

Instructions

This document is a template, and you are not required to follow it exactly. However, the kinds of questions we ask here are the kinds of questions we want you to focus on. While you might have answered similar questions to these in your project presentations, we want you to go into a lot more detail in this write-up; you can refer to the Lab homeworks for ideas on how to present your data or results.

You don't have to answer every question in this template, but you should answer roughly this many questions. Your answers to such questions should be paragraph-length, not just a bullet point. You likely still have questions of your own -- that's okay! We want you to convey what you've learned, how you've learned it, and demonstrate that the content from the course has influenced how you've thought about this project.

Project Name

Project mentor: Gordon Sun

Louise Lu ylu106@jhu.edu, Xihan Zhao xzhao77@jhu.edu, Zhichen Sha zsha2@jhu.edu Peiyuan Xu pxu11@jhu.edu

Github repo link: https://github.com/XavierZhao7/CS475_FacialExpressionRecognition

Outline and Deliverables

List the deliverables from your project proposal. For each uncompleted deliverable, please include a sentence or two on why you weren't able to complete it (e.g. "decided to use an existing implementation instead" or "ran out of time"). For each completed deliverable, indicate which section of this notebook covers what you did.

If you spent substantial time on any aspects that weren't deliverables in your proposal, please list those under "Additional Work" and indicate where in the notebook you discuss them.

Uncompleted Deliverables

1. "Would like to complete #2": Add attention mechanism to recognize partially occluded faces. (not enough time)
2. "Would like to complete #3": Introduce Faster R-CNN to improve facial recognition. (not enough time)
3. "Expected to complete #2": An extended neural network to classify both 7 basic expressions and 12 additional compound facial expressions. (completed classification of 7 basic expressions, but ran out of time to complete the classification of compound expressions.)

Completed Deliverables

1. "Must complete #1": An image reprocessing program that crops image around face areas and converts into gray-scaled JPEG format images.
2. "Must complete #2": A face capturing program that read in video and output face (if any) on specific frames.
3. "Must complete #3": A neural network to classify 7 basic facial expressions
4. "Would like to complete #1": Adaptation of the model to real-time video feed and create real-time classification output
5. "Would like to complete #2": Attained around 60% test accuracy for the classification using the extended neural network

Additional Deliverables

1. We included additional dataset from extended Cohn-Kanade facial expression database and create a new dataset.
2. In addition to simple CNN model presented in our presentation, we have improved our simple CNN and implemented Gabor-CNN and Gabor-VGG models for the solution.

Preliminaries

What problem were you trying to solve or understand?

We are trying to solve the problem of recognizing real-time facial expressions (both basic and compound emotions) from video clips and real-time colored images. This is a supervised learning classification problem where CNNs with RELU activation layers and various pooling filters are being applied.

What are the real-world implications of this data and task?

With the rise of Deep Learning technology, facial recognition has become an important application broadly used in our current society. Facial expression is used at many places, such as Zoom meetings, behavioral interviews, as well as for identity recognition at government facilities. However, during talks and behavioral interviews, it's hard to control people's emotions, especially when he or she is nervous. To address this situation, we plan to develop an application that can recognize real-time facial expressions from videos and photographs to aid preparation for the behavior interviews alone. Our application, if successfully, can be used by professional practicing for a behavioral interview, as well as professors and students and other working professionals giving presentations.

How is this problem similar to others we've seen in lectures, breakouts, and homeworks?

In lecture, we discussed Deep Learning, particularly with applications of image classification using Convolutional Neural Networks (CNN). Facial recognition is certainly an essential application of Deep Learning, and for this project we decided to use Convolution Neural Networks (trying different hyperparameters and architectures) on the real-world colored images of faces.

What makes this problem unique?

Researchers have recently been working on the facial expression recognition problem. However, our work expect to focus on predicting multi-class emotions (such as happily surprised, etc). We will also go beyond colored images to real-time video clips. Our project follows the FATE (Fairness, Accountability, Transparency, and Ethics) principles very closely and is very inclusive of race and gender. In other words, our application does not discriminate across race and gender.

What ethical implications does this problem have?

This problem does not seem to have ethical implications as of current state (in particularly, our application will certainly not discriminate across race and gender). However, as one of our feedback suggests, if this turns into application for preparations for behavioral interviews in the future, it is important to ensure that there is as little age, gender or racial bias in this program as possible.

