

# Technical Research

MLP data sets

# Table of Contents

- Open Source Data sets
  - [MNIST](#)
  - [Kaggle's Mathematical Symbols](#)
  - [CROHME data set](#)
  - [HASY data set](#)
- OCR technology
  - [Tesseract OCR](#)
- Character Recognition
  - [Deep Columnar Convolutional Neural Network](#)

# Open Source Data Sets

## MNIST Data Set:

The MNIST Data Set is a Data Set containing over 60 000 samples of training data on handwritten digits, and a test set of over 10 000 samples. It contains images of the digits as well as the corresponding label. This can be really useful when training a model to classify handwritten digits.

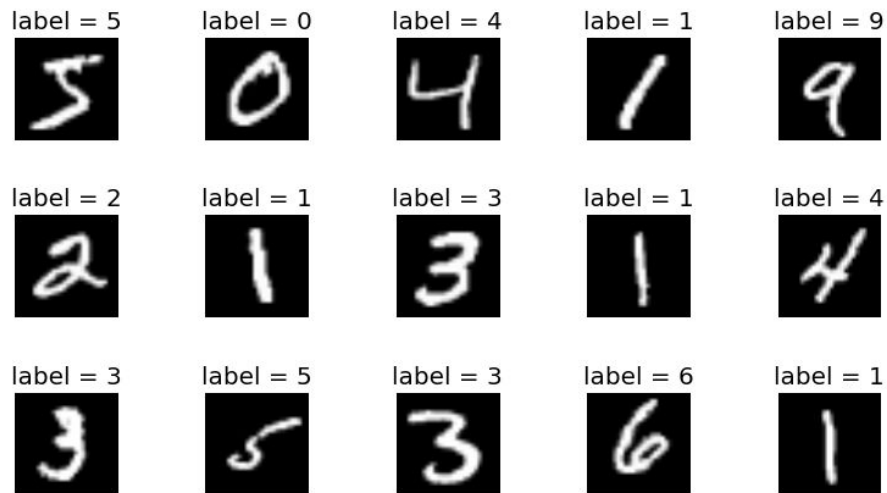
Extracting Data Set:

<http://yann.lecun.com/exdb/mnist/>

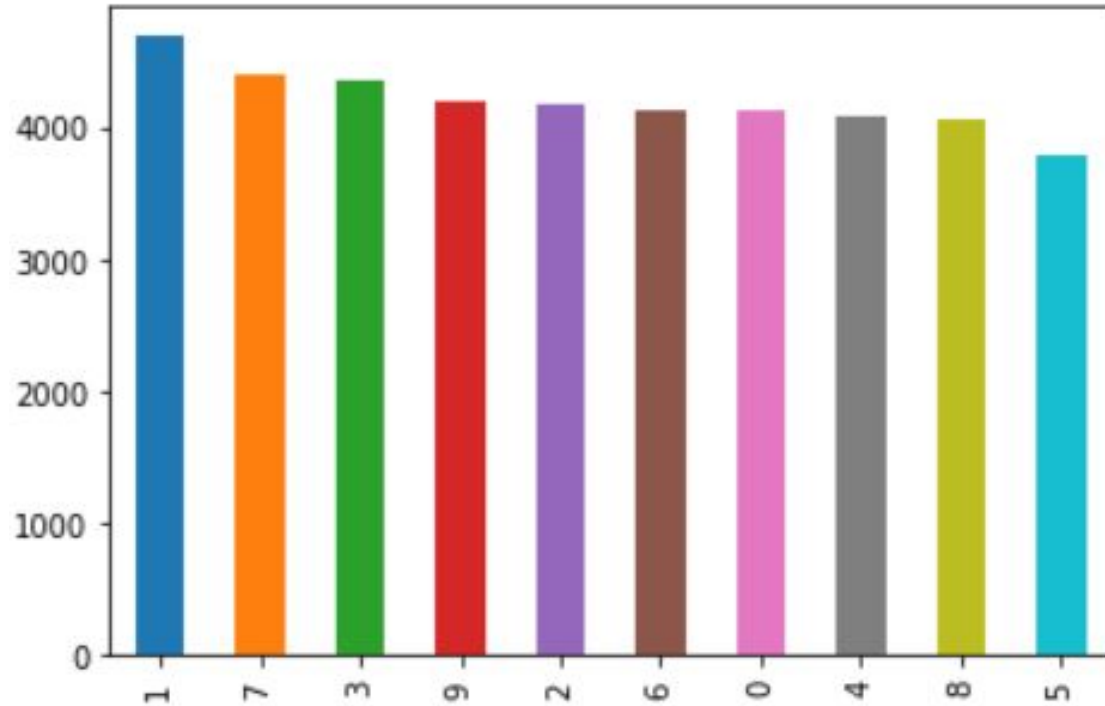
Article:

<https://towardsdatascience.com/image-classification-in-10-minutes-with-mnist-dataset-54c35b77a38d>

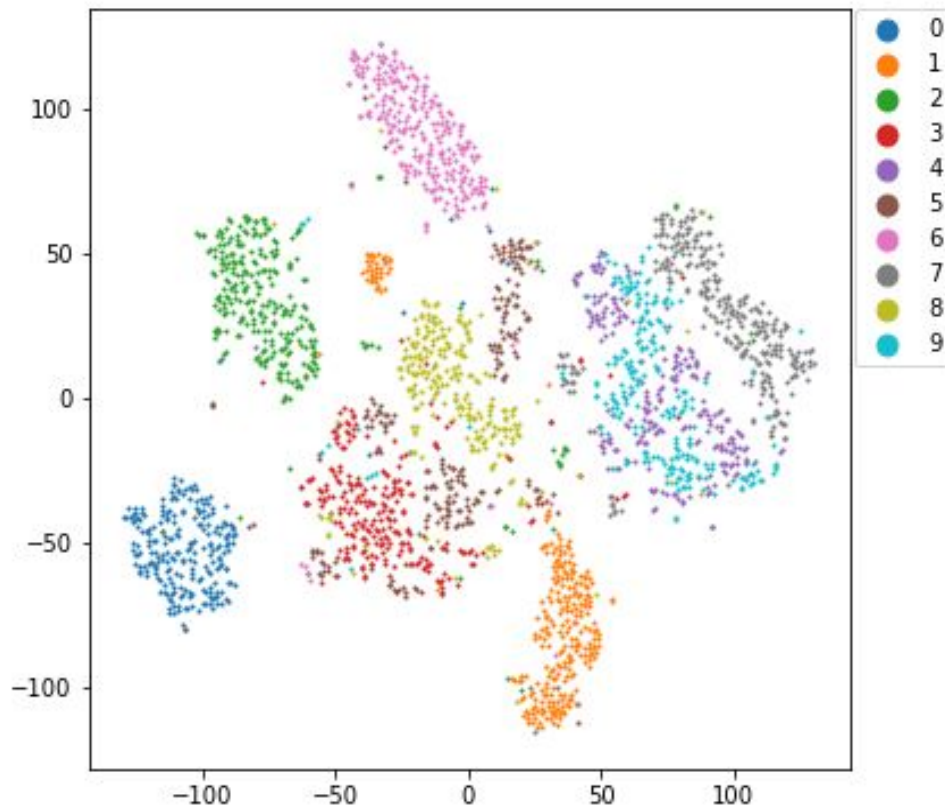
# MNIST Example Data:



# Distribution of Digits in the MNIST Dataset



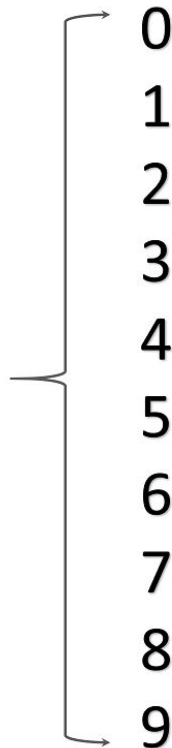
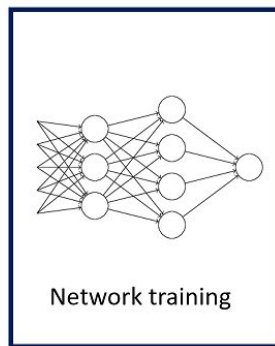
# tsne representation of the MNIST dataset



# MNIST Classification process



Data & Labels



# Kaggle's Mathematical Symbols:

The Kaggle Mathematical Symbols Data Set contains jpg files of: Basic Greek alphabet symbols, English alphanumeric symbols are included, All math operators, Set operators, Basic predefined math functions, Math symbols.

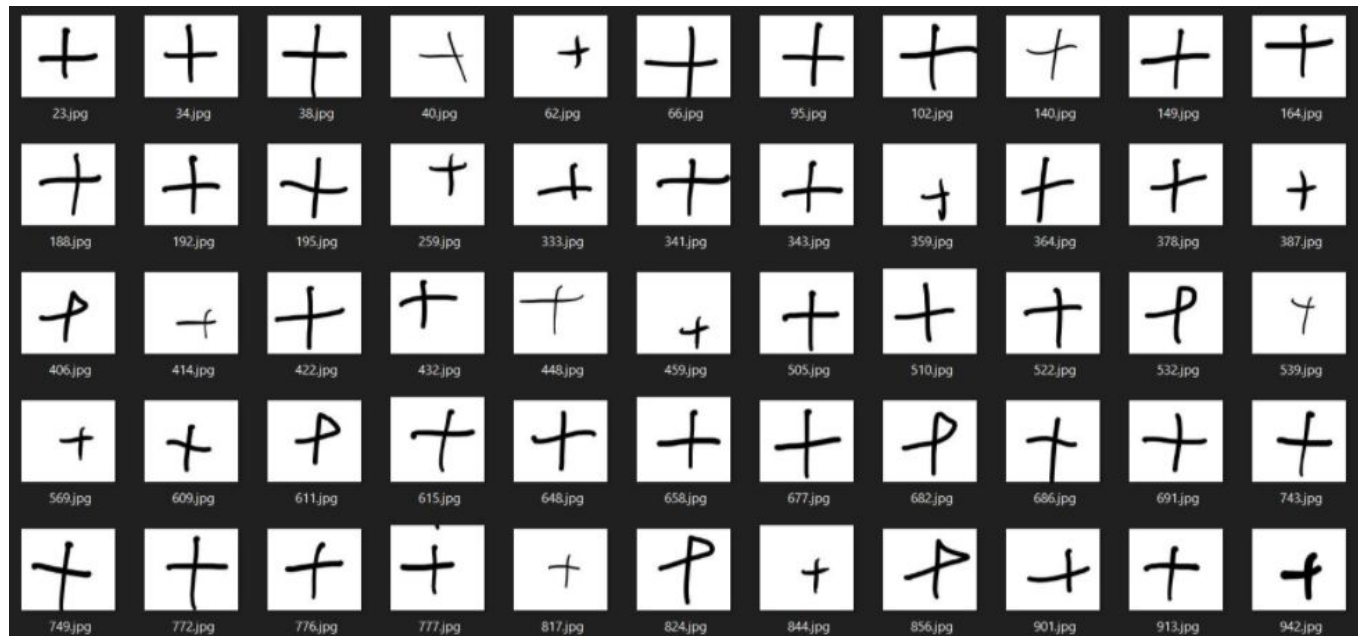
This Data Set can be useful for identifying Math symbols rather than just the digits as in the MNIST Data Set.

Extracting Data Set (410 MB):

<https://www.kaggle.com/xainano/handwrittenmathsymbols>



# Kaggle's Mathematical Symbols Example Data:



