

# Is This Restaurant Good for Kids ?

Christina Bogdan

Center for Data Science, NYU  
New York City  
ceb545@nyu.edu

Vincent Chabot

Center for Data Science, NYU  
New York City  
vec241@nyu.edu

Urjit Patel

Center for Data Science, NYU  
New York City  
up276@nyu.edu

## I. INTRODUCTION

Our goal throughout this project was to predict high level attributes related to venues based on pictures uploaded by customers on the Yelp platform. Typically, we were interested in predicting whether a restaurant would be good for dinner, good for groups and/or good for kids. Each business could potentially have several true labels. There were mainly two points of interest we had to deal with to solve this question :

- 1) Inferring very general attributes from images/caption analysis (rather than the most usual case of simply classifying the images based on the objects represented)
- 2) The labels to predict were only business-level labels. However, each image might represent a different business-level level (for example, a picture of a cocktail might indicate that a place was good for groups, but an image of a meal might indicate that that business was also 'Good for Dinner'. This lack of image-level labels made our task semi-supervised.

## II. BUSINESS UNDERSTANDING



COULD you tell us, by looking at this picture, if this restaurant looks rather good for kids? Would you rather recommend it as takeaway or for dinner? This is typically the kind of high level, context information customers would like to know before choosing the next restaurant to go to. The kind of questions you could ask a friend that you know already went to this place for example. But what if your favorite recommendation platform / research engine / social network could answer all your questions? Of course, it is not easy to find data that carry such high level information. But both the large amount of pictures uploaded by customers and possibilities offered by machine learning and deep learning techniques makes it interesting to try to answer such questions.

## III. DATA UNDERSTANDING

Our initial data consisted of two JSON files and a dataset of images. We highlight information about the datasets below:

### YELP BUSINESS FILE

- JSON file of 77445 rows corresponding to 77445 businesses
- Contains fifteen columns. Full details can be found at the bottom of [this](#) page under the *business* section. We used three columns in our project:
  - *business id*: Unique identifier for each business
  - *categories*: Describes the type of business. We did not directly use this field as a feature, but used it to filter down our data.
  - *attributes*: Our label. The attributes field contains many general subcategories indicating whether the business accepts credit cards, the ambiance of the business, whether it is good for breakfast/lunch/dinner, etc. From these, we wanted to predict all attributes indicating what the business was "Good for". There were five binary attributes relating this - **Good for Breakfast**, **Good for Lunch**, **Good for Dinner**, **Good for Kids**, **Good for Groups**, and **Good for Dancing**. Originally, we attempted to predict all five of these attributes. Upon seeing the skew of the data, however, we decided to only predict three: **Good for Dinner**, **Good for Groups**, **Good for Kids**

### YELP PHOTO FILE

- JSON file of 200143 rows and 4 columns
- Maps business IDs to photo IDs
- Contains two photo-level features:
  - *label*: Either food, drink, menu, outside, inside. Relates to the content of the image.
  - *caption*: The caption, if any, attached to the image. We used this as a feature in one of our classifiers.

### IMAGE DATASET

- 200,003 images in JPG format

#### A. Data Cleaning & Engineering

We parsed our two JSON files in Python using Pandas. This caused some issues for us - when parsing the business JSON file, the *attributes* field was represented by a sub-dictionary in our original dataframe:

	attributes	business_id	categories	city
0	{u'Take-out': True, u'Drive-thru': False, u'...}	5UmKMjUEUNdYWqAnhGckJw	[Fast Food, Restaurants]	Dravosburg
1	{u'Happy Hour': True, u'Accepts Credit Cards': ...}	UsFtqoBl7naz8AVUBZMjQQ	[Nightlife]	Dravosburg

To get our data into the right format, we needed to discard all of the unneeded attributes and create columns out of the ones that were relevant to us:

