CNNs: Classification and Verification

Recitation 6 Jinhyung Park

Today we will talk about

- Problem Statement for Part 2
 - Classification and Verification
 - Transfer Learning
- Model Architectures
- Types of Losses
- Run through the starter notebook

Problem Statement

- Face Classification:
 - Classifying the person(ID) based on the image of the person's face
- Face Verification:
 - How would you use the same network you trained for Classification to do face verification, Ideas??
 - You identify the most important features in the image which capture the identity of the face
 - The extracted features will be represented by a fixed length vector, known as an embedding
 - In order to do verification, we need to identify if a given embedding is similar to a reference embedding of that person using a distance metric like the Cosine Distance

Classification vs Verification

- Classification
 - An N way classification task, predicting from a fixed set of possible output classes
- Verification
 - It is a matching operation, where you match the given sample to the closest sample from a reference of N other samples
 - Can also be a 1 to 1 task, where we want to verify if the two embeddings are similar (belong to the same class)
 - N way classification using cross entropy loss



Verification

25 · 50 · 75 · 100 · 125 ·

175

25 -50 -75 -100 -125 -150 -

175



Open set vs Closed Set

Classification versus Verification & Open versus Closed set for the facial recognition problem.



Training for Classification task



Training for Verification task



Transfer Learning

- Transfer Learning is a method where a model developed for a task is reused as the starting point for a model on another task
- Two ways to do Transfer Learning
 - Use trained model weights to initialize the network and train the entire network again
 - Freeze the weights of the network and fine tune the last few layers on the target dataset



CNN Architectures: AlexNet

Published in 2012 and was the winner of the Imagenet LSVRC-2012

- Uses ReLU as it makes training faster
- Implements dropout between layers
- Uses data augmentation methods (Random Crop, Random Flip in PyTorch)
- The network is trained using SGD with momentum



https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

VGG Net

VGG (2014) architecture reduces the size of each layer yet increases the overall depth of it. reinforces the idea that CNNs must be deep in order to work well on visual data. $224 \times 224 \times 3 = 224 \times 224 \times 64$

- Conv filters (3*3)
- Max Pool (2*2)

https://arxiv.org/pdf/1409.1556.pdf

224 x 224 x 3	224 x 224 x 64
	112 x 112 x 128
	56 x 56 x 256
	7 x 7 x 512
	1 x 1 x 4096 1 x 1 x 1000
201	C convolution+ReLU
	max pooling
	fully nected + ReLU
	softmax

ResNet

- Introduced in 2015, utilizes bottleneck architectures efficiently and learns them as residual functions
- Easier to optimize and can gain accuracy from increased depth due to skip connections



https://arxiv.org/pdf/1512.03385.pdf

Dense Nets

- It includes a stack of dense blocks followed by a transition layers
- Strengthen feature propagation, encourage feature reuse and decrease the number of parameters



https://arxiv.org/pdf/1608.06993.pdf

Grouped Convolutions



Depthwise Separable Convolution is a special case of Grouped Convolutions where g = 1

ShuffleNet is an extension to the Group Convolution which shuffles the output channels C'

Figure 0: Grouped Convolutions with $N_i = C$ input Channels, g groups, and $N_o = C'$ output Channels

MobileNet

Introduced in 2017, intended to run neural networks efficiently on mobile devices

- V1 introduces Depth wise separable convolution
- V2 introduces Inverted Residual block



Residual and Inverted Residual block



Figure 1: MobileNetV2 Architecture, BN is used after every conv and is omitted for brevity

ConvNext

ConvNeXt Block



- A 7x7 depth-wise convolution.
- A point-wise convolution increasing # of channels
- A point-wise convolution decreasing # of channels
- Residual Connection

Other Interesting Papers

- ResNeXt (2016)
 - <u>https://arxiv.org/pdf/1611.05431.pdf</u>
 - Generally a strict improvement to ResNet, but slower. It's like 3 lines of code changed.
- SENet (2017)
 - <u>https://arxiv.org/pdf/1709.01507.pdf</u>
 - Channel-wise attention in CNNs. It's like 20 lines of code.
- EfficientNet (2019)
 - <u>https://arxiv.org/pdf/1905.11946.pdf</u>
 - Optimized model scaling. Probably can hard code this with some effort.
- RegNet (2020)
 - https://arxiv.org/pdf/2003.13678.pdf
 - ResNet with optimized layer sizes. It's probably... 10 lines changed?
- ResNeSt (2020)
 - <u>https://arxiv.org/pdf/2004.08955.pdf</u>
 - ResNeXt on steroids + attention. I (we?) will be really impressed ©
- NFNet (2021)
 - <u>https://arxiv.org/pdf/2102.06171v1.pdf</u>
 - Quite doable actually

