribution is rare), if that happened in split for cross-validation, hyperparameters would be independent of split
- Since this is theoretical, even with cross-validation the best hyperparameters for model trained on full training data could be a little different

• Various ways to find best hyperparameters (with hyperparameter tuning)

• Grid search - cycle through each combo of hyperparameters and see how well they perform on validation set

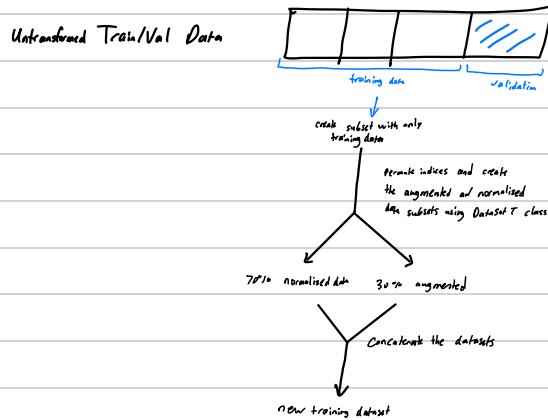
• Choose ones that lead to consistent performance on each fold (would probably handle outliers distribution well!)

Properly Creating Train Dataset with Normal and Augmented Examples

- Can help model differentiate between various classes
- Shown below are two approaches I came up with to create a **new dataset** that contains both normal and augmented data
 - One reason this was done is because more training data resulted in a better model
- These approaches to create the **new dataset** for training were used in both k-folds cross validation and the final training of the model
- Importantly, neither of these approaches augment examples of the validation set during k-folds cross validation
- The visualizations of each approach use the following split: 70% normalized images, 30% augmented

Approach 1 - uses custom Dataset class to define transformations per image

K-folds:

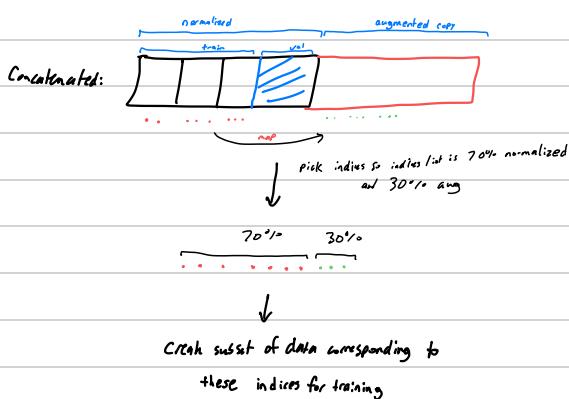


• For full training data, the same applies except there is no validation set

Approach 2

K-folds: Each fold...

- Generate train and val indices on the normalized dataset
- Concatenate an augmented dataset same size as normalized
- Map the train indices to the augmented copy
- Choose 70% of all original train indices to be from normalized and the other 30% from augmented
- Create a subset of data with these indices; this dataset is the same size as original normalized train subset



• For final training with all training data, use similar approach except...

- To maximize number of training examples in dataset, use 100% of the normalized data and an additional amount of augmented data so the proportion is still 70% normal, 30% augmented

$$(\text{total training examples})_{\text{prop-normal}} = |\text{train normalized-data}| \Rightarrow \text{total training examples} = \frac{|\text{train normalized-data}|}{\text{prop-normal}}$$

• Additionally, for the augmented data, after analyzing confusion matrix, I replaced some examples of classes the model performed well on with classes the model performed poorly on

Approach 2 performed better, and I believe it did due to the following:

- Mapping the indices onto a new dataset provided greater similarity between the normalized and augmented training examples
- These similar examples in Approach 2 train dataset could have made the model learn better since it is seeing more similar data as opposed to an augmentation of an image it has never seen in the batch before (as in Approach 1)
- Also improved performance in final training due to more training data

The Perceptron Algorithm

- It's a supervised binary classification model whose goal is to find the hyperplane that separates the two classes of data points in a feature space
- Used when a dataset is linearly separable (falls XOR problem)
- Single layer neural network (parameters that are adjusted in the algorithm are weights and biases)
- Called "perceptron" because the model stimulates human perception and perceptron perceives data

$0 = \vec{w} \cdot \vec{x} + b$ is hyperplane

w_1, w_2, \dots, w_n

$$y = \vec{w} \cdot \vec{x} + b$$

↑
weights
↑
input vector
↓
bias

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad \vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

If you have n features/dimensions, you have n components in \vec{w} and \vec{x}

- The weight components adjust hyperplane orientation/slope

- Step function - makes input abruptly "step" from one value to another at a specified threshold (a piecewise function)

This perceptron contains...

- An input layer, where each neuron contains a feature of the data point fed in
- An output layer, consisting of one neuron with the binary output

For each data point i ...

- Calculates the weighted sum $\vec{w} \cdot \vec{x}_i$ → adds the bias → passes result through step function → updates \vec{w} and b when perceptron misclassifies the data → finish when $\hat{y}_i = y_i$ for all training examples (i.e. the algorithm converges)
- Step function: $\hat{y} = \begin{cases} 1 & \vec{w} \cdot \vec{x} + b > 0 \\ 0 & \vec{w} \cdot \vec{x} + b \leq 0 \end{cases}$

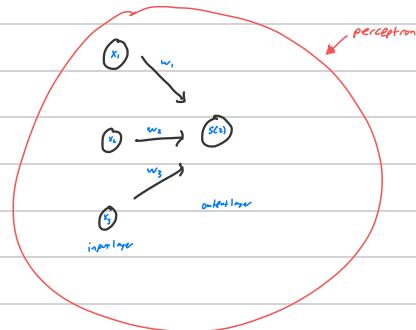
• Updating parameters

$$\vec{w} \leftarrow \vec{w} + \eta (y - \hat{y}) \vec{x}$$

$$b \leftarrow b + \eta (y - \hat{y})$$

↑
learning rate

$$S = \text{step}, z = \sum_{i=0}^n w_i x_i$$



- Cannot separate nonlinear classes or solve XOR problem, multiple perceptrons needed with non-linear activation

- With perceptron algorithm the hyperplane cannot separate the classes

Parameter Initialization

• He/Kaiming initialization

- Technique used to randomly initialize parameters when ReLU is used
- Weights are initialized as random numbers from Gaussian/normal distribution (then multiplied by scaling factor)
 - Mean at 0 \Rightarrow center of bellcurve at 0 \Rightarrow most weights are zero since are small positive and negative
 \downarrow
so is spread of curve
- Scaling factor - $\sqrt{\frac{2}{\text{in} \cdot \text{out}}}$ $\xrightarrow{\text{narrow distribution width?}}$
- Scaling factor helps set variance of weights to proper value and prevent vanishing/exploding gradients
- We also want inputs to be normalized because larger inputs \Rightarrow numerical instability and $\frac{d}{dx}[\text{activation}]$ is extremely small \Rightarrow vanishing gradient problem

Multilayer Perceptrons and Forward Propagation

• Convolutional neural network - good for image recognition

• Long short-term memory network - good for speech recognition

• MNIST dataset - popular dataset of handwritten digits 0 through 9 in black and white

• Neural networks have input layer, hidden layers, and output layers

◦ For image recognition it is common for each neuron of the inner layer to represent a pixel attribute (e.g. brightness/color)

◦ A neuron is like $f(\vec{x}, \vec{w}) = \sigma(\sum_i w_i x_i + b)$

↓
input
layer
activation
function

• The "basic" neural network is a multilayer perceptron (MLP)

• Learning - finding the right parameters

• For classification, neuron in output layer with highest activation is the predicted class

• Superscripts indicate the layer; subscripts indicate the neuron in the layer

◦ All connections have a weight, and there is 1 bias in each neuron

$$\vec{z}_j^e = \sum_k w_{jk}^e a_k^{e-1} + b_j^e$$

• $\vec{z}_j^e = (\sum_i w_{ij}^e a_i^{e-1}) + b_j^e$

◦ \vec{z}_j^e is the weighted sum for neuron j in layer e

◦ w_{ij}^e is the weight for connection from neuron i in layer $e-1$ to neuron j in layer e (flows backward)