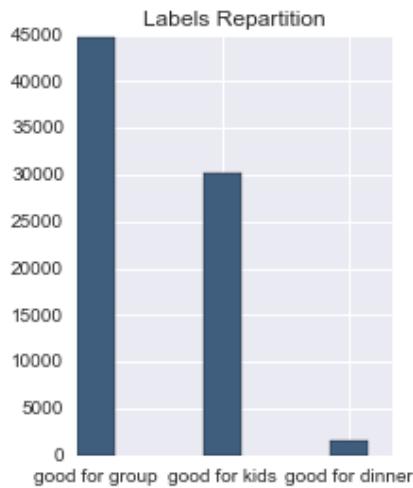
	business_id	Good_for_kids	Good_for_groups	dinner
100	D8y00MIChtJmhP4zMDJdw	1	1	0
101	0r61RRNbGqCFvqNV_vp9bw	0	1	0
102	wddLqUkQAbsEaoRwD4Mgg	1	1	1
103	JUIR59o5IJ1LaeR3NtiOw	1	1	0
104	PZjml2VgR0aEWBx3C7d8Jg	1	1	0

Once this was done, we joined our parsed business dataset to the auxiliary photos dataset, to associate each set of labels with a photo. After this was done, we filtered our photos to only include images that came from businesses with the category **Restaurants, Food, Bars, Nightlife, Fast Food, and Coffee & Tea**. Our reasoning was that the labels that we were trying to predict were only relevant to businesses within these categories.

Originally, our plan was to predict all five 'Good For' fields. However, most of our data fell into the three attributes 'Good for Groups', 'Good for Kids', and 'Good for Dinner'.

	photo_id	business_id	caption	label
0	...L4adlY4jEWo3cnWPcw	WltPLuH8Np_btPnoGaF5XA	NaN	none
1	--0Uxeaz2kZn3aUjOWw	1-One6SzxeQfKfZeVKJWW	Traditional cribs, convertible cribs, twin or ...	none
2	--5rFJBzhXlpGKgMb51w	PrYKua8LhcYYjMmniPPkhA	hmmm... Thai iced tea w/ simple syrup on the... drink	drink
3	--852IDNHxJls5FTlzbMCQ	9iT4dxLPgAXSGiSQVdMLpA	NaN	none
4	--8N4jEwQwGFYpc3F589A	g9USj3AgG0SqCS6hgXiA	Robyn feeding a mountain lion!	none

Even among these three labels, the distribution of the data was highly skewed:



To prepare our image data, we rescaled each image to be of size **96 x 96** and converted it from RGB to YUV format before storing it in matrix format. After this was done, we performed entry-wise normalization of the pixels on the train, validation, and test data such that for each pixel  $x_{ij}$ :

$$x'_{ij} = \frac{x_{ij} - \mu_{tr}}{\sigma_{tr}}$$

Where  $\mu_{tr}$  and  $\sigma_{tr}$  are the mean and standard deviations of the mean and standard deviation across the training set.

## B. Data Insights

Examining the correlation between the three target labels we chose (Figure 1) reveals some interesting information about our data. It may be noted in particular that the 'Good for Kids' and 'Good for Groups' labels are the most correlated, and their correlation is negative. Indeed, we'll see later in our analysis of the image captions that the most relevant words for 'Good for Groups' labeled places correspond to attributes such as "night", "alcohol", "pool", ... which of course do not correspond to places we would recommend for kids.

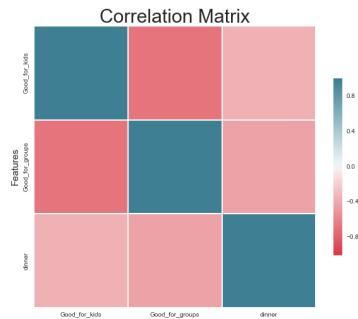


Fig. 1: Correlation matrix between the labels

## IV. GLOBAL FRAMEWORK

The global framework we designed is illustrated on Figure 2. In order to solve the fact that most of the pictures we had had several labels (the same labels as the restaurant they were associated to), we first isolated the images from businesses that had only one true label. We trained a model on the captions of these images in order to predict individual labels on the remaining images associated with multi-label businesses. This would turn our initial problem into a fully supervised problem with every picture only having now one label, that we believed would be the most relevant for that picture. We would then solve that supervised problem by training independently both an image and a caption-based model. We then averaged the predicted probability from each model and attributed the label that had the maximum probability. From the single image level predictions, we could eventually predict the business level multiple attributes as follows: as soon as a business would have a picture with a given predicted label, we would predict that label for the business as well.

