



**INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY
ADVANCED SCIENTIFIC RESEARCH AND INNOVATION
(IJMASRI)**

ISSN: 2582-9130

IBI IMPACT FACTOR 1.5

DOI: 10.53633/IJMASRI

RESEARCH ARTICLE

REAL-TIME AMERICAN SIGN LANGUAGE RECOGNITION WITH NEURAL NETWORKS

¹Mohd Arifullah ¹Fais Khan and ¹Yash Handa

¹Maharaja Agrasen Institute of Technology, GGSIPU, Delhi

Abstract

Actual-time signal language translator is a crucial milestone in facilitating communication among the deaf community and the general public. Introducing the development and use of yanked sign Language Spelling Translator (ASL) based on the convolutional neural network. We use the pre-skilled Google Net architecture educated inside the ILSVRC2012 database, in addition to the ASL database for Surrey University and Massey university ASL to apply gaining knowledge of switch in this task. We have developed a sturdy version that constantly separates the letters a-e from the original users and any other that separates the spaced characters in maximum cases. Given the limitations of the information sets and the encouraging consequences acquired, we are assured that with similarly studies and further facts, we can produce a totally customized translator for all ASL characters.

Keywords: Sign Language, Image Recognition, American Sign Language, Expressions signals, CNN

Introduction

American Sign Language (ASL) is used for communication in the deaf community. However, there are only ~ 2.5-5 lakh speakers which significantly lowers the number of people that they can anyway communicate with. The other option of

written means to communicate is impersonal and even non-useful when an emergency occurs. In order to reduce this obstacle and to enable dynamic communication, we present an ASL finding system that uses Convolutional Neural Networks (CNN) in time to translate a video of a user's ASL signs into

text our problem consists of three tasks to be done in real time:

1. Obtaining video of the user signing (input)
2. Classifying each frame in the video to a letter
3. Reconstructing and displaying the foremost likely word from classification scores (output)

represents a big challenge because of variety of considerations, including:

- Natural concerns (e.g., lighting sensitivity, background, and camera position)
- occlusion (e.g., hand being out of the field of view)
- Sign layer detection (when a sign ends and the next begins)
- Co-articulation (when a sign is disturbed by the preceding or succeeding sign)

While Neural Networks are applied to ASL letter recognition (Appendix A) within the past with accuracies that are consistently over 91% , many of them require a 3-D capture element with motion-tracking gloves or a Microsoft Kinect, and just one of them provides real-time classifications. The constraints imposed by the additional requirements reduce the scalability and feasibility of those solutions.

Our system features a pipeline that takes video of a user signing a word as input through an internet application. We then extract individual frames of the video and generate letter probabilities for every employing a CNN (letters a through y, excluding j and z since they require movement). With the employment of a range of heuristics, we group the frames supported the character index that every frame is suspected to correspond to. Finally, we use a language model so as to output a possible word to the user.

Similar Work

ASL recognition isn't a replacement computer vision problem. Over the past twenty years,

researchers have used classifiers from a spread of categories that we are able to group roughly into linear classifiers, neural networks and Bayesian networks.

While linear classifiers are easy to figure with because they're relatively simple models, they require sophisticated feature extraction and preprocessing methods to achieve success (Mitchell et al.,2006; Singha and Das, 2003; Sharma. 2013). Singha and Das obtained accuracy of 96% on 10 classes for images of gestures of 1 hand using Karhunen-Loeve Transforms. These translate and rotate the axes to ascertain a brand new organization supported the variance of the info. This transformation is applied after employing a skin filter, hand cropping and edge detection on the pictures. They use a linear classifier to differentiate between hand gestures including thumbs up, finger pointing left and right, and numbers (no ASL). Sharma et al. use piece-wise classifiers (Support Vector Machines and k-Nearest Neighbors) to characterize each color channel after background subtraction and noise removal . Their innovation comes from employing a contour trace, which is an efficient representation of hand contours. They attain an accuracy of 62.3% using an SVM on the segmented color channel model.