◦ b_j^e is bias for neuron j in layer e

• $a_j^e = \sigma(\vec{z}_j^e)$

↓
activation

• O is output layer, h is hidden layer

• Feed forward calculated more efficiently with linear algebra

Shorthand for $a_j^e = \sigma(\sum_i w_{ij}^e a_i^{e-1} + b_j^e)$

$$W = \begin{bmatrix} w_{00} & w_{01} & \dots & w_{0n} \\ \vdots & & & \\ w_{m0} & w_{m1} & \dots & w_{mn} \end{bmatrix} \quad \vec{a}^{e-1} = \begin{bmatrix} a_0^{e-1} \\ \vdots \\ a_n^{e-1} \end{bmatrix}$$

matrix entries between two layers

• W_k is all connections from layer $e-1$ to neuron k in layer e

• $W_k \vec{a}^e = \vec{w}_k^T \cdot \vec{a}^e = \vec{z}_k$, the weighted sum for that neuron

$$\vec{z}^e = W^e \vec{a}^{e-1} + \vec{b}^e = \begin{bmatrix} w_{00} & w_{01} & \dots & w_{0n} \\ \vdots & & & \\ w_{m0} & w_{m1} & \dots & w_{mn} \end{bmatrix} \begin{bmatrix} a_0^{e-1} \\ \vdots \\ a_n^{e-1} \end{bmatrix} + \begin{bmatrix} b_0 \\ \vdots \\ b_n \end{bmatrix} = \begin{bmatrix} \sum_i w_{0i}^e a_i^{e-1} \\ \vdots \\ \sum_i w_{ni}^e a_i^{e-1} \end{bmatrix} + \begin{bmatrix} b_0 \\ \vdots \\ b_n \end{bmatrix} = \begin{bmatrix} z_0 \\ \vdots \\ z_n \end{bmatrix}$$

weighted sum + bias for neuron 0 in layer e

$$\sigma(\vec{z}) = \vec{a}^e = \begin{bmatrix} a_0^e \\ \vdots \\ a_n^e \end{bmatrix}$$

• Each component of the matrix vector product is a weighted sum for that corresponding neuron in layer e

• There is a weight matrix and a bias vector for each layer of NN

• Z usually matrix (each row is output of an example; # of rows is batch size) ; $Z = AW + b$

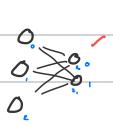
↑
for parallelism

↑
only adding this to each row

batch size = input x output

Analyzing the Forward Propagation Process

Between Layer 0 and 1



$$A_0 = X$$

$$z'_0 = w_{00}x_0 + w_{01}x_1 + w_{02}x_2 = \begin{bmatrix} w_{00} \\ w_{01} \\ w_{02} \end{bmatrix} \cdot \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} w_{00} & w_{01} & w_{02} \end{bmatrix} \vec{x}$$

$$W_1 X + B_0 = Z$$

↑ Parameterization

$$z'_1 = w_{10}x_0 + w_{11}x_1 + w_{12}x_2 = \begin{bmatrix} w_{10} \\ w_{11} \\ w_{12} \end{bmatrix} \vec{x}$$

$$\vec{Z}' = \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = W \vec{x}$$

Parameter matrix

$$z'_i = w_i' \vec{x} = w_i^T \vec{x}$$

↓ row vector

• i-th row of matrix W contains weight of connecting from i-th neuron in layer 0 to all neurons in layer $L-1$

• For batch size > 1 there are multiple \vec{x} 's with same calculation, which corresponds to $X = \begin{bmatrix} \vec{x}_1 & \vec{x}_2 \end{bmatrix}$

• More generally for connections between layer $L-1$ and layer L , $A^{L-1} = \begin{bmatrix} \vec{z}_1 & \vec{z}_2 \end{bmatrix}$

↑ features of one example

↑
batch size
Each row corresponds to all neurons in layer L-1
of certain example

$$\begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \end{bmatrix} \begin{bmatrix} x_{00} & x_{01} \\ x_{10} & x_{11} \\ x_{20} & x_{21} \end{bmatrix}^{L-1} + \begin{bmatrix} b_0 & b_1 \\ b_1 & b_2 \end{bmatrix}^L = \begin{bmatrix} z_{00} & z_{01} \\ z_{10} & z_{11} \end{bmatrix}^L$$

↑
Same as size
Size times of
neurons don't change
for example

↑
outputs in
layer L
for first example

$$o(z^L) = A^L$$

↑
activation
vector

The Weight Matrix - Fundamental Subspaces

• Basis vector j is all connections from neurons i to neuron j

• $C(W)$ contains all possible outputs of $W\vec{x} = \vec{z}$ before bias and activation

• Orthogonal initialization for W columns ensures network can map to larger, more varied set of inputs \Rightarrow can fit complex data better

• $N(W)$ shows what layers of activations give $\vec{z} = 0$ (which directions of input space are "ignored")

• High dimensional $N(W) \Rightarrow$ more of input space being ignored

• $R(W)$ related to how much all the features matter to each neuron j ; linearly independent rows \Rightarrow each neuron has different weights for features/inputs \Rightarrow each neuron encodes more info

• $N(W^T)$ less interpretable

• Optimizing $C(W)$ or $R(W)$ can lead to simpler, effective model

• W^T used for backpropagation

Cross-Entropy Loss

- Viewing why cross entropy is a good loss function for classification

e.g. Good prediction

$$\hat{y} = \text{softmax}(\text{logit})$$

Class 0	0.99	1
Class 1	0.01	0

$$-\sum_i y_i \log(\hat{y}_i) = -(1) \log(0.99) + (0) \log(0.01) = -\log(0.99) = 0.004 \leftarrow \text{low loss for good classification}$$

e.g. Bad prediction

$$\hat{y} = \text{softmax}(\text{logit})$$

Class 0	0.01	1
Class 1	0.99	0

$$-\sum_i y_i \log(\hat{y}_i) = -(1) \log(0.01) + (0) \log(0.99) = -\log(0.01) = -(-2) = 2 \leftarrow \text{high loss for bad classification}$$

- In ML, \log_e is commonly used to due its simple derivative, but the same properties as above apply

Backpropagation

↳ Backward propagation of errors

- Backpropagation: the algorithm for determining how a single training example would like to nudge the weights and biases in terms of what changes these parameters cause the most rapid decrease to the loss (calculates gradient)
 - Based on output you want, you know how to adjust the parameters in the layer before it
 - Prioritize adjusting the neurons whose activation is not close to what you want

To increase activation of neuron j in layer L ...

- Increase b_j
- Increase w_{ij} in proportion to a_i
 - Higher activations and weights complement each other in changing the output since $\sigma(a)$ for higher $w_{ij} \rightarrow$ higher $\sigma(a)$
- We care about which of these actions decreases the cost the most (most cost effect step)
 - Tiny change in w_{ij} changes loss the most $\Rightarrow \left| \frac{\partial L}{\partial w_{ij}} \right| = \text{High} \Rightarrow w_{ij} - n \frac{\partial L}{\partial w_{ij}}$ causes a bigger step (one reason why we use derivative for gradient descent)
- Doing this can "strengthen connections" / make it more likely for that output neuron to be activated for specific training example

• Weights can be negative

work

- Basically to make o_i brighter you can decrease activations for the negative weights and increase them for positive ones to have greatest effect on L
- But you have to take all o_i into account and how each of them want to change a_i
 - Add all these "wanted" changes are added together
- This is what propagating backwards gives you list of nudges to apply to layer $L-1$

• Each training example will have how it wants to nudge each weight and bias, so take the average

- This nudge is proportional to negative gradient (i.e. larger nudge means larger partial derivative)

$$\theta_i \leftarrow \theta_i - n \frac{\partial L}{\partial \theta_i} \Rightarrow \theta_i \leftarrow \theta_i + \Delta \theta_i \Rightarrow \Delta \theta_i = -n \frac{\partial L}{\partial \theta_i}$$

$$\frac{\Delta \theta_i}{n} = -\frac{\partial L}{\partial \theta_i}$$

↑ gradient = slope down = layer change

Using the chain rule to find the gradient

- Start by propagating/transmit the loss backwards throughout the network; it computes the partial derivatives layer by layer

Simple NN: $C(w_1, b_1, \dots, w_3, b_3)$

Loss function: MSE

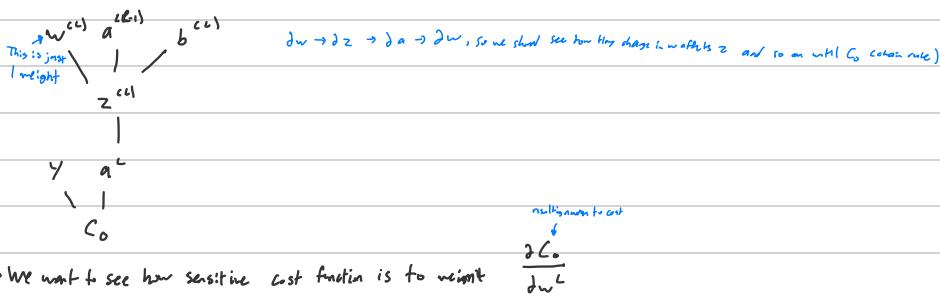


$$C_o(\dots) = (a^l - y)^2$$

↑
cost for single
training example (and
is avg of C_k)

$$a^l = \sigma(z^l)$$

$$z^l = w^l a^{l-1} + b^l$$



- We want to see how sensitive cost function is to weight

$$\frac{\partial C_o}{\partial w^{l+1}} = \frac{\partial z^{l+1}}{\partial w^L} \cdot \frac{\partial a^l}{\partial z^l} \cdot \frac{\partial C_o}{\partial a^l}$$

for specific training example

- Chain rule is much more efficient than plugging in $\sum w_i a_i + b$ backwards into each a_i for network to derive cost function then taking the partials
- Chain rule allows for reusing your previous computation, and it is recursive so it can be written more easily as an algorithm

$$\frac{\partial C}{\partial w^L} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^L}$$