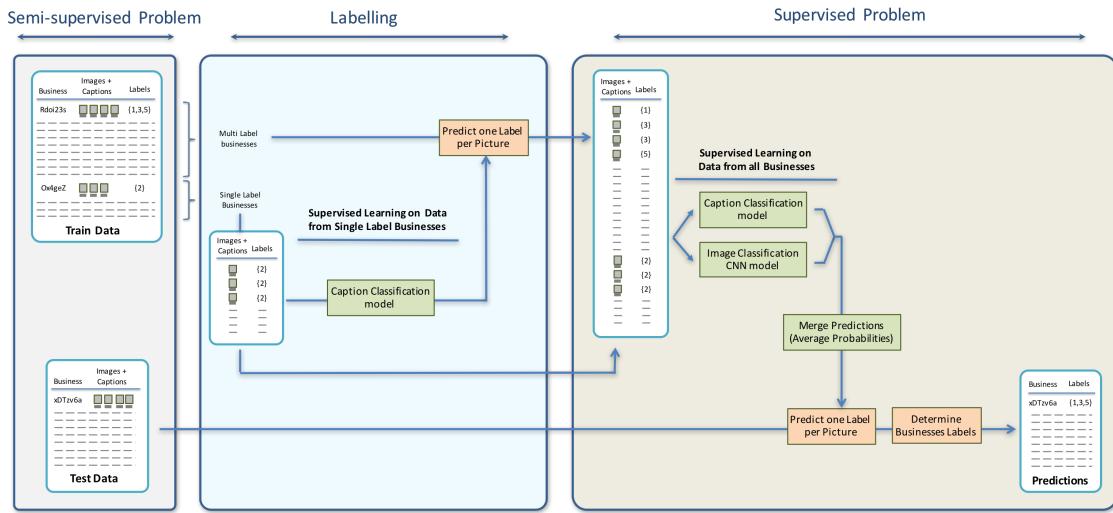
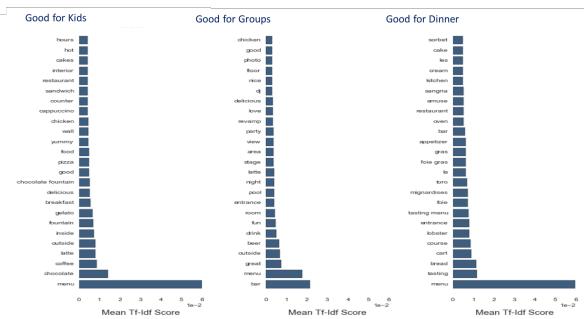


Fig. 2: Global Framework

## V. RETRIEVING INDIVIDUAL IMAGE LABELS

By only using images that were associated to single label businesses, we could perform fully supervised classification. We used the captions on these images to train a model using tf-idf. We then used this model to determine a unique picture-level label for the pictures associated to businesses with multiple labels. The figure below displays top contributing words for each class (stop words removed):



We also tried various classifiers and among them, Linear SVC, SGD and Multinomial models gave the best accuracy:

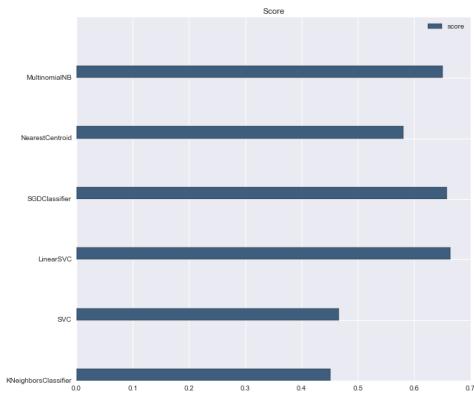


Fig. 3: Performance with various models

Our final results are shown below:

Classifier	MultinomialNB
Train Accuracy	0.8943
Test Accuracy	0.6662

## VI. IMAGE CLASSIFICATION VIA CAPTION & IMAGE ANALYSIS

Once we recovered individual labels for each image, we trained two fully supervised models separately before aggregating predictions from both to produce a final classification:

- 1) A text classifier on the images captions, using the exact same methodology as before
- 2) A Convolutional Neural Network

We outline this stage of our model below:

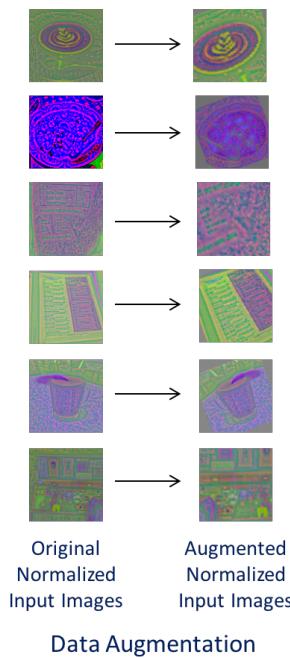
### A. CNN Methodology

We implemented our convolutional neural network using the Torch package in the Lua programming language. We ended up using a ten layer model (7 convolutional layers + 3 linear layers), though we originally attempted a larger model and found that it did not work well with our data. A visualization of our network architecture can be found in the **appendix**.

The main issue that we faced when building our CNN was dealing with the skew of the data. After running our 70,000 images through the text classifier to get image-level labels, we had approximately 30,000 images in the 'Good for Groups' class, 22,000 in 'Good for Kids', and 700 in 'Good for Dinner'. We attempted to overcome this by performing selective data augmentation on the two classes with reduced size. That is, we randomly applied visual transformations to the smaller classes in our training data and added these transform images to our data set in order to increase the number of images for those classes. Our augmentation techniques are summarized below:

- **Rotation:** Rotate the input by some random angle
- **Scaling:** Randomly rescale image up or down before re-cropping to fit model input size
- **Translation:** Randomly shift image left, right, up, or down by some number of pixels
- **Contrast:** Modify the contrast of the input image
- **Horizontal Flip:** Horizontally flip input image

For each image, we set a 50-80% probability of applying a particular transformation. We experimented with many different augmentation sizes, and ultimately decided to add **1000** images to the 'Good for Kids' class and **600** images to 'Good for Dinner'. We did not want to add too many augmented images to the data set for fear of overfitting, and the large volume of images in the 'Good for Groups' and 'Good for Dinner' classes made the augmentation only have moderate success - if we had less total images than maybe it would have worked better. Some samples of our data augmentation are shown below:



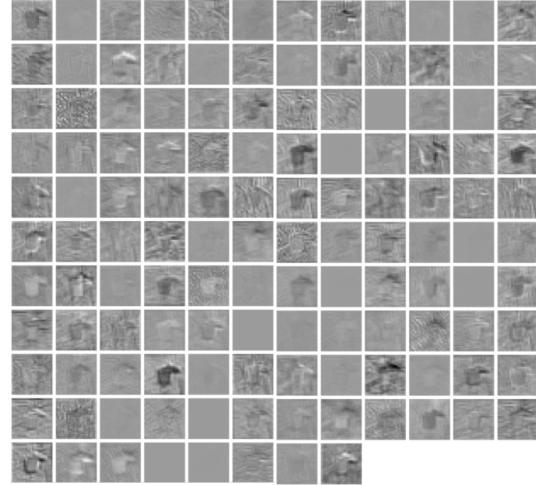
Our aim with the CNN was to extract spatially invariant, high-level features from the images that represent the very general classification of being good for groups, kids, or dinner. Each convolutional layer has three stages that work to accomplish this: convolution, application of a non-linearity, and pooling. The work done in each of these three stages is summarized below:

- 1) **Convolutional Stage:** In this stage, we compute the product of a learned "filter" (represented by a weights matrix) with local regions of the image. This product allows us to identify what sections of an image are relevant to what we are attempting to predict. At each convolutional layer we define the number of feature maps we want to output - this is how many filters we are applying over each image in this layer.
- 2) **Non-Linearity:** The application of a non-linearity to the

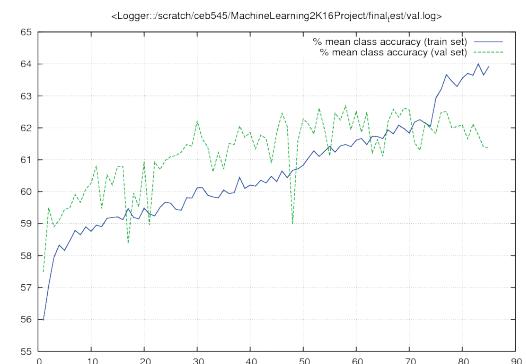
data allows the network to learn beyond linear combinations of the data. Without applying a non-linearity to the output of the convolutional stage, we could represent all layers in the CNN as one linear combination.

- 3) **Pooling:** The pooling stage downsamples the output of the previous two stages - this is important for introducing spatial invariance and reducing computational costs.

The output of one of the convolutional layers for one of our images is shown below:



With our best model, we achieved **67.11%** accuracy on our training data and **60.21%** accuracy on our test data, which resulted after running our model for approximately 90 epochs. The convergence of the model's accuracy over time is shown in the figure below:



## VII. FINAL PREDICTION OF PHOTO LABELS

Both our caption model and CNN assigned class probabilities for each photo. To create a final prediction for a particular photo, we average these two scores, and then took the maximum score of this average as the label for the photo.

## VIII. FINAL PREDICTION OF BUSINESS LABELS

We assigned labels to each business based on the predicted labels of the images associated to that business. We formed the prediction for a business by taking the union of predictions across all images.

## IX. EVALUATION

We use Precision, Recall and F1 score to evaluate the performance of our model. Indeed, Precision is very important for our use case because customer satisfaction is crucial for recommendation platform. The satisfaction of customers won't be particularly affected if you do not recommend a certain place, whereas it would highly affect their trust in your service if you were for instance to recommend a good place for kids when in fact it is not.

The chart below shows our results in terms of true/false positives/negatives obtained with our model. As we can see, we were able to have a low number of false positives. Nevertheless, this was performed at the cost of a quite high number of false negatives

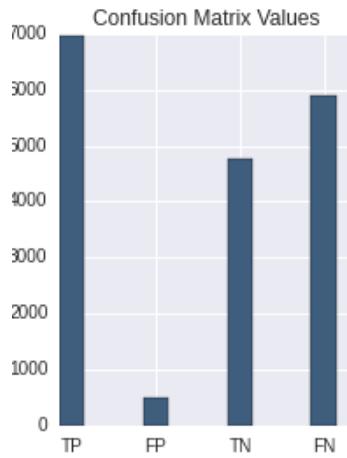


Fig. 4: Summary of Confusion Matrix values

Finally we used the F1 score to assess the overall performance of our model, which is computed as follow:

$$F1score = 2 \frac{Precision \cdot Recall}{Precision + Recall}$$

With:

$$Precision = \frac{TruePositive}{TruePositives + FalsePositives}$$

$$Recall = \frac{TruePositive}{TruePositives + FalseNegatives}$$

We summarize our final results below:

Precision	<b>0.9318</b>
Recall	<b>0.5405</b>
F1 score	<b>0.6841</b>

## X. CONCLUSION & FUTURE WORK

Ultimately, we were able to improve on our random baseline with our joint CNN and caption-based model. If we had more time, we would have used another decision function to infer

business labels from our prediction photo labels, instead of just taking the union of predicted labels across all photos. We would have also explored other methods to combine our predictions from our image and caption based models. Finally, we would have included our implementation of  $\alpha$ SVM [1], which we believe would have helped us in our attempts to generate image-level labels.