Bayesian networks like Hidden Markov Models have also achieved high accuracies (Starner and Pentland. 1997; Jeballi, 2013; Suk. 2010). These are particularly good at capturing temporal patterns, but they require clearly defined models that are defined before learning. Starner and Pentland used a Hidden Markov Model (HMM) and a 3-D glove that tracks hand movement. Since the glove is ready to get 3-D information from the hand no matter spatial orientation, they were ready to achieve a formidable accuracy of 99.2% on the test set. Their HMM uses timeseries data to trace hand movements and classify supported where the hand has been in recent frames. (Suk 2010) propose a way for recognizing hand gestures in an exceedingly continuous video stream

employing a dynamic Bayesian network or DBN model . They try to classify moving hand gestures, like making a circle round the body or waving. They achieve an accuracy of over 99%, but it's worth noting that everyone gestures are markedly different from one another which they're not American linguistic communication. However, the motion-tracking feature would be relevant for classifying the dynamic letters of ASL: j and z.

Some neural networks are wont to tackle ASL translation (Mekala, 2011; Admasu and Raimond, 2010; Atwood, 2012). Arguably, the foremost significant advantage of neural networks is that they learn the foremost important classification features. However, they require considerably longer and data to coach. To date, most are relatively shallow. Classified video of ASL letters into text using advanced feature extraction and a 3-layer Neural Network (Mekala, 2011). They extracted features in two categories: hand position and movement. Before ASL classification, they identify the presence and placement of 6 “points of interest” within the hand: each of the fingertips and also the center of the palm. (Mekala, 2011) also take Fourier Transforms of the pictures and identify what section of the frame the hand is found in. While they claim to be able to correctly classify 100% of images with this framework, there's no mention of whether this result was achieved within the training, validation or test set.

Admasu classified Ethiopian Sign Language correctly in 98.5% of cases employing a feedforward Neural Network (Admasu and Raimond, 2010). They use a big amount of image preprocessing, including image size normalization, image background subtraction, contrast adjustment, and image segmentation. Admasu and Raimond extracted features with a Gabor Filter and Principal Component Analysis.

The most relevant work to this point is L. Pigou et Al's application of CNN's to classify 20 Italian gestures from the Cha Learn 2014 gazing People gesture spotting competition. They use a

Microsoft Kinect on full body images of individuals performing the gestures and achieve a cross-validation accuracy of 91.7%. As within the case with the aforementioned 3-D glove, the Kinect allows capture of depth features, which aids significantly in classifying ASL signs.

Approaches

Classification

Our ASL letter classification is completed employing a convolutional neural network (CNN or Convent). CNNs are machine learning algorithms that have seen incredible success in handling a spread of tasks associated with processing videos and pictures. Since 2012, the sphere has experienced an explosion of growth and applications in image classification, object localization, and object detection.

A primary advantage of utilizing such techniques stems from CNNs abilities to be told features furthermore because the weights akin to each feature. Like other machine learning algorithms, CNNs seek to optimize some objective function, specifically the loss function. We utilized a SoftMax-based loss function:

Equation (2) is that the SoftMax function. It takes a feature vector z for a given training example, and squashes its values to a vector of $[0,1]$ -valued real numbers summing to 1. Equation (1) takes the mean loss for every training example, x_i , to provide the total SoftMax loss.

$$Loss = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_{i,y_i}}}{\sum_{j=1}^C e^{f_{i,j}}} \right) \quad (1)$$

$$f_j(z) = \frac{e^{z_j}}{\sum_{k=1}^C e^{z_k}} \quad (2)$$

N = total number of training examples
 C = total number of classes

Using a SoftMax-based classification head allows us to output values like probabilities for every ASL letter. This differs from another popular choice: the SVM loss.

Using an SVM classification head would end in scores for every ASL letter that may ultimately map to probabilities. These probabilities afforded to us by the SoftMax loss allow us to more intuitively interpret our results and prove useful when running our classifications through a language model.

Learning Transfer

Transfer Learning may be a machine learning technique where models are trained on (usually) larger data sets and refactored to suit more specific or niche data. This can be done by recycling some of the weights from the pre-trained model and reinitializing or otherwise altering weights at shallower layers. The foremost basic example of ready to this be a completely trained network whose final classification layer weights are reinitialized to be able to classify some new set of information. The first benefits of such a method are its less demanding time and data requirements. However, the challenge in transfer learning stems from the differences between the initial data won't to train and therefore the new data being classified. Larger differences within these data sets often require re-initializing or increasing learning rates for deeper layers in the net.