↑
amount of C_k
1 training example; need to average over
all training examples

$$\frac{\partial C_o}{\partial a^l} = 2(a^l - y) \quad \text{Cost is proportional to error}$$

$$\frac{\partial a^l}{\partial z^l} = \frac{da^l}{dz^l} = \sigma'(z^l)$$

$$\frac{\partial z^l}{\partial w^L} = a^{l-1} \quad \text{amount the weight changes } z \text{ is dependent on } a^{l-1}$$

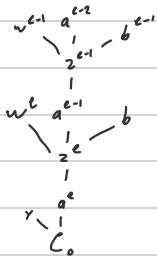
$$\frac{\partial C_o}{\partial w^L} = a^{l-1} \sigma'(z^l) 2(a^l - y)$$

$$\frac{\partial C_o}{\partial b^L} = \frac{\partial z^l}{\partial b^L} \cdot \frac{\partial a^l}{\partial z^l} \cdot \frac{\partial C_o}{\partial a^l} = \sigma'(z^l) 2(a^l - y)$$

$$\frac{\partial C_o}{\partial a^{l-1}} = \frac{\partial z^l}{\partial a^{l-1}} \cdot \frac{\partial a^l}{\partial z^l} \cdot \frac{\partial C_o}{\partial a^l} = a^{l-1} \sigma'(z^l) 2(a^l - y)$$

Propagating Backwards

Going Further Back in Network



$$z^l = w^l a^{l-1} + b^l \Rightarrow z^{l-1} = w^{l-1} a^{l-2} + b^{l-1} \Rightarrow \frac{\partial z^{l-1}}{\partial w^{l-1}} = a^{l-2} \text{ or } \frac{\partial z^{l-1}}{\partial w^l} = a^{l-2} \frac{\partial w^l}{\partial w^{l-1}}$$

⋮
and so on

making L single neuron layers

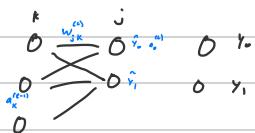
$$\frac{\partial C_0}{\partial w^{(l-1)}} = \frac{\partial z^{l-1}}{\partial w^{l-1}} \cdot \frac{\partial a^{l-1}}{\partial z^{l-1}} \cdot \frac{\partial z^l}{\partial a^{l-1}} \cdot \frac{\partial a^l}{\partial z^l} \cdot \frac{\partial C_0}{\partial a^l} = \sim$$

Process

Replace ℓ with $\ell-k$ in your equation → find unknown partial → take its average across all training examples in mini-batch → use formula for gradient descent to adjust parameters

With Indices

n_ℓ - number of neurons in layer ℓ



$$C_0 = \sum_{j=0}^{n_\ell} (a_j^\ell - y_j)^2 \quad z_j^\ell = \sum_{k=0}^{n_{\ell-1}} w_{jk}^\ell a_k^{\ell-1} + b_j^\ell \quad a_j^\ell = \sigma(z_j^\ell)$$

Weight indices feed loss backwards for backprop

$$\frac{\partial C_0}{\partial w_{jk}^\ell} = \frac{\partial z_j^\ell}{\partial w_{jk}^\ell} \cdot \frac{\partial a_j^\ell}{\partial z_j^\ell} \cdot \frac{\partial C_0}{\partial a_j^\ell}$$

↑ follow "path" from that weight to the cost

Same as before:

$$\frac{\partial z_j^\ell}{\partial w_{jk}^\ell} = a_k^{\ell-1} \quad \frac{\partial a_j^\ell}{\partial z_j^\ell} = \sigma'(z_j^\ell)$$

Different:

$\frac{\partial C_0}{\partial a_k^{\ell-1}}$: the activation of a neuron in another layer has multiple paths to output neurons, so you have to trace each of these outputs and how they each affect the cost function

This has multiple paths, so you have to sum of partials for each path to get total effect of $a_k^{\ell-1}$ on C_0 .

$$\frac{\partial C_0}{\partial a_k^{\ell-1}} = \sum_{j=0}^{n_\ell} \frac{\partial z_j^\ell}{\partial a_k^{\ell-1}} \cdot \frac{\partial a_j^\ell}{\partial z_j^\ell} \cdot \frac{\partial C_0}{\partial a_j^\ell}$$

Sum partials & account for total amount that a_j^l affects $z_j^l \rightarrow a_j^l \rightarrow \text{cost}$

$$\frac{\partial z_j^{e+1}}{\partial a_j^e} \cdot \frac{\partial a_j^{e+1}}{\partial z_j^{e+1}}$$

solved

$$\frac{\partial C_0}{\partial a_j^e} = 2(a_j^e - y_j) \quad \text{or} \quad \sum_{j=1}^{n_{e+1}} w_{jk}^{e+1} \sigma'(z_j^{e+1}) \frac{\partial C}{\partial a_j^{e+1}}$$

summing up all paths from activation neuron

\downarrow recursive until end with $2(a_j^e - y_j)$

s_j^e

$$\frac{\partial C_0}{\partial w_{jk}^{e+1}} = \frac{\partial z_j^{e+1}}{\partial w_{jk}^{e+1}} \cdot \frac{\partial a_j^{e+1}}{\partial z_j^{e+1}} \cdot \frac{\partial z_j^e}{\partial a_j^{e+1}} \cdot \frac{\partial a_j^e}{\partial z_j^e} \cdot \frac{\partial C_0}{\partial a_j^e}$$

$\frac{\partial C_0}{\partial a_j^e}$, already calculated

• Same idea of summing the paths until the end

$$\delta_j^e = \frac{\partial C_0}{\partial a_j^e} \cdot \frac{\partial a_j^e}{\partial z_j^e}$$

↑
summation step

is defined as the error of neuron j in layer e ; determines how much that neuron contributes to the error in output

Backpropagation Derivations

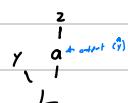
- Optin 1: Calculate loss for each training example \Rightarrow calculate ∇C_i for each training example \Rightarrow gradient is $\frac{1}{m} \sum_i \nabla C_i$ more common, less computationally expensive
- Optin 2: Calculate loss for each training example and average it with $-\frac{1}{m} \sum_i y_i \log(a_i) \Rightarrow$ find partial derivatives with this avg loss

Backpropagating Error from Loss \rightarrow Final Layer

Loss Function: Softmax categorical cross-entropy

$$\text{Loss} = -\sum_{k=1}^K y_k \log(a_k)$$

$$a_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}} = \frac{e^{z_k}}{e^{z_1} + e^{z_2} + \dots + e^{z_K}}$$



$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = ?$$

Single example

$$\frac{\partial a}{\partial z} = \frac{\partial}{\partial z} \begin{bmatrix} \frac{a_1}{\sum_{k=1}^K a_k} & \frac{a_2}{\sum_{k=1}^K a_k} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

$$\frac{\partial a_i}{\partial z_j}$$

\downarrow even when $i \neq j$ and mostly a_i depends only on z_i , softmax activation and all components of \vec{z}

$$\frac{\partial a_i}{\partial z_j} = \frac{\partial}{\partial z_j} \left(\frac{e^{z_i}}{\sum_{k \neq i} e^{z_k}} \right) = \frac{-e^{z_i} (e^{z_j})}{(\sum_{k \neq i} e^{z_k})^2} = \frac{-e^{z_i} (e^{z_j})}{(\sum_{k \neq i} e^{z_k})^2} = \frac{-e^{z_i}}{\sum_{k \neq i} e^{z_k}} \cdot \frac{e^{z_j}}{e^{z_k}} = -a_i a_j$$

$j = i$

$$\frac{\partial a_i}{\partial z_i} = \frac{\partial}{\partial z_i} a_i = \frac{\partial}{\partial z_i} \left(\frac{e^{z_i}}{\sum_{k \neq i} e^{z_k}} \right) = \frac{e^{z_i} (\sum_{k \neq i} e^{z_k}) - e^{z_i} (e^{z_i})}{(\sum_{k \neq i} e^{z_k})^2} = \frac{e^{z_i} - e^{z_i} + e^{z_i} \sum_{k \neq i} e^{z_k}}{(\sum_{k \neq i} e^{z_k})^2} = \frac{e^{z_i}}{\sum_{k \neq i} e^{z_k}} \cdot \frac{\sum_{k \neq i} e^{z_k}}{e^{z_i}} = a_i (1 - a_i)$$

$$\frac{\sum_{k \neq i} e^{z_k}}{\sum_{k \neq i} e^{z_k}} = \frac{\sum_{k \neq i} e^{z_k} - e^{z_i}}{\sum_{k \neq i} e^{z_k}} = 1 - \frac{e^{z_i}}{\sum_{k \neq i} e^{z_k}} = 1 - a_i$$

$$\frac{\partial a_i}{\partial z_j} = \begin{cases} -a_i a_j & i \neq j \\ a_i (1 - a_i) & i = j \end{cases}$$

$$\text{Loss} = -\sum_{k=1}^K y_k \log(a_k) = -[y_i \log(a_i) + \sum_{k \neq i} y_k \log(a_k)]$$