## XI. ACKNOWLEDGMENTS

We would like to thank our advisor Kush R. Varshney for the strong support and always very helpful advice he provided us throughout our project. He always provided us very thoughtful responses to our questions and was had great suggestions for us throughout our project.

## REFERENCES

- [1]  $\alpha$ SVM for Learning with Label Proportions  
[http://felixyu.org/pdf/felix\\_yu\\_pSVM\\_icml2013.pdf](http://felixyu.org/pdf/felix_yu_pSVM_icml2013.pdf)

## XII. APPENDIX

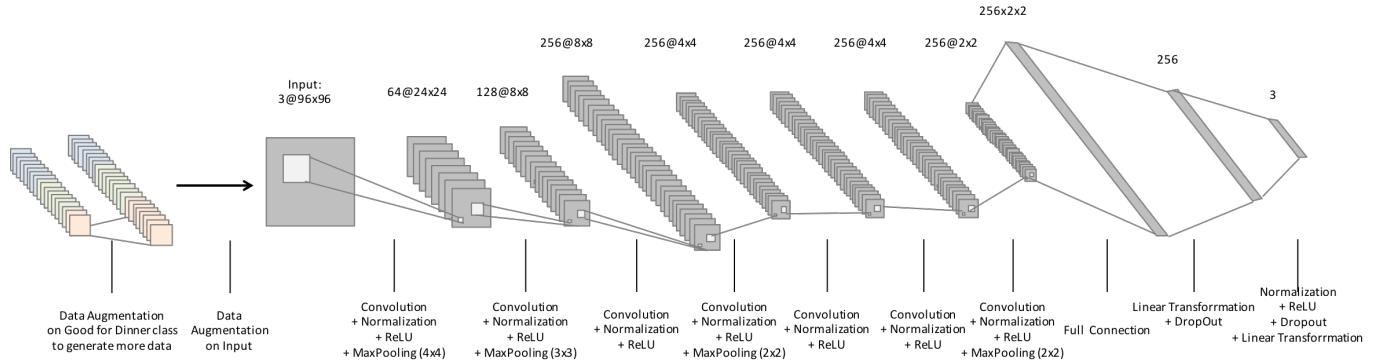


Fig. 5: Architecture of the Convolutional Neural Network we used



**Business Case**

**CHOCOLATE FOUNTAIN!**

**Christina Bogdan**   **Vincent Chabot**   **Urjit Patel**

**Overview**

Our goal throughout this project was to predict high level attributes related to venues based on pictures uploaded by customers on the Yelp platform. Typically, we were interested in predicting whether a restaurant would be 'good for dinner', 'good for groups' and/or 'good for kids'. Each business could potentially have several true labels. There were mainly two major points of interest we had to deal with to solve this question :

- 1) Inferring very general attributes from images/caption analysis (rather than the most usual case of simply classifying the images based on the objects represented)
- 2) The labels to predict were only business level labels. From the fact that a business could have several true labels, we would not have direct access to individual image level labels, making the task semi-supervised.

To solve this, we first isolated the images from businesses that had only one true label. We trained both an image based model and a caption based model in order to predict individual labels on the remaining images associated with multi-label businesses. This would turn our initial problem into a supervised model, that we would solve by training another two image and caption based models. From the single image level predictions, we could eventually predict the business level multiple attributes.

**Data Understanding**

Our initial data consisted of two files. The first mapped each business to its various attributes we had to predict (possibly multiple true labels per business) :

Business ID	Name	Address	Attributes
1	Good for dinner	123 Main St	Good for dinner
2	Good for groups	456 Main St	Good for groups
3	Good for kids	789 Main St	Good for kids
4	Good for dinner, good for groups	123 Main St	Good for dinner, Good for groups
5	Good for dinner, good for kids	456 Main St	Good for dinner, Good for kids
6	Good for groups, good for kids	789 Main St	Good for groups, Good for kids

In the second data set, each row consisted of a picture id with its caption if any and the business it was related to :

Image ID	Caption	Business ID
1	Chocolate fountain at a restaurant	1
2	People eating at a restaurant	1
3	Food on a plate	1
4	Chocolate fountain at a restaurant	2
5	People eating at a restaurant	2
6	Food on a plate	2
7	Chocolate fountain at a restaurant	3
8	People eating at a restaurant	3
9	Food on a plate	3
10	Chocolate fountain at a restaurant	4
11	People eating at a restaurant	4
12	Food on a plate	4
13	Chocolate fountain at a restaurant	5
14	People eating at a restaurant	5
15	Food on a plate	5
16	Chocolate fountain at a restaurant	6
17	People eating at a restaurant	6
18	Food on a plate	6

We can observe below that our classes were not evenly distributed. More particularly, we had very few "good for dinner" businesses. Furthermore, displaying the correlation matrix between the attributes, we see that businesses that are good for kids and those that are good for groups are relatively negatively correlated, as we could expect. On the other hand, good for dinner part, good for group business are more related to key words as "tear", "bar", "night", ...:

Labels Distribution : 

Correlation Matrix : 

**Image individual label retrieving via captions of images with unique label**

Good for Kids   Good for Groups   Good for Dinner

From the analysis of the words contributing the most to each class (chart above), we add additional stop words such as: 'page', 'restaurant', 'area', 'place', 'interior', 'outside', 'like', 'saucer', 'view'

Across various classifiers, Linear SVC, SGD and Multinomial models give the best accuracy.

Final Results

Classifier	Train Accuracy	Test Accuracy
Linear SVC	0.8943	0.6662
Multinomial		

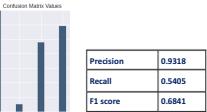
**Conclusion & Future Work**

Ultimately, we were able to improve on our random baseline with our joint CNN and caption-based model. At the same time, we also tried out another decision function using business labels from our prediction photo labels, instead of just taking the most frequent label from all the photos. We would have also explored other methods to combine our predictions from our image and caption baselines. Finally, we also tried to extend our implementation of SVM, which we believe would have helped us in our attempts to generate image-level labels.

**Business multi labels prediction & Final results**

Finally, we predicted the business labels based on the predicted labels of the images associated to each business: if a business had at least one image with a given label, we would predict that label for the business as well.

We expose below our final results :

Confusion Matrix Values : 

Prediction	Recall	F1 score
Good for dinner	0.9318	
Good for groups	0.5405	
Good for kids	0.6841	

**Images individual labels predictions**

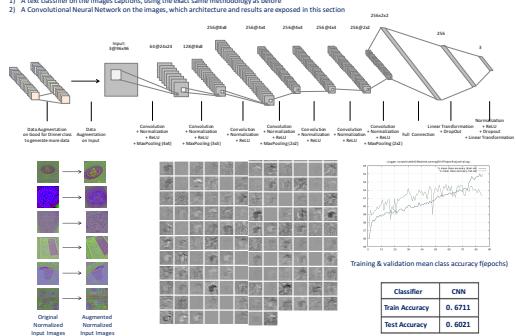
From the merge predictions of both caption based and image based models (by simply averaging probability predictions), we predicted individual labels for each image of the test set by taking the label corresponding to the maximum predicted merged probability.

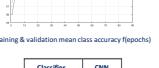
**Images classification via both caption and image analysis**

Once concerned what we feel confident in being image unique labels, we trained 2 fully supervised models separately :

- 1) A text classifier on the image captions, using the exact same methodology as before
- 2) A Convolutional Neural Network on the images, which architecture and results are exposed in this section

Architecture of the CNN used :



Training & validation mean class accuracy (Epochs) : 

Classifier	Train Accuracy	Test Accuracy
CNN	0.6711	0.6002

Final Results (after 92 epochs)

Visualization of the 128 feature maps of the second layer : 

Fig. 6: Final Poster for the NYU CDS Projects Presentation Event