Caffe and Google Net

We employed Caffe, a deep learning framework, so as to develop, test, and run our CNNs. Specifically, we used Berkeley Vision and Learning Center's Google Net pre-trained on the 2012 ILSVRC dataset. This net output three different losses at various depths of the web and combines these losses before computing a gradient at training time.

Basic Technique

Our overarching approach was to fine-tune a pre-trained Google Net. Our data consists exclusively

of hands in 24 different orientations, while the ILSVRC data consists of 1000 uniquely different objects or classes. Although the ILSVRC data has some classes that are quite just like one another (e.g., various breeds of dogs), our data was comprised of the identical object merely positioned and oriented in numerous ways. Considering the stark differences between our data and therefore the data that Google Net was pre-trained on, we decided to check the effectiveness of altering a range of the pre-trained weights at different depths. We amplified learning rates by an element of ten and completely reset weights within the first one, two or three layers of the online using Xavier initialization.

Developing our Structure

In order to get images of the user signing in Realtime, we created an internet application that's able to access a native camera on a laptop through the browser solely using HTML and JavaScript. This was done using an API created by the globe Wide Web Consortium (W3C). Image capture rate was a large problem we struggled with. We looked to balance network request speeds with computation speeds of our neural network (discussed below). The latter proved to be the first bottleneck, forcing us to decide on a picture capture rate of 1 frame per second. Increasing our capture rate created larger delays in giving the user feedback than we were satisfied with.

Our web application sends images to our server one by one. Each time, the server classifies the image and presents probabilities for every letter. It keeps a running cache of classified images. When it feels confident about the sign being made by the user, it records the top-5 possibly letters supported the cache. It then clears the cache, and lets the user know to maneuver on to the following letter.

Once the user has indicated that they're finished signing, the top-5 letters for every position within the spelled word are passed to a custom unigram language model supported the Brown Corpus. The model takes under consideration substitution costs for similar-looking letters that are

often confused by the classifier additionally because the probabilities for the top-5 letters at each position within the word. Using these, single word unigram probabilities and some of other heuristics, it returns the optimal word to the user.

Datasets



Fig. 1. Dataset examples. Top left: *y*. Bottom left: *k*. Top right: *i*. Bottom right: *e*.

Description

The ASL Fingerspelling Dataset from the University of Surrey's Center for Vision, Speech and Signal Processing is split into two types: color images and depth pictures to create our Translator accessible through a straightforward web-app and laptop with a camera, we chose to only use the color images. they're close-ups of hands that span the bulk of the image surface (Fig. 1).

Dataset A comprises the 24 static signs of ASL captured from 5 users in several sessions with similar lighting and backgrounds. the pictures are in color. Their height-to width ratios vary significantly but average approximately 150x150 pixels. The dataset contains over 65k images.

The Massey University Gesture Dataset 2012 contains 2,524 close-up, color images that are cropped such the hands touch all four edges of the frame. The hands have all been tightly cropped with little to no negative space and placed over a consistent black background. Coincidentally, this dataset was also captured from five users. The frames average approximately 500x500 pixels.

Since there was little to no variation between the photographs for the identical class of every signer, we separated the datasets into training and validation by volunteer. Four of the five volunteers from each dataset were accustomed train, and therefore the remaining volunteer from each was wont to validate. We opted to not separate a test set since that might require us to get rid of one in every of four hands from the training set and thus significantly affect generalizability. Instead, we tested the classifier on the net application by signing ourselves and observing the resulting classification probabilities outputted by the models.

Pre-processing

Both datasets contain images with unequal heights and weights. Hence, we resize them to 256x256 and take random crops of 224x224 to match the expected input of the Google Net. We also zero-center the information by subtracting the mean image from ILSRVC 2012. Since the possible values within the image tensors only span 0-255, we don't normalize them. Furthermore, we make horizontal flips of the photographs since signs is performed with either the left of the proper hand, and our datasets have samples of both cases.

Since the difference between any two classes in our datasets is subtle compared to ILSRVC classes, we attempted padding the pictures with black pixels specified they preserved their ratio upon resizing. This padding also allows us to get rid of fewer relevant pixels upon taking random crops.

Experiments

Evaluation

We evaluate two metrics so as to check our results with those of other papers. the foremost popular criterion within the literature is accuracy within the validation set, i.e. the share of correctly classified examples. One other popular metric is top-5 accuracy, which is that the percentage of classifications where the right label appears within the 5 classes with the very best scores.

Additionally, we use a confusion matrix, which may be a specific table layout that enables visualization of the performance of the classification model by class. this permits us to gauge which letters are the foremost misclassified and draw insights for future improvement.

Experiments

For each of the subsequent experiments below, we trained our model on letters a-y (excluding j). After some initial testing, we found that using an initial base learning rate of 1e-6 worked fairly well in fitting the training data - it provided a gradual increase in accuracy and perceived to successfully converge. Once the improvements within the loss stagnated, we manually stopped the method and decreased the educational rate so as to undertake and increase our optimization of our loss function. We cut our learning rate by factors starting from 2 to 100.

Furthermore, we used the training routine that performed best with real users on our web application ('2_init'). We also built models to only classify letters a-k (excluding j) or a-e to judge if we attained better results with fewer classes.

1 layer reinitialization and learning rate multiple increase ('1_init'):

We initialized this model with the pre-trained weights from the Google Net trained on ILSVRC 2012. We then reinitialized all the classification layers with Xavier initialization and increased the training rate multiple of only this layer so as to assist it learn faster than the remainder of the net's pre-trained layers.

2 layer reinitialization and learning rate multiple increase ('2_init')

We initialized this model with the pre-trained weights from the Google Net trained on ILSVRC 2012. We then reinitialized all the classification layers with Xavier initialization and appropriately adjusted their dimensions to match our number of classes. We increased the educational rate multiples for the highest three layers beneath each classification head.

1 layer reinitialization, learning rate multiple increase, and increased batch size ('1_init')

We initialized this model with the pre-trained weights from the Google Net trained on ILSVRC 2012. We then reinitialized all the classification layers with Xavier initialization and increased the educational rate multiple of only this layer so as to assist it learn faster than the remainder of the net's pre-trained layers. Finally, we drastically increased the batch size from four to twenty.

1 layer reinitialization, uniform learning rate ('full_lr')

We initialized this model with the pre-trained weights from the Google Net trained on ILSVRC 2012. We then reset all the classification layers with Xavier initialization but kept learning rate constant through the web

Results

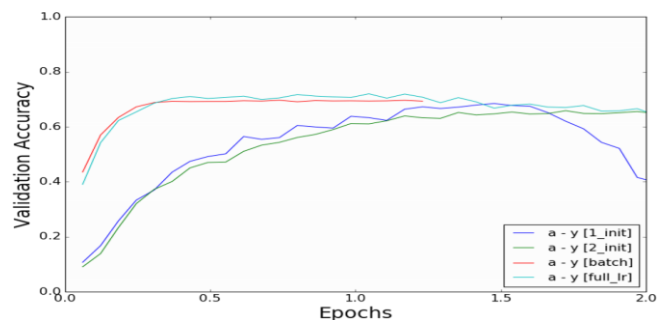


Fig. 2: Epochs vs. validation accuracy for all models trained on letters a-y

| Model | Top-1 Val Accuracy | Top-5 Val Accuracy |
|-----------------|--------------------|--------------------|
| a - y [1_init] | 0.6847 | 0.9163 |
| a - y [2_init] | 0.6585 | 0.9043 |
| a - y [batch] | 0.6965 | 0.9076 |
| a - y [full_lr] | 0.7200 | 0.9098 |
| a - k [2_init] | 0.7430 | 0.897 |
| a - e [2_init] | 0.9782 | 1.000 |

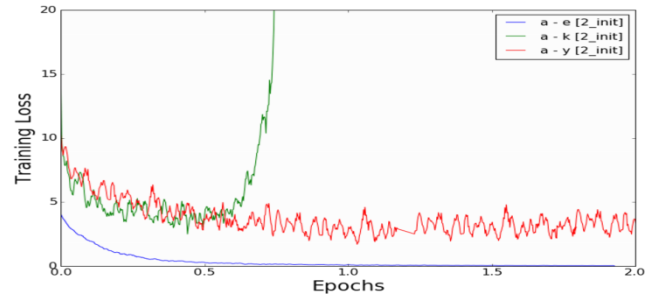


Fig. 5: Epochs vs. training loss for the 2_init models trained on each letter subset

Table 1. Optimal accuracy ranges for all models trained on each letter subset.