$$\frac{\partial L}{\partial z_i} = \frac{\partial a_i}{\partial z_i} \left(\frac{\partial L}{\partial a_i} [L] \right) = \frac{\partial a_i}{\partial z_i} \left[\frac{\partial}{\partial a_i} \left[-y_i \log(a_i) + \sum_{k \neq i} y_k \log(a_k) \right] \right] = \frac{\partial a_i}{\partial z_i} \left(-\frac{y_i}{a_i} - \sum_{k \neq i} \frac{y_k}{a_k} \right) = -\frac{\partial a_i}{\partial z_i} \cdot \frac{y_i}{a_i} - \frac{\partial a_i}{\partial z_i} \cdot \sum_{k \neq i} \frac{y_k}{a_k} = -\frac{\partial a_i}{\partial z_i} \cdot \frac{y_i}{a_i} - \sum_{k \neq i} \frac{\partial a_k}{\partial z_i} \cdot \frac{y_k}{a_k}$$

$$= \frac{-a_i (1 - a_i) y_i}{a_i} - \sum_{k \neq i} \frac{-a_k a_i y_k}{a_k} = -(1 - a_i) y_i + \sum_{k \neq i} a_i y_k = -y_i + a_i y_i + a_i \sum_{k \neq i} y_k = -y_i + a_i (y_i + \sum_{k \neq i} y_k) = -y_i + a_i (\sum_{k \neq i} y_k) = -y_i + a_i = a_i - y_i$$

$$\delta_i^L = \frac{\partial L}{\partial z_i} = a_i - y_i \Rightarrow \frac{\partial L}{\partial z_i} = \delta_i^L \cdot \vec{a}_i - \vec{y} \quad \text{or } \frac{\partial L}{\partial z_i} = \frac{\partial L}{\partial z_i} = A_i - Y \quad \text{for batches}$$

$$\delta_i^L = \nabla_i L \circ \sigma'(z_i^L) = \begin{bmatrix} \frac{\partial L}{\partial z_1} & \frac{\partial L}{\partial z_2} & \dots \\ \frac{\partial L}{\partial z_n} & \frac{\partial L}{\partial z_n} & \dots \end{bmatrix}$$

Explains the L O(i) $\xrightarrow{\text{if true}}$ $\frac{\partial L}{\partial z_i}$

Tips for intuition

- Just as $W\vec{a}_i$ maps info from layer $i-1$ to i , \vec{w}^T maps error backwards from layer i to $i-1$
- This is because $w = j \sum_{k \in K^{i-1}}^K \vec{a}_k$ and $w^T = i \sum_{j \in J^i}$

- Some formulas have A^T and O both so matrix shape output is as expected and so the derivatives are scaled properly

E.g.

$$f'(z) = \frac{\partial z}{\partial z_i} \text{ is low} \Rightarrow \text{neuron activation not that sensitive to input} \Rightarrow \text{not as important to take steps in weight direction with this derivative} \Rightarrow \text{element-wise multiplication of } \frac{\partial L}{\partial \vec{a}_i} \odot \frac{\partial z}{\partial \vec{a}_i}$$

will be far corresponding derivatives \Rightarrow smaller steps in certain direction

- Hadamard product also correctly applies chain rule with vectors/matrices

For NN with 1 hidden layer



- For batches use matrices and scale down the forward loss with $\frac{1}{m}$ for W and B ; don't scale down $\frac{\partial \text{Cost}}{\partial z}$ since it is intermediary

$$\text{Activation function} = f \quad z_i = w_i A_{i-1} + b_i \quad A_0 = X$$

$$\frac{\partial z_i}{\partial w_i} = \frac{\partial}{\partial w_i} [w_i A_{i-1} + b_i] = A_{i-1} \quad \frac{\partial z_i}{\partial b_i} = 1$$

By previous page, $\frac{\partial L}{\partial z_i} = \vec{a}_i - \vec{y}_i$. For single example for last layer $\Rightarrow \frac{\partial L}{\partial z_2} = \vec{a}_2 - \vec{y}$

$$\frac{\partial L}{\partial w_2} = \underbrace{\frac{\partial L}{\partial z_2}}_{\substack{n \times n \\ \text{hidden layer 2} \\ \text{of network}}} \cdot \underbrace{\frac{\partial z_2}{\partial \vec{a}_2}}_{\substack{(n_2 \times 1) \\ \text{and example}}} \cdot \underbrace{\frac{\partial \vec{a}_2}{\partial w_2}}_{(n_2 \times n_1)} = \underbrace{(\vec{a}_2 - \vec{y})}_{(n_2 \times 1)} \underbrace{\vec{a}_1^T}_{(n_1 \times 1)}$$

$$\text{For a batch, } \frac{\partial \text{Cost}}{\partial w_2} = \underbrace{\frac{1}{m} (\vec{A}_2 - \vec{Y})}_{(n_2 \times n_1)} \underbrace{\vec{A}_1^T}_{(n_1 \times m)}$$

$$\frac{\partial L}{\partial \vec{a}_2} = \underbrace{\frac{\partial L}{\partial z_2}}_{\substack{\text{different sample} \\ \text{without a} \\ \text{different } \vec{a}_2}} \cdot \underbrace{\frac{\partial z_2}{\partial \vec{a}_2}}_{\substack{\text{different } \vec{a}_2 \\ \text{and columns}}} = d\vec{z}_2 \cdot \frac{\partial z_2}{\partial \vec{a}_2} = \underbrace{\frac{\partial L}{\partial z_2}}_{\substack{\text{1 training example}}} \cdot \underbrace{\frac{\partial z_2}{\partial \vec{a}_2}}_{\substack{\text{1 training example}}} = d\vec{z}_2 = \vec{a}_2 - \vec{y}_2$$

$$\cancel{\frac{\partial \text{Cost}}{\partial B_2} = \frac{1}{m} d\vec{z}_2} \quad \begin{array}{l} \text{different sample} \\ \text{without a} \\ \text{different } \vec{B}_2 \end{array} \quad \begin{array}{l} \text{different } \vec{B}_2 \\ \text{and columns} \end{array}$$

$$\checkmark \frac{\partial \text{Cost}}{\partial B_2} = \frac{1}{m} \sum_{m=1}^M (d\vec{z}_2, 1) \quad \begin{array}{l} \text{sum along length 1 (2nd dimension - columns); sums all columns and average them} \\ \text{m=1 to M} \end{array}$$

Now we can sum and take avg to find average $d\vec{B}$ across batch and get correct size

Why $\frac{1}{m}$?

$A^T A \vec{a}_i = a_i^T A^T \vec{a}_i = \vec{a}_i^T \vec{a}_i$. Basically any matrix $A\vec{x}$ involves summing with all elements (samples of A), so we should balance this out with $\frac{1}{m}$

$$\left[\begin{array}{c|c|c} a_1 & a_2 & a_3 \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right]$$

Output layer \rightarrow Hidden layer

$$\text{Layer 1: } \delta_i^{(1)} = \frac{\partial L}{\partial z_i^{(1)}} = \frac{\partial L}{\partial a_i^{(1)}} \cdot \frac{\partial a_i^{(1)}}{\partial z_i^{(1)}} = \frac{\partial L}{\partial a_i^{(1)}} \cdot \frac{\partial z_i^{(1)}}{\partial z_i^{(1)}} \stackrel{\text{relu}}{=} \frac{\partial L}{\partial a_i^{(1)}} \cdot w_2 \cdot f'(z_i^{(1)})$$

Single neuron in single example

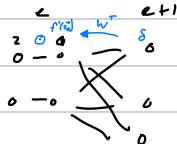
• This works for a single z_j ; if for \vec{z}^c (a vector of neurons), we need to adjust the formula to properly apply the chain rule

$$d\vec{z}_2 = \vec{\delta} = \frac{\partial L}{\partial \vec{z}_2} = \frac{\partial L}{\partial a_2} \cdot \frac{\partial a_2}{\partial \vec{z}_2} = \frac{\partial L}{\partial a_2} \cdot \frac{\partial \vec{z}_2}{\partial \vec{a}_1} \cdot \frac{\partial \vec{a}_1}{\partial \vec{z}_2} \stackrel{\text{relu}}{=} d\vec{z}_2 \cdot w_2 \cdot f'(\vec{z}_2) \quad X$$

$$\begin{array}{c} \text{adjust for proper} \\ \text{matrix size} \end{array}$$

$$= (W_2^T d\vec{z}_2) \odot f'(\vec{z}_2) \quad \checkmark$$

$$\text{Layer } c+1: \delta^{c+1} = \frac{\partial L}{\partial \vec{z}_c} = d\vec{z}_c^T W_{c+1}^T \vec{\delta}^{c+1} \odot f'(\vec{z}_c)$$



• Hadamard product properly applies chain rule (no matrix multiplication since \vec{z} and \vec{a} are not fully connected) after W^T maps the error backwards