Dataset(s)

Describe the dataset(s) you used.

We used real-world colored images from the Real-World Affective Face Database (RAF-DB) dataset. Within the RAF-DB dataset, there are two major subsets: one subset with single label each indicating 7 classe of basic facial emotions (e.g., Surprised, Fear, Disgust, Happy, Sad, Angry, and Neutral) and one with two labels each indicating compound emotions (such as Fearfully Surprised, Fearfully Disgusted, Sadly Fearful). We also combine the extended Cohn-Kanade (CK+) dataset with RAF-DB dataset to improve the accuracy. The CK+ dataset contains 593 video sequences from a total of 123 different subjects, ranging from 18 to 50 years of age with a variety of genders and heritage. The data are label with one of 7 classes: anger, neutral, disgust, fear, happiness, sadness, and surprise.

How were they collected?

The RAF-DB dataset is available online as they were pre-collected from publicly availble real-world images. The extended Cohn-Kanade dataset is also publicly avaialble online and they were obtained by experiementers of the study.

Why did you choose them?

Because they include both images and videos of facial expressions, there is also various gender, age or race to fight against potential bias in the datasets.

How many examples in each?

The RAF-DB dataset contains 29672 real-world colored images with 7 classes of basic emotions as labels and 12 classes of compound emotions. The CK+ dataset contains 593 video sequences from a total of 123 different subjects, ranging from 18 to 50 years of age with a variety of gender and heritage.

```
In [ ]: from google.colab import drive
        from tensorflow.keras.preprocessing.image import img_to_array, load_img
        import matplotlib.pyplot as plt
        import os
        from PIL import Image, ImageOps
        import numpy as np

        drive.mount('/content/gdrive/')
        %cd "gdrive/MyDrive/"
        %cd "CS475_FinalProject/"
        %cd "CS475_proj_Xihan/"
```

```
In [ ]: # Load your data and print 2-3 examples
        curr_directory = os.getcwd()
        import IPython.display as display

        ## Example 1
        print("Example of angry image from ck+ dataset")
        file_path = os.path.join(curr_directory, 'data/ck+/anger/S999_001_00000018.png')
        img1 = load_img(file_path)
        img1 = img_to_array(img1)
        im1 = Image.open(file_path)
        display.display(im1)

        ## Example 2
        print("\nExample of happy image from combined dataset")
        file_path = os.path.join(curr_directory, 'data/combined/happy/S138_005_00000016.png')
        img2 = load_img(file_path)
        img2 = img_to_array(img2)
        im2 = Image.open(file_path)
        display.display(im2)

        ## Example 3
        print('\nExample of an image from RAF-DB dataset')
        file_path = os.path.join(curr_directory, 'data/others/aligned/train_00569_aligned.jpg')
        #images = np.load(file_path)
        img3 = load_img(file_path)
        img3 = img_to_array(img3)
        im3 = Image.open(file_path)
        display.display(im3)
```

Example of angry image from ck+ dataset



Example of happy image from combined dataset



Example of an image from RAF-DB dataset



Pre-processing

What features did you use or choose not to use? Why?

We only use the image itself as our feature because we believe that the face image contains everything we need for the model and neural network is capable of handling single image as the input.

If you have categorical labels, were your datasets class-balanced?

Our datasets is not class-balanced that we have more "happy" class than all the other classes.

How did you deal with missing data? What about outliers?

We don't have missing data, but we do have outlier. We filtered the data with the haar cascade classifier to get rid of the outliers.

What approach(es) did you use to pre-process your data? Why?

We converted our image data into greyscaled and normalize data. Then, we only choose frontal face image and combine filtered RAF data with CK+ data to build a new dataset for our model.

Are your features continuous or categorical? How do you treat these features differently?

Our features include just the image itself, so it's hard to say if they are continuous or discrete. In our opinion, we believe that the image data is continuous as it the pixel values are floating points. Each image is represented as an array of pixels. The question as to how we treat these features differently does not really apply as all our features are continuous.