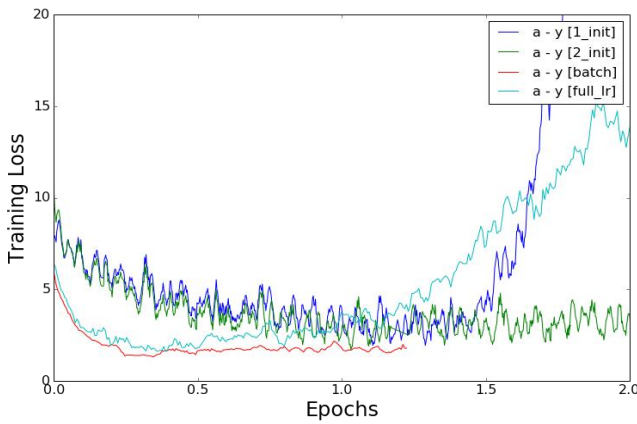


Fig. 3: Epochs vs. training loss for all models trained on letters ay

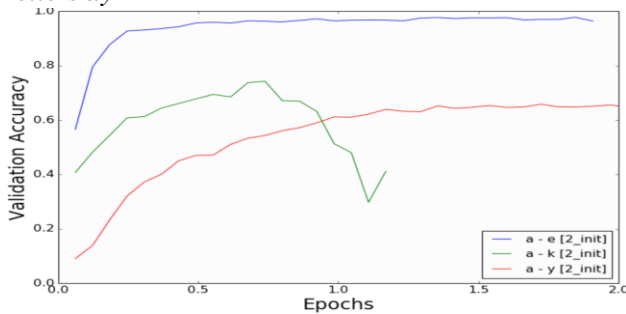


Fig. 4: Epochs vs. validation accuracy for the 2_init models trained on each letter subset