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \vec{z}_1} \cdot \frac{\partial \vec{z}_1}{\partial a_2} \cdot \frac{\partial a_2}{\partial \vec{z}_2} \cdot \frac{\partial \vec{z}_2}{\partial \vec{a}_1} \cdot \frac{\partial \vec{a}_1}{\partial w_i} = \frac{\partial L}{\partial \vec{z}_1} \cdot \frac{\partial \vec{z}_1}{\partial w_i} \stackrel{\text{Single example}}{=} (W_2^T d\vec{z}_2) \odot f'(\vec{z}_1) \cdot \vec{a}_0^T \quad \text{or } d\vec{z}_1 \vec{a}_0^T$$

$$\frac{\partial \text{Cost}}{\partial w_i} = \frac{1}{m} d\vec{z}_1 A_0^T \quad \text{Batch}$$

$$\cdot \text{ Generally, } \frac{\partial L}{\partial w_i} = \vec{\delta}^c (\vec{a}^{c+1})^T$$

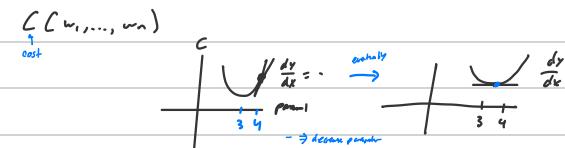
$$\frac{\partial L}{\partial b_i} = \frac{\partial L}{\partial \vec{z}_1} \cdot \frac{\partial \vec{z}_1}{\partial \vec{b}_1} = \frac{\partial L}{\partial \vec{b}_1} = d\vec{z}_1 \quad \text{Generally, } \frac{\partial L}{\partial b_i} = \vec{\delta}^c$$

$$\frac{\partial \text{Cost}}{\partial \vec{z}_1} = \frac{1}{m} \sum (d\vec{z}_1, 1)$$

• No need to backpropagate error to input layer

Optimizing the Loss Function with Gradient Descent.

- Gradient descent is how the model knows the amount it should update its parameters
- Take average of $\sum (\hat{y} - \hat{y})^2$ for all training data to get loss/cost \rightarrow MSE
 - Output of cost function dependent on all the parameters to neural network (hyperparameters aren't direct input)
- Computing minimum is computationally expensive
 - using single variable
 - unless you have identity risk
- For each step, find the minimum on that step \rightarrow and go that size size, you usually can only converge to a local minimum
- We want to minimize the cost, so we can figure out which step to take (which raw parameters to take on)



- Similar for higher dimensions; for partial $\frac{\partial f}{\partial y}$ look at y-axis plane and take single variable derivative from there

$$\frac{\partial f}{\partial y} = \lim_{\Delta y \rightarrow 0} \frac{f(x, y + \Delta y) - f(x, y)}{\Delta y}$$

- You are isolating the change to one variable at a time

- Backpropagation is the algorithm that computes the gradient

- Learning rate is the step size

$$\theta_t = \theta_{t-1} - \eta \frac{\partial L}{\partial \theta_{t-1}}$$

or $\theta_t = \theta_{t-1} - \eta \frac{\Delta L}{\Delta \theta}$

Makes sense: $\frac{\partial L}{\partial \theta_{t-1}} = \pm \rightarrow$ You want to decrease the parameter so you subtract it by this positive derivative

Also written as

$$\begin{aligned} W &= W - \eta dW = w_{ij} - \eta \frac{\partial L}{\partial w_{ij}} \\ b &= b - \eta db = b_i - \eta \frac{\partial L}{\partial b_i} \end{aligned}$$

- You have a weight matrix and a bias vector for each pair of layers in a neural network

- for all of the training data
 - Depends on type of gradient descent
- Feed forward info \rightarrow calculate cost $\rightarrow C(\hat{w})$ is "starting point" \rightarrow calculate $\frac{\Delta C}{\Delta w}$ \rightarrow Update parameters which causes a "step" \rightarrow stop after certain number of iterations or convergence

$$-\nabla C(\hat{w}) = \begin{bmatrix} 2 \\ -2 \\ 0.21 \\ \vdots \\ 1 \end{bmatrix}$$

all gradients

- step to descend to minima
+ this weight change doesn't matter as much since that weight is close to a minimum

- Higher partial (just like normal derivative) means that the cost is very sensitive to that weight changing

- Doesn't seem like there's any pattern for changing the weights/biases

- Batch Gradient Descent (BGD) - take 1 "more accurate" step per epoch (batch); converges slowly

- Stochastic gradient descent (SGD) - rather than computing the function and gradient after using all training data you do it after each training example (lot of noise)

Mini-batch gradient descent

- Full batch gradient descent involves passing X through, calculating loss with Y , calculating derivatives for each example and averaging, then taking one step

For computational speedup, we can divide X and Y into X_{mini} and Y_{mini} , where $m = \text{batch size}$

$$\frac{\text{num-train-samples}}{m} = n \text{, the number of mini batches, where } t \in \{1, \dots, n\}$$

- Per epoch, do * with mini batch matrices, then move to next mini batch matrices and repeat until the full batch of training examples have been used

* More steps per epoch \Rightarrow less epochs \Rightarrow faster convergence \Rightarrow SGD is more efficient

$$\nabla \hat{L}(\theta_t) = \nabla L(\theta_t) + \varepsilon_t$$

↑
estimated gradient of
batch size t
↑
noise

Variance of noise: $\text{Var}[\varepsilon_t] \propto \frac{1}{B}$

increasing
importance

If at local minimum...

$$\theta_{t+1} = \theta_t - \eta (\nabla \hat{L}_g(\theta_t)) = \theta_t - \eta (\nabla L(\theta_t) + \varepsilon_t) \quad r = \theta_t - \eta (\varepsilon_t)$$

local minimum $\Rightarrow \nabla L(\theta_t) = 0$

- For full batch size, $\theta_{t+1} = \theta_t$

- This is why the noise of smaller batch sizes can "kick" you out of local minima

Theoretically, increase batch size \Rightarrow more accurate avg of each partial \Rightarrow more accurate true gradient of loss function \Rightarrow can take a bigger, more confident step \Rightarrow safer to have high learning rate

- However, ? batch size \Rightarrow less exploration of loss landscape

Using SGD in PyTorch with differing batch sizes determines the gradient descent used

Convex function - only one global minima and no local minima

- You ideally want convex loss function

Vanishing gradient problem - Activation function has small derivative (e.g. sigmoid) \Rightarrow small numbers in chain rule being multiplied \Rightarrow small $\frac{\partial \mathcal{E}}{\partial \theta_{ij}} \Rightarrow$ small step size \Rightarrow learns very too slow

Exploding gradient problem - Activation function has small derivative (e.g. ReLU) \Rightarrow large numbers in chain rule being multiplied \Rightarrow large $\frac{\partial \mathcal{E}}{\partial \theta_{ij}} \Rightarrow$ large step size \Rightarrow may skip minimum

Other Optimization Algorithms

Exponentially weighted moving average (EWMA)

- EWMA - "Exponentiating β " results in prioritizing newer data points/moving average because each average takes a step

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t \quad \text{calculates moving average of } \theta \quad v_0 \text{ is arbitrary, non-zero for EWMA}$$

EWMA at time $t+1$ current value of a parameter

- β ↑ weight to recent observations; $0 \leq \beta \leq 1$ usually 0.9

This applies to all parameters and can be stored in a vector \vec{v}_{t+1}

$$v_0 = 0$$

$$v_1 = (1-\beta) \theta_1^*$$

$$v_2 = \beta v_1 + (1-\beta) \theta_2 = \beta(1-\beta) \theta_1 + (1-\beta) \theta_2 = (1-\beta)(\beta \theta_1 + \theta_2)$$

$$v_3 = \beta v_2 + (1-\beta) \theta_3 = (1-\beta)(\beta^2 \theta_1 + \beta \theta_2 + \theta_3)$$

lower weight &
old data point

$$v_6 = (1-\beta) \left(\sum_{i=0}^5 \theta_i \cdot \beta^i \right)$$

- * is too low if β is high (and since there is no average), can mess up future calculations

• Using **bias correction**, compute $\frac{v_t}{(1-\beta)^t}$ instead

• low $t \Rightarrow$ higher start (in denominator dividing)

• higher $t \Rightarrow \lim_{t \rightarrow \infty} \frac{v_t}{(1-\beta)^t} = v^* \Rightarrow$ negligible effect

SGD with Momentum

- Smooths out the zig-zag path
- In regions with consistent gradients you subtract mean of these "velocities" to get bigger one in same direction, leading to bigger step and faster convergence
- If loss landscape has gradients oscillating frequently (opposite signs), the previously averaged gradients result in taking a smaller step

$$W = W - \eta V_{dw}$$

EWMA of dw

$$w = w - \eta V_{dw}$$

$$B = B - \eta V_{db}$$

t iterations

$$V_{dw_t} = \beta * V_{dw_{t-1}} + (1-\beta) dw \quad dw \text{ means EWMA of } dw, S_{dw} \text{ means } (dw)^2; S \text{ starts for second}$$

$$V_{db_t} = \beta * V_{db_{t-1}} + (1-\beta) db$$

- EWMA will cancel out $\beta \delta^2 \eta$ due to oscillations and converge faster

- No momentum for the momentum because that has unnecessary complexity and isn't really very beneficial

Root Mean Square Propagation (RMSprop)