In []:

```
# For those same examples above, what do they look like after being pre-processed?
from sklearn.preprocessing import StandardScaler

# image from CK+ looks the same after preprocessing

# image from RAF-DB

print('RAF-DB images')

## originally
print('Before preprocessing ')
display.display(im3)
```

```
## preprocessing grayscale
print('\nAfter preprocessing')
im4 = ImageOps.grayscale(im3)
im4 = im4.resize([48,48])
im4 = StandardScaler(im4)
display.display(im4)
```

RAF-DB images

Before preprocessing



After preprocessing



In []:

```
# Visualize the distribution of your data before and after pre-processing.
# You may borrow from how we visualized data in the Lab homeworks.
# It's hard to show the distribution of image data so we leave this part blank
```

Models and Evaluation

Experimental Setup

How did you evaluate your methods? Why is that a reasonable evaluation metric for the task?

We evaluated our methods by comparing the loss function and accuracy of all of our approaches together.

What did you use for your loss function to train your models? Did you try multiple loss functions? Why or why not?

For our loss function to train the model, we used categorical cross-entropy function in the TensorFlow.keras. This is a loss function that is commonly used in multi-class classification tasks and the model will decide which one class of the data is belong to. This loss function will automatically do the tie-breaking for us.

How did you split your data into train and test sets? Why?

We used split_test_train in the sklearn package and we split our data to 30% testing and 70% training data. In machine learning literatures, researchers commonly follow the 70% training and the

30% testing split. We thus decided to go with that approach to have enough data to train the CNN but also leave out enough data for evaluation down the road.

```
In [ ]: # Code for loss functions, evaluation metrics or link to Git repo

# model.compile(optimizer = SGD, loss='categorical_crossentropy', metrics=['accuracy'])

# The codes for the loss functions, and evaluation metrics are found in the train.py fi
# Our loss function, as suggested above, is the categorical cross-entropy function in t
```

Baselines

What baselines did you compare against? Why are these reasonable?

We set our baseline to be VGG16/19 whose accuracy on the RAF data is 58%. Since VGG16/19 is the most common neural network used for face recognition problems, which is also suitable for our study.

Did you look at related work to contextualize how others methods or baselines have performed on this dataset/task? If so, how did those methods do?

Yes we looked at [past literatures](#) that used architecture such as VGG16/19, AlexNet, baseDCNN, center loss and DLP-CNN. The average accuracy is around 60% for these models implemented on RAF dataset.

Methods

What methods did you choose? Why did you choose them?

We implemented three models, namely Simple CNN, Gabor CNN, and Gabor VGG. Simple CNN is the basic Convolutional Neural Network, thus we use this as the basic model. The Gabor layer we added in our neural network is delivered from Gabor filter, a linear feature used for analyzing whether there is any specific frequency content in the image in specific directions in a localized region around the point or region of analysis. Using Gabor in our model is because using Gabor compared with using HOG or LBP in mSVM or LDA method from the baseline paper achieved the highest accuracy of RAF dataset. So we chose to use Gabor features in the CNN and VGG model.

How did you train these methods, and how did you evaluate them? Why?

Using `skimage.filters.gabor_kernel` can help us build the Gabor layer. This layer was coded in this [py file](#). We trained these methods on different dataset and observed the loss curves and accuracy curves.

Which methods were easy/difficult to implement and train? Why?

Our Gabor-CNN model are easier to impletment and train because it is not as deep as our Gabor-VGG model. Our Gabor-VGG is more difficult to train because it is a 17-layer model with large network on each layer and we have to use GPU power to enable training is within reasonable time.

For each method, what hyperparameters did you evaluate? How sensitive was your model's performance to different hyperparameter settings?

We evaluate the learning rate, batch size weight-decays and random state. Our model is more sensitve to learning rate than all the other hyperparameters. Each parameter setting for different model and dataset can be seen from the [train.py file](#)

Please refer to `src/model.py` in our github repo [Github model py file](#)

Results

Show tables comparing your methods to the baselines.

Dataset	Simple CNN	Gabor CNN	Gabor VGG	VGG16 (baseline)
CK+	58.85%	96.87%	95.83%	88.92%
Combined	50%	57.08%	57.45%	57.88%
RAF	37.98%	40.40%	51.46%	55.98%

What about these results surprised you? Why?

We are surprised by our Gabor-VGG model tends to overfit in the Combined dataset. We tried different optimizer such as SDG and Adam with different learning rate, weight decay factor, and batch size. Generally, overfitting happend when the network is complicated and the training process prefer to learn the training set well. Reducing batch size and adding data augmentation did slow down the model to over-fit. We finally use the early stopping and reducing the epochs to prevent the overfitting in the training process.