Discriminative Features

- Classification optimizes learning separable features
- Optimally we wish to learn discriminative features
 - Maximum inter class distance
 - Minimum intra class distance



Center Loss

- Tries to minimize the intra class distance by adding a euclidean distance loss term
- If you use this, YOU MUST USE CENTER LOSS FROM THE BEGINNING OF TRAINING CLASSIFICATION!

$$egin{aligned} \mathcal{L} &= \mathcal{L}_S + \lambda \mathcal{L}_C \ &= -\sum_{i=1}^m \log rac{e^{W_{y_i}^T oldsymbol{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T oldsymbol{x}_i + b_j}} + rac{\lambda}{2} \sum_{i=1}^m \|oldsymbol{x}_i - oldsymbol{c}_{y_i}\|_2^2 \end{aligned}$$

Softmax only. Left: training set. Right: test set.



Sommax + center loss. Left: training set. Right: test set.

https://github.com/KaiyangZhou/pytorch-center-loss

Triplet Loss



Here, x is the anchor image, x+ is the positive sample and x- is the negative sample. It is crucial to select challenging/hard negative samples to truly improve the model, known as **Hard Negative Mining**.

ArcFace

ArcFace calculates the geodesic distance on a hypersphere as opposed to Euclidean distance and tries to maximize the angular margin between two classes.



$$-\frac{1}{N}\sum_{i=1}^{N}\log\frac{e^{s*(\cos(\theta_{y_i}+m))}}{e^{s*(\cos(\theta_{y_i}+m))}+\sum_{j=1,j\neq y_i}^{n}e^{s*\cos\theta_{y_i}}}$$

where θ_j is the angle between the weight W_j and the feature x_i s - feature scale, the hypersphere radius m - angular margin penalty Class:

https://arxiv.org/pdf/1801.07698.pdf

Other types of Losses

- Contrastive Loss (maintain margin between classes)
- Pair Wise Loss (separate distributions of similarity scores)
- Angular Softmax Loss



https://github.com/adambielski/siamese-triplet/blob/master/losses.py

Regularization Techniques

- Label Smoothing
- DropBlock

Label Smoothing

- Deep Learning Models undergo a problem of Overfitting and Overconfidence.
- Label Smoothing is a technique that can help us solve the problem of Overconfidence.

What is Overconfidence?

For each sample, the model predicts outcomes with higher probabilities than the accuracy over the entire dataset.

This is a poorly calibrated model.

For example, it may predict 0.9 for inputs where the accuracy is only 0.6.

Example

Without Label Smoothing

Suppose we have K = 3 classes, and our label belongs to the 1st class. Logit Vector z = [a, b, c]Label vector y = [1, 0, 0] (one-hot encoded)

Gradient of Loss = softmax(z) – y Our model will make $a \gg b$ and $a \gg c$

z = [10, 0, 0]softmax(z) = [0.9999, 0, 0]

Label Smoothing

 $y_{ls} = (1 - \alpha) * y_{hot} + (\alpha / K)$

Example: With Label Smoothing

 $\alpha = 0.1$ $y_{ls} = [0.9333, 0.0333, 0.0333]$ This would result into the logits z = [3.3332, 0, 0] softmax(z) = [0.9333, 0.0333, 0.0333]

https://arxiv.org/pdf/1812.01187.pdf https://arxiv.org/abs/1812.01187 https://github.com/ankandrew/online-label-smoothing-pt

DropBlock

- Dropout usually works better with Fully Connected Networks
- Dropout has not proven to be useful in CNNs because of the spatial correlation between the activation outputs
- DropBlock is a regularization technique that has proven to be useful for CNNs
- It is a structured form of dropout that drops contiguous regions and not just random pixels





References

- https://arxiv.org/pdf/1512.03385.pdf
- https://arxiv.org/pdf/1608.06993v3.pdf
- https://arxiv.org/pdf/1409.1556.pdf
- <u>https://arxiv.org/pdf/1704.08063.pdf</u>
- https://arxiv.org/pdf/1503.03832v3.pdf
- <u>http://ydwen.github.io/papers/WenECCV16.pdf</u>
- <u>https://www.cs.cmu.edu/~rsalakhu/papers/oneshot1.pdf</u>
- <u>https://towardsdatascience.com/densenet-2810936aeebb</u>
- <u>http://yann.lecun.com/exdb/publis/pdf/hadsell-chopra-lecun-06.pdf</u>
- https://towardsdatascience.com/what-is-label-smoothing-108debd7ef06
- https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convoluti onal-neural-networks.pdf