- Without RMS large step taken in steeper planes while smaller steps for smaller gradients
- Uses an adaptive learning rate
 - For large gradients, "reduce" η to avoid overshooting
 - For small gradients, "increase" η to converge more quickly

- Involves squaring gradient components to focus on magnitudes
- Uses exponentially weighted moving average to prioritize more recent gradients

$$W = W - \eta \frac{dW}{\sqrt{S_{dW} + \epsilon}}$$

$$S_{dW_t} = \beta S_{dW_{t-1}} + (1-\beta) (dW)^2$$

$\frac{\partial L}{\partial w}$ / S_{dW}

Exponential moving average
decay rate

$$\beta = \beta - \eta \frac{d\beta}{\sqrt{S_{d\beta} + \epsilon}}$$

for usually 10^{-3}
points to same denominator
for β don't have big steps
our overshoot

$$S_{d\beta_t} = \beta S_{d\beta_{t-1}} + (1-\beta) (d\beta)^2$$

initially
at $t=1$

$$W = W - \frac{\eta}{\sqrt{S_{dW} + \epsilon}} \nabla L(w)$$

adaptive η

$$S_{dW_t} = \beta S_{dW_{t-1}} + (1-\beta) (\nabla L(w))^2$$

- $\eta dW \propto S_{dW} \propto \frac{dW}{\sqrt{S_{dW}}}$
- Steep direction \Rightarrow large $\frac{dL}{dW}$ \Rightarrow high moving average \Rightarrow smaller step in this axis
- RMSprop focuses on gradient magnitude/step size while momentum focuses on the sign/direction of step

Adam (Adaptive Moment Estimation)

$$W = W - \frac{\eta}{\sqrt{S_{dW} + \epsilon}} V_{dW}$$

$\frac{\partial L}{\partial w}$
momentum

$$V_{dW_t} = \beta + V_{dW_{t-1}} + (1-\beta) dW$$

$$S_{dW_t} = \beta S_{dW_{t-1}} + (1-\beta) (dW)^2$$

$\frac{\partial L}{\partial w}$
decay rate

$$\beta = \beta - \eta \frac{d\beta}{\sqrt{S_{d\beta} + \epsilon}} V_{d\beta}$$

- β_1 is the momentum parameter, β_2 is the RMSprop parameter
- V_{dW_t} could m_t , the first moment, and S_{dW_t} is second moment
 - Moment captures characteristic of gradient distribution overtime

Both of these use bias correction so they don't favor zero

$$\sigma m_{1,t} = \frac{m_t}{1-\beta_1^t}$$

- Adam usually better than SGD with momentum or RMS prop alone

Autograd engine computes tensor gradients; it traverses through computation graph

Learning Rate Optimizer/Scheduler

- Optimization algorithms like Adam work to adjust learning rate differently for each parameter (Local adjustments)
- Adam uses moments for the momentum and RMSprop to adjust the step size taken in direction of each parameter
- There are also changes that should apply globally to the step size at certain points in the training process
 - Learning rate schedulers will change the learning rate efficiently over n epochs to converge faster and not get stuck in local minima/saddle points
 - Can complement optimization algorithms well

• Step decay - learning rate reduced by fixed amount (γ) every n epochs (step size)

- Makes sure you don't overshoot and take smaller steps for all parameters as they get closer to minimum

- $\gamma = 0.1 \Rightarrow$ learning rate is 10% of its previous value

$$\lfloor L \times \underbrace{J}_{\text{epoch \#}} \rfloor = \text{floor}(x)$$
$$n_n = n_0 \times \gamma^{\lfloor \frac{n}{\text{step-size}} \rfloor}$$

Say $n_0 = 1$ and step-size = 10

$$\text{Epoch 0-9: } n_9 = 1 \cdot (0.1)^{\lfloor \frac{9}{10} \rfloor} = 1 \cdot (0.1)^0 = 1$$

↓ 10% of previous value

$$\text{Next 10 epochs: } n_{11} = 1 \cdot (0.1)^{\lfloor \frac{11}{10} \rfloor} = 0.1$$

• Exponential decay decreases exponentially over epochs

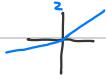
Functions

$$\text{ReLU}(z) = \max(0, z) = \begin{cases} 0 & z \leq 0 \\ z & z > 0 \end{cases} \quad \text{relu}'(z) = \begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

• Most people use ReLU instead of sigmoid now because it is simple and $\theta'(z) = 1$, which mitigates vanishing gradient problem

◦ Dying ReLU problem: many negative $z \Rightarrow \theta'(z) = 0 \Rightarrow$ impacts parts of previous layers/parameters \Rightarrow stops learning/stops

To solve this, leaky ReLU $= \max(\alpha z, z) = \begin{cases} \alpha z & z \leq 0 \\ z & z > 0 \end{cases}$ $\theta'(z) = \begin{cases} \alpha & z \leq 0 \\ 1 & z > 0 \end{cases}$



◦ Small $\alpha \Rightarrow$ small $\theta'(z)$ \Rightarrow neurons don't die

$$\text{Cross-entropy loss } (y, \hat{y}) = - \sum_{i=1}^C y_i \log(\hat{y}_i) \quad \text{only term left is } \frac{y_i}{\hat{y}_i} \text{ of the correct classification}$$

severely punishes this

◦ Cross-entropy good for classification because it focuses on correct class and wrong classification \Rightarrow high loss; just one activation effects cost, not all

◦ Because $\log_e(x)$ when $x \leq e$ is negative so you have to make the loss non-negative function to minimize it with gd

◦ e and \ln used because they simplify derivatives for chain rule

◦ For whole batch take average of the losses

$$\text{Mean Squared Error (MSE)} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

• For mini-batch GD, it is common to track average loss over each epoch instead of final loss; epoch loss defined as average batch loss in the epoch

◦ This can stabilize the loss resulting from step of earlier batch; avg-loss per epoch = $\frac{\text{total epoch loss}}{\text{num batches}} \leftarrow \Sigma \text{batch loss}$

◦ Even though average loss per epoch doesn't go down as fast as loss at end of epoch, it can be more robust, smooth, and generalizable results can be generalized on whole training data not just last batches of epoch

Softmax activation function

◦ For output neurons before computing loss they're passed through softmax function that maps results between 0 and 1 (probabilities)

◦ The raw "scores" passed to softmax function are called logits (not related to stats)

$$\text{probability of this classification from output layer} = \text{softmax}\left[\begin{matrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{matrix}\right] = \frac{e^{z_1}}{\sum_{j=1}^n e^{z_j}} \quad \text{softmax value}$$

◦ Ensures the outputs add to 1, larger values are closer to 1; $(\text{positive number})^n > 0$, ensures all the probabilities are positive

◦ Sigmoid function, $\sigma(z) = \frac{1}{1+e^{-z}}$ does this more efficiently for binary classification

◦ Sigmoid does not rely on other neuron outputs from the output layer

$$\text{E.g. } \vec{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$\text{Softmax}\left(\begin{bmatrix} z_1 \\ z_2 \end{bmatrix}\right) = \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} \quad p_1 = \text{prob of positive class}$$

$$\text{Sigmoid}(z_1) = p_1$$

Convolutional Neural Network (CNN)

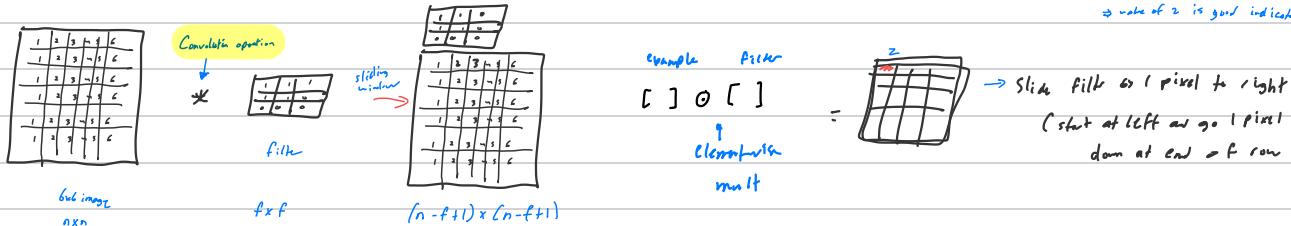
- Useful for image recognition
- Large images \Rightarrow many features \Rightarrow too many parameters \Rightarrow computationally expensive and overfitting

• **Filter/condition Kernel** - sliding window consisting of certain H of pixels for detecting a feature/pattern

- Each  is a parameter

Low bias \Rightarrow good raw input \Rightarrow each z_i is good \Rightarrow good weights in filter

\Rightarrow value of z is good indicator of feature



• The convolution operation produces one pixel of output feature map at a time

- Convolution layer directly takes in an image as input, no need to flatten

• Different filters detect different things (discovered by mathematicians);

- Edge detection; transpose of horizontal filter is the vertical filter

• Many filters in each layer of the neural network that each produce an image

$$(n \times n \times 1) \times (f \times f \times c) \rightarrow (n-f+1) \times (n-f+1) \times c \quad \text{when stride=1}$$

grayscale \uparrow
size \uparrow
number of filters \uparrow
number of images produced / per filter