And we are supprised that Gabor CNN has a good performance on CK+ dataset as the Gabor VGG network. Adding Gabor layer to CNN helps the CNN extracted important Gabor features from the input layers, and hence helps the learning of our data.

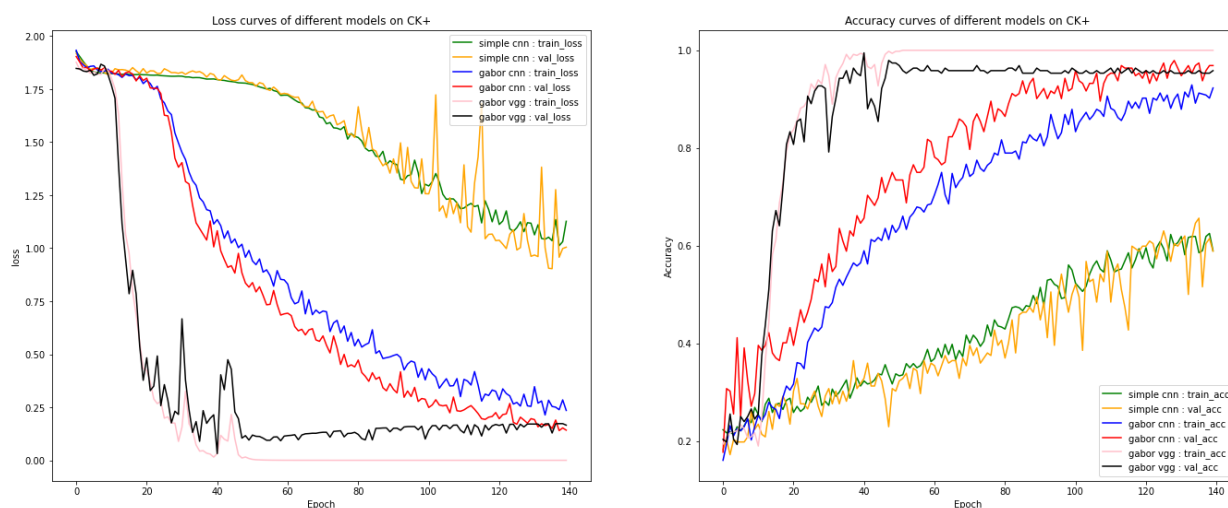
Did your models over- or under-fit? How can you tell? What did you do to address these issues?

Our Gabor-VGG model tend to over-fit as the figure below shown: while the training accuracy continuously increases, the testing accuracy first goes up and then goes down, indicating overfitting. We address this issue using an earlystopper mechanism to stop overfitting.

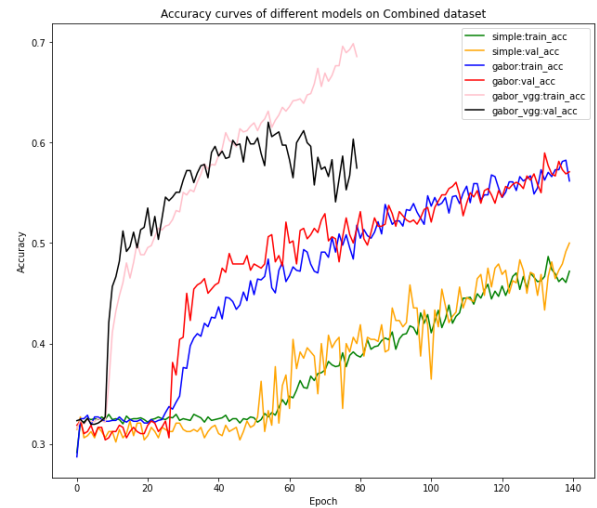
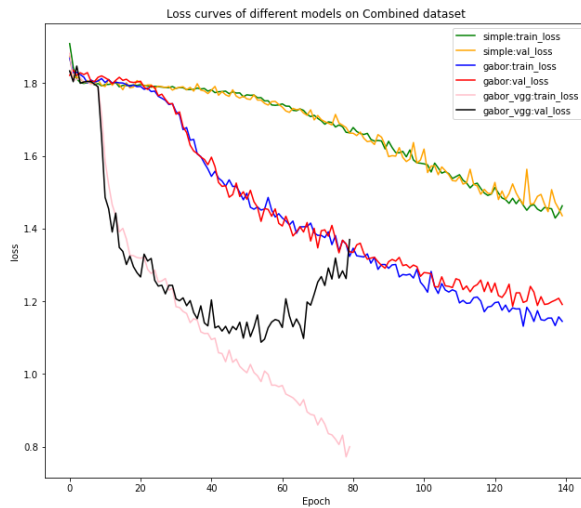
What does the evaluation of your trained models tell you about your data? How do you expect these models might behave differently on different data?