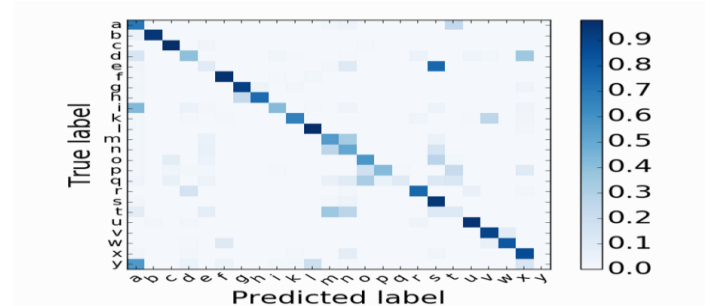


Fig. 6: Confusion matrix for the 2_init model trained on letters ay

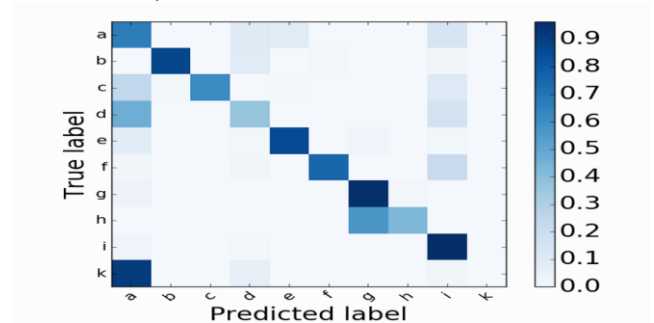


Fig. 7: Confusion matrix for the 2_init model trained on letters ak

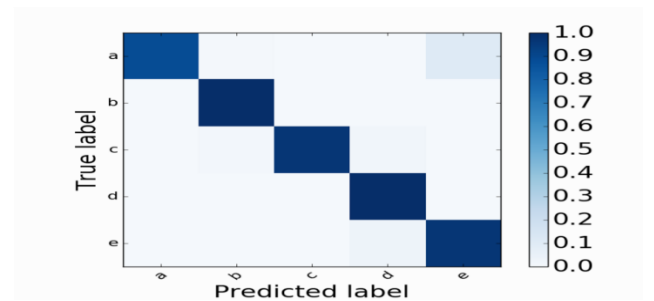


Fig.8: Confusion matrix with test result in training model using matplotlib

Discussion

Accuracy

Our losses (Fig. 3) were very noisy within the '1_init' and '2_init' models. Our space and time constraints initially required us to settle on a less-than-optimal batch size value of 4, leading to the noisy loss. However, after seeing these results, we trained a net employing a Lighting Memory Mapped Database (LMDB) and were ready to increase the batch size to twenty. This allowed us to cut back our loss more smoothly and monotonically, additionally to more quickly converging on a validation accuracy. Similarly, the 'full train' model had uniform learning rates throughout each layer within the net, allowing us to find out more quickly.

This is likely because of the actual fact that we are ready to more easily alter the pre-trained weights within the Google Net and proper for the numerous differences between the datasets. While this results in a tighter fitting of the training data, it didn't seem to end in lower validation accuracy than those attained with other models. After analysing many of the pictures from our dataset, we concluded that they were possibly produced by taking video frames of people making ASL signs within the same room and in an exceedingly single sitting. the dearth of variation in our data set explains the similar validation accuracy in our 'full train' model relative to others.

Interestingly, changing our re-initialization scheme and learning rates had little effect on the ultimate top-1 and top-5 accuracies: the biggest difference between two models was but 7% on top-1 and just over 1% on top-5 accuracy. We did notice, however, that the models utilized that only re-initialized weights at the classification layer strictly outperformed the 2-layer re-initialization model. Notwithstanding the very fact that the excellence between classes in our dataset and also the ILSVRC 2012 dataset are starkly different, this can be not unexpected due to the top quality of the features extracted by Google Net. it had been pre-trained on

significantly more images than we had available in our dataset.

There was little difference between the '1_init' and '2_init' models. Given the very fact that GoogLeNet is 22 layers deep, intuition would lead us to believe that reinitializing 2 layers (versus reinitializing 1 and increasing the training rate multiple on another) wouldn't sway be incredibly advantageous in fine-tuning our model to our validation set. This was confirmed in our experiments.

We qualitatively tested the four models on real users with our web application (see below). We decided to require the '2_init' model and make classifiers for letters a – e and a – k (excluding j) since we expected it might be easier to tell apart between fewer classes.

Not surprisingly, there was a right away, indirect correlation between the validation accuracies attained using the '2_init' model, and also the number of letters we sought to classify (Fig. 4). We attained a validation accuracy of nearly 98% with five letters and 74% with ten.

Confusion matrices

The confusion matrices reveal that our accuracy suffers primarily because of the misclassification of specific letters (e.g. k and d in Fig. 7). Often the classifier gets confused between two or three similar letters or heavily prefers one in all the 2 in such a pair (e.g. g/h in Fig. 7 and m/n/s/t in Fig. 6).

The confusion matrix for the ten-letter model reveals that with the exception of classifying k, it performed reasonably well. We believe that there are two main reasons for this result. First, within the dataset, k is signed from various different perspectives - from the front to the rear of the hand facing the camera - and rotations - with the fingers pointing up and to the perimeters. the sole consistent aspect is that the center mass of the hand altogether pictures, which

is precisely what a resembles. Second, within the range a-k, k it's features that are very the same as letters d, g, h, and that i – it will be constructed by a mixture of segments of those letters. Hence, if the classifiers over-relied on any of those features individually, it might be easy to misclassify k as any of them and would make it harder to be told the latter.

Real-time testing

As aforementioned, we initially tested our four a – y classification models on real users through our web application. This consisted of testing on images during a big selection of environments and hands. A key observation was that there's no significant correlation between the ultimate validation accuracies of the models and their real-time performance on our web application. as an example, 'full_lr' classified an input as a convince half the cases with > 0.99 probability, almost independently of the letter we showed it.

It was very apparent that the '2_init' model outshined the remainder although it produced all-time low validation accuracy. However, it still did not consistently predict the proper letter within the top-5. For this reason, we created the models alluded to above with five (a – e) and ten (a – k, excluding j) classes.

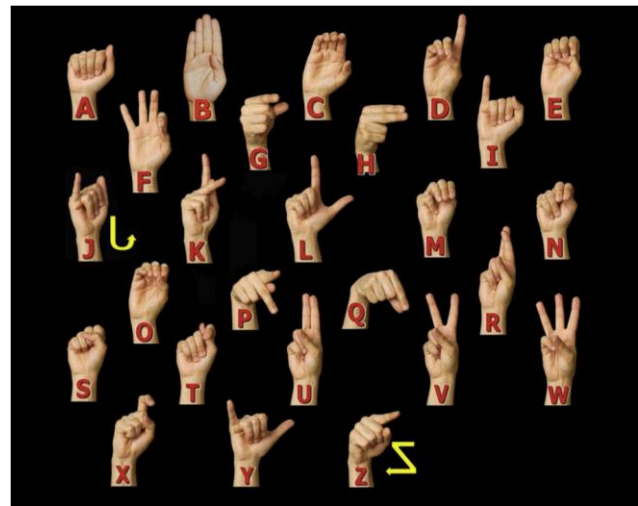
Testing the classifiers with fewer classes on the live web application also yielded significantly different results from testing on the validation set. On ten-letter classification, some letters that were classified correctly in > 70% of cases (e.g. b, c, Fig. 7) were nearly always absent within the top-5 predictions. Moreover, some letters were noticeably overrepresented within the predictions, sort of a and e (Fig. 7). We attribute this to the very fact that neither of those has any fingers protruding from the middle of the hand. they'll thus share the central mass of the hand within the pictures with every other letter without having a contour that might exclude them.

Conclusions and Future Work

Conclusions

We implemented and trained an American language translator on an internet application supported a CNN classifier. We are ready to produce a sturdy model for letters a-e, and a modest one for letters a-k (excluding j). Due to the dearth of variation in our datasets, the validation accuracies we observed during training weren't directly reproducible upon testing on the net application. We hypothesize that with additional data taken in several environmental conditions; the models would be ready to generalize with considerably higher efficacy and would produce a strong model for all letters.

Appendix



American Sign Language finger-spelling alphabet

Reference

1. Mitchell, Ross., Young, Travas, Bachleda, Bellamie, Karchmer and Michael. 2006. "How Many People Use ASL in the United States? Why Estimates Need Updating" (PDF). Sign Language Studies (Gallaudet University Press.) 6 (3). ISSN 0302-1475. Retrieved November 27, 2012.
2. Singha, J. and Das, K. 2003. "Hand Gesture Recognition Based on Karhunen-Loeve Transform", Mobile and Embedded Technology

- International Conference (MECON), January 17-18. India. 365-371.
recognition using k-Nearest Neighbors Classifier. 3rd International Conference on Information and Communication Technology. 533-536.
4. Sharma. R. 2013. Recognition of Single Handed Sign Language Gestures using Contour Tracing descriptor. Proceedings of the World Congress on Engineering 2013 Vol. II, WCE 2013, July 3 - 5, London, U.K.
 5. Starner, T and A. Pentland. 1997. Real-Time American Sign Language Recognition from Video Using Hidden Markov Models. Computational Imaging and Vision, 9(1); 227-243.
 6. Jeballi. M. 2013. Extension of Hidden Markov Model for Recognizing Large Vocabulary of Sign Language. International Journal of Artificial Intelligence & Applications 4(2); 35-42, 2013
 3. Aryanie. D and Y. Heryadi. 2015. American Sign Language-Based Finger-spelling
 7. Suk, H. 2010. Hand gesture recognition based on dynamic Bayesian network framework. Patter Recognition 43 (9); 3059-3072.
 8. Mekala. P. 2011. Real-time Sign Language Recognition based on Neural Network Architecture. System Theory(SSST), 2011 IEEE 43rd Southeastern Symposium 14-16 March.
 9. Admasu, Y. F and K. Raimond. 2010. Ethiopian Sign Language Recognition Using Artificial Neural Network. 10th, International Conference on Intelligent Systems Design and Applications, 2010. 995-1000.
 10. Atwood, J., Eicholtz, M and J. Farrell. 2012. American Sign Language Recognition System. Artificial Intelligence and Machine Learning for Engineering Design. Dept. of Mechanical Engineering, Carnegie Mellon University.