• Colored image is RGB, $(n \times n \times 3)$ size \Rightarrow filter is $f \times f \times 3$

- Same sliding window concept applies except the sliding window is a cube (to output one feature map element)
- Applies to higher dimensions

• Each element of filter is a parameter that is trained

• Each convolutional layer processes outputs/feature maps of previous layer

- Each layer can detect different/more complex features

• After convolution, bias is added to raw image and you pass through activation function

Padding

- Each convolution reduces image size
 - If you have too many layers where the output is passed to another filter, it may get too small and lose info
- Certain pixels (like corner) only pass through filter once
- Padding - surrounding image with certain width of zeros
- Valid convolution - no padding
- Same convolution - image $\xrightarrow{\text{convolve}}$ feature map of same size since $n^i \cdot f + 1 = n$

$$n^i = n + 2p_{\text{padding}}$$

$$n^i \cdot f + 1 = n$$

$$n + 2p_{\text{padding}} \cdot f + 1 = n$$

$$p = \frac{cf-1}{2}$$
 must be true for same convolution $\Rightarrow f > p$ to be whole number f must be odd $\Rightarrow f$ usually odd for same convolution

Strides

- Stride - the amount of pixels the filter moves right/down
- Convolution called **strided convolution**
- This could discard all the pixels at the bottom
- Size of final image = $\text{Floor} \left[\frac{n-f}{s} + 1 \right] \times \text{Floor} \left[\frac{n-f}{s} + 1 \right]$
- If you double the stride it halves the dimension of output feature map

Pooling Layer

- Purpose of pooling layer is to reduce dimensions of image while preserving features
 - Smaller image \Rightarrow less computational cost
 - Enhances/sharpen features

Image \rightarrow convolutional layer \rightarrow pooling layer (can have 1 or more convolutional layers)

- For max pooling matrix $n = \text{stride}$ so pooling windows don't overlap \Rightarrow no redundancy in extracted features
 - Max pooling - takes max value in the sliding window
 - Average pooling - takes average of values in sliding window
 - Doubling stride halves n of resulting feature map
- No parameters involved in pooling (no training required)

Multilayer Perceptron

- The end of a CNN has a fully connected (FC) multilayer perceptron
- Fully connected layer - every neuron in layer C connects to every neuron in layer $C+1$
- All transformed feature maps from convolving/pooling flattened into a vector \rightarrow passed as input to NN
- This associates the extracted features to a label
- Conv + Pooling = 1 layer
- Backprop very similar except now you have filter parameters

Backprop

$X = \text{Image}$

$K = \text{Filter / kernel}$

$B = \text{Bias}$

$$X^{\text{image}} \\ K^{\frac{26}{26}} \quad Z^{\text{analytic}} \quad P^{\text{MLP - classification}} \\ K^{\frac{26}{26}}, B^{\frac{26}{26}}$$

- Since the filter in one section of image is an element of Z , you can think of that element of Z , once CNN trained, as a "feature" detector for that part of the image
- Adjusting the weights in the filter so Z can be a feature detector

Confusion Matrix and Model Analysis

- Element c_{ij} of confusion matrix C is equal to the number of examples labeled as class i predicted to be in class j

- Macro averaging** - involves calculating metric for each class independently then taking average ; $\frac{1}{\text{# classes}} \sum$ stat;
→ Better for treating classes equally regardless of size

- Micro averaging** - aggregate/sum up results from all classes at once to get metric representative of entire model; often for imbalanced classes
o Bigger categories have bigger impact; **weighted averaging** accounts for this for imbalanced classes also

$$\sum \frac{x_i}{x_i+y_i} \neq \frac{\sum x_i}{\sum x_i+y_i} \quad \text{If formula } = \frac{x_i}{x_i+y_i}, \text{ microaveraging is } \frac{\sum x_i}{\sum x_i+y_i}$$

- Support** - # of instances per class used in weighted averaging

- Multiclass score calculations applied per label and then some sort of average is taken

$TP+FN$ = All examples in positive class

$TN+FP$ = All examples in negative class

- Accuracy** - measures overall class of the model

• Binary

$$\text{Accuracy} = \frac{\text{correct classification of positive and negative class}}{\text{total}} = \frac{TP+TN}{TP+TN+FP+FN}$$

• Multiclass

$$\text{Accuracy} = \frac{\sum \text{# correct classification for class } i}{\sum \text{ total classifications for class } i}$$

focuses on just positive class, multiclass accuracy

- Precision / Positive Predictive Value** - measures accuracy of the positive predictions "how many minors that you classified as 2 are actually a 2"

◦ Binary

$$\text{Precision}_n = \frac{TP}{TP+FP} \quad \text{Precision} \Rightarrow \text{every example that the model classified as positive class was correct, but it may have misclassified other examples of that positive class (FN)}$$

◦ Multiclass

- Macro Precision** - Calculate precision for each class then take the average

$$\text{Precision}_{\text{macro}} = \frac{1}{\text{# classes}} \sum_{i=1}^{\text{# classes}} \text{Precision}_i$$

- Micro Precision** - Aggregate the contributions of all classes to compute average precision

$$\text{Precision}_{\text{micro}} = \frac{\sum TP_i}{\sum (TP_i + FN_i)}$$

• Recall / True Positive Rate / Sensitivity - measures ability of the model to correctly identify actual positives; out of all examples in positive class, what proportion is correctly identified

$TP + FN =$ Total examples labeled as positive class

• Binary

$$\text{Recall} = \frac{TP}{TP + FN}$$

Total classified by model
as positive correctly

• Multiclass

◦ Micro and macro recall

• Micro Precision = Micro Recall due to a FP in one class is FN in another class

• False Negative Rate / Miss Rate - complement of recall; out of all examples in positive class, what proportion is missed / incorrectly identified

$$FNR = \frac{\text{Incorrectly identified class}}{\text{Total # of class}} = \frac{FN}{TP + FN}$$

$$FNR + TPR = 1$$

• F1 Score - harmonic mean of precision and recall; provides a balance between the two

◦ Harmonic mean = $\frac{2ab}{a+b}$; gives more weight to smaller numbers; if a or b is small, the harmonic mean is also small

$$F_1 \text{ Score} = 2 \cdot \frac{(\text{Precision})(\text{Recall})}{\text{Precision} + \text{Recall}}$$

• In general F-measures evaluate performance of classification models by combining precision and recall with different weights; some parameterized by β

$$\circ F_\beta = (1+\beta)^2 \cdot \frac{\text{Precision} \times \text{Recall}}{\beta^2 \text{Precision} + \text{Recall}}$$

β is beta
 $\beta > 0$

• Multiclass

◦ Macro F1 Score and Micro F1 Score

• True Negative Rate / Specificity - out of all the examples in negative class, what proportion is correctly identified

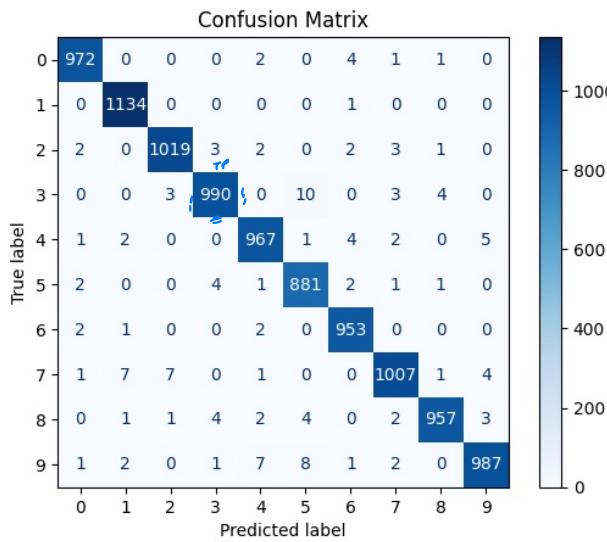
$$TNR = \frac{TN}{TN + FP}$$

• False Positive Rate - complement of TNR; out of all the negative class, how many are falsely identified

$$FPR = \frac{FP}{TN + FP}$$

$$TNR + FPR = 1$$

Example Calculations



• Per class, look at the column for accuracy and row for recall

Class 3

$$10,000 = \# \text{ total examples in dataset} = 3\bar{3} + \text{non } 3\bar{3} = (TP + FN) + (TN + FP)$$

$$\begin{aligned} TP &= 990 && \text{"how many actual 3\bar{3} did model correctly identify/classify"} \\ &\uparrow && c_{ii} \text{ in diagonal} \end{aligned}$$

$$\begin{aligned} FN &= 3 + 10 + 3 + 4 = 20 && \text{"how many 3\bar{3} did the model fail to identify"} \\ &\uparrow && \text{the row} \end{aligned}$$

$$TP + FN = 1010 \quad \text{"how many 3\bar{3} were in the dataset"}$$

$$\begin{aligned} FP &= 1 + 4 + 4 + 3 = 12 && \text{"how many non 3\bar{3} were incorrectly identified as a 3\bar{3}"} \\ &\uparrow && \text{column} \end{aligned}$$

$$TN = Total - (TP + FN + FP) = 8,978 \quad \text{"how many non 3\bar{3} were correctly identified as non 3\bar{3}"} \quad \text{doesn't mean they were classified as correct number}$$

$$FP + TN = 8,990 \quad \text{"how many non 3\bar{3} are in the dataset"}$$

$$\text{Accuracy} = \frac{\text{Correct # of identifications of 3\bar{3} and not 3\bar{3}}}{\text{total # of 3\bar{3} and not 3\bar{3}}} = \frac{TP + TN}{TP + TN + FP + FN} = \frac{990 + 8,978}{10,000} = \frac{9,968}{10,000} = 0.9968$$