We do expect our model to behave in the following performance: CK+ > Combined > RAF, and our evaluation is consistent with our expectation. Our evaluation shows that the RAF data, our original selection, contains too many outliers or noise since the accuracy of RAF in all of our model or in all of the previous published work related to it are all below 60%.

The first result plots show the loss curves and accuracy curves for different models (Simple CNN, Gabor CNN and Gabor VGG) on the CK+ dataset. This plot does not show slightly over-fitting for the gabor vgg model on the ck+ dataset, but the overall performance is good. This picture shows that Gabor VGG model has the highest loss decline rate and accuracy incline rate. Both Gabor CNN and Gabor VGG model have better performances than the Simple CNN on the CK+ dataset. These curves are generated from [this ipynb](#)



This result plots show the loss curves and accuracy curves for different models (Simple CNN, Gabor CNN and Gabor VGG) on the Combined dataset. It shows that the Gabor VGG model performed over-fitting on this dataset. We used Early stopping mechanism in the training process and only have the first 80 epochs on file. This plot has the similar trend as the previous one, that Gabor CNN and Gabor VGG have a better accuracy than simple CNN, and Gabor VGG achieved certain accuracy or loss quickly. These curves are generated from [the same notebook](#)



We also have implemented taking video input or real-time camera to do the facial expression analysis. The .py file (app) can find in this [github page](#), and some result shown in [the read me file](#)

Discussion

What you've learned

Note: you don't have to answer all of these, and you can answer other questions if you'd like. We just want you to demonstrate what you've learned from the project.

What concepts from lecture/breakout were most relevant to your project? How so?

The concept that was most relevant to the project is deep learning and convolutional network. Through the lectures on the Deep Learning and Convolutional Neural networks, we are able to learn about how CNNs extract features from images (which can be applied to our problem). The lectures also cover training CNN models and other Deep Learning networks through hyperparameter tuning such as learning rate, weight decay and etc., which we applied towards this project. In fact, this project builds on top of the facial recognition of CNN discussed in lecture.

What aspects of your project did you find most surprising?

We were more surprised about how difficult it is to tune and train a model for real-life application. The majority time of our project has been spent on building different models, tuning parameters to increase the accuracy.

What lessons did you take from this project that you want to remember for the next ML project you work on? Do you think those lessons would transfer to other datasets and/or models? Why or why not?

We learned that using the right tools would greatly help with building and training the model. Particularly, we switched from using PyTorch which learned from class to Tensorflow in the middle of the project, and we can finally use the GPU to power our model training. This gives us better efficiency and result for the model and we will probably use this for the next similar project in neural network.

What was the most helpful feedback you received during your presentation? Why?

We received a lot of useful feedback from other groups. Some of which are particularly helpful for our project. First of all, three out of four groups raised the question on how we should increase the accuracy of our model (which at the time had around 36%) to our baseline target of 60%. Secondly, two out of four groups asks us how long does it take to analyze each frame/picture with our model since for real-time detection, since it is very important to process fast. One group also raised the concern of ethical implications for the application of real-time sentiment identification for factors such as racial or gender bias. This is something we have not implemented but need to think about in the future for similar project. We also got inspiration from other group's project's performance improvement strategy so we switched from pytorch to tensorflow.

If you had two more weeks to work on this project, what would you do next? Why?

If we had two more weeks to work on this project, we would complete the model for compound emotion. In addition, we would like to use some cloud services such as Azuer ML or AWS with more GPU power as currently we train our model on our local machines. This would also enable us to try more complex architecture. Lastly, we would also like to introduce other mechanisms such as attention or RCNN which helps us to achieve faster recognition for frames in the video.