• Model good at making classification of 3 or not 3

$$\text{Precision}_3 = \frac{\text{Correct # of identifications of 3's}}{\text{Total identification of 3's}} = \frac{TP}{TP+FP} = \frac{990}{990+12} = 0.9880$$

- Of the examples the model classifies as 3, 98.8% are actually a 3 and 1.12% of those are false positives

- Model overpredicting class 3 \Rightarrow ↑ FP \Rightarrow ↓ Precision

$$\text{TAR: Recall}_3 = \frac{\text{Correctly identified 3's}}{\text{Total 3's}} = \frac{TP}{TP+FN} = \frac{990}{1010} = 0.9802$$

- Out of all the threes, the model correctly classifies 98.02% of them

$\text{FNR} = .0198$ Out of all 3's being classified by model 1.98% are incorrectly classified as not 3

- False Positive Rate - out of all the not threes, how many does the model falsely classify as a 3

$$\text{FPR} = \frac{FP}{TN+FP} = \frac{12}{8990} = 0.001$$

$\text{TNR} = 0.999$ model rarely misclassifies other numbers as 3

Not want important stuff sent to spam
(e.g. spam)

- We prioritize low number of false positives \Rightarrow low FPR = $\frac{FP}{TN+FP}$ \Rightarrow precision decreases \Rightarrow higher precision
- We prioritize low number of false negatives \Rightarrow low FNR = $\frac{FN}{TP+FN}$ \Rightarrow high recall

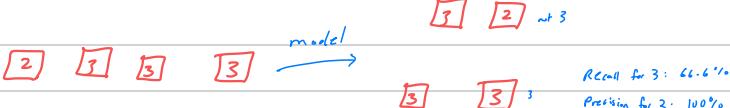
- For something like medicine you prioritize lower # of false negatives; for number classification it's not as important

$$\text{F}_1 \text{ Score} = \frac{2(\text{Precision})(\text{Recall})}{\text{Precision} + \text{Recall}} = 0.984$$

model has good balance of correctly identifying the number 3 (recall) and making sure its predictions are accurate (precision)

Ex.

Data:



Recall for 3: 66.6%
Precision for 3: 100%

- In addition to each class analysis, aggregate results for whole model performance with macro, micro, or weighted averaging

Colab/Metal

- Servers - Computer that provides resources, data, and services to other computers (clients) over a network

- If server breaks you are screwed

Cloud

- When you start something it is saved in various servers in data centers so there are backups
- When you "request resources" to do something from cloud infrastructure, you are given part of a CPU
- Hypervisor allocates certain amount of physical resources like CPU per person running

multiple computer devices

- TPU (Tensor Processor Unit) - a hardware accelerator unit (does something specific very fast)

- Open Colab → one HTTP request sent to servers and answered → request to cloud for TPU sent & allocate resources → I get results back

- TPU fast for parallelism, matrix multiplication

- Since it's very parallel you should mark end of calculation step (mark-step); Could otherwise lead to race conditions and poor memory management
- Has to warm up TPU so it establishes data paths so it runs quickly

operation starts when previous
are finished

- xm.xla-device() sets up communication between CPU and TPU

- .to(device) moves tensor from virtual machine memory (Copper & RAM of server) to TPU memory (or whatever backend)

CPU memory: registers, cache, RAM

TM memory: Scratchpad/SRAM, unifiler buffers, high bandwidth memory



CPU

- Metal API - software that directly interacts with the hardware (contains the functions to utilize Metal framework)

Backend - system that performs the actual computation (operations on tensors)

- You choose a device/backend to execute your code on

- This process uses the Metal API, which schedules how to send/start the GPU operations

Command Queue holds Command Buffers (sequences of commands)

- torch.device("metal") creates MTLDevice in Metal and creates these metal objects to transfer tensor

- PyTorch encodes commands for computation with Metal Compute pipeline

Compute Pipeline

- Compute shader/kernel - a program written in Metal Shading Language (MSL) for data-parallel GPU computations

- Metal Performance Shaders (MPS) include pre-optimized kernels for common ML operations like matrix multiplication

- Pipeline state object (PSO) properly configures shader

- Command buffer sequences with resources

- GPUs and TPUs benefit from large batch size since all examples in batch size passed through once at once (parallelism)

- Sometimes overhead of transferring data to GPU/TPU doesn't pay off (can mitigate by sending next batch to memory while calculating current one)

- Simple dataset → simple model with few layers/neurons → not many computations to be done in parallel → not much benefit from GPU/TPU

less overhead/latency
small weight matrices/bias nodes;
calculating less components in parallel

form of regularization

- Small dataset → small batch size to add noise to steps to prevent overfitting ⇒ less parallelism ⇒ less benefit from GPU/TPU

more frequent steps + faster convergence

Multiprocessing

- When getting an item from dataset, all transformations/preprocessing is applied
 - Dataloader loads batch from dataset → for each example the dataset's `__getitem__` method is called → transformations applied → data points batched together
 - data moved to GPU ↑ frameworks → move data to GPU
 - more of them ready with prefetch
 - does this at same time ↑ an advantage
- The goal is to make sure the GPU is being fed data and isn't idle while CPU preprocesses the data
- Page-locked/pinned memory - region of memory that OS guarantees will not be swapped out to disk
 - speeds up transferring data from CPU to GPU if tensors put here
 - Allows CPU to load/preprocess next batch of data while current batch being processed by GPU; good for asynchronous otherwise GPU is waiting
 - Metal makes sure this synchronization is handled properly and CPU doesn't give GPU more data until it's done
- Non-blocking = True means CPU to GPU data transfer is asynchronous
 - CPU will preprocess data of new batch while the previous batch is being transferred to GPU; CPU could technically return function values or move on before GPU is done calculating so you should synchronize
- 0 workers ⇒ data loading/preprocessing occurs in main process ⇒ the CPU core that runs the main script preprocesses the data
 - No inter-process communication
 - The main process always handles transfer to GPU
- > 0 workers ⇒ workers independently access data and utilize CPU cores ⇒ combine batch into a queue shared with the main process
 - GlobalWorkers: one keeps workers throughout all epochs rather than remaking them
 - Each worker process initialized with copy of dataloader and works independently
 - Queues a queue for inter-process communication (IPC)
- Pre-fetch factor - the amount of batches each worker process preloads for the main process to get
 - Workers need to serialize and deserialize dataset to create separate worker processes
 - Putting code in separate script makes it easier for worker processes to find the class definition
 - Warning: Dataloader creates workers and gives them the initial data to prepare and cache more data into memory; communication channels also have to be set up
- Global Interpreter Lock (GIL) prevents true parallelism since Python memory management isn't threadsafe; instead uses concurrency
 - Can also bypass by giving each process its own Python interpreter/GIL; multiprocessing module does this
 - taking turns
 - Slightly different than multithreading since it's not one process with multiple threads sharing memory

Varieties in Results

- Floating point operations (especially lots of them in neural network) can introduce tiny error in rounding over many iterations which can result in slight differences
 - when rounding
- Parallel computing can sometimes affect results since some floating point arithmetic is not associative (i.e. $(a+b)+c \neq a+(b+c)$ sometimes due to rounding)
- Even with setting seeds this error can still occur
 - Same data and code ↘ same errors

Compilers

- JIT compiler translates high-level code into machine code, and it can replace function call with whole function
 - First time you JIT-compile a function, it "traces" who initial parameters end up and caches it
 - In PyTorch, TorchScript has a JIT compiler
 - If you call function again with different parameters w/ same type(shape), it can use the cached machine code
 - Usually you just turn model to Torchscript model; more optimization seen
- Torchscript provides intermediate representation of model and can optimize computation graph, which is good for GPU
 - Graph-based representation of model means PyTorch can turn different operations into 1 kernel (more kernels \Rightarrow overhead to launch them on GPU)
 - Python computation graph marks dynamically (not at once) so how to tell what kernels can be combined
 - Removes parts of computation graph that don't effect output (dead code)
 - Reuses memory allocations
- Python interpreter reads and executes one line at a time; the bytecode doesn't map to CPU registers
 - Some code \rightarrow bytecode \rightarrow interpreted reads \rightarrow machine code \rightarrow executes
- **Module** is base class for all neural network layers and models contained in torch.nn module; classes that extend this own module
 - Can be a single layer or container of them (sequential)
 - Module contains parameters of model and methods to initialize them; also implements forward pass of network
 - Any nested Module in a larger module is called a **sub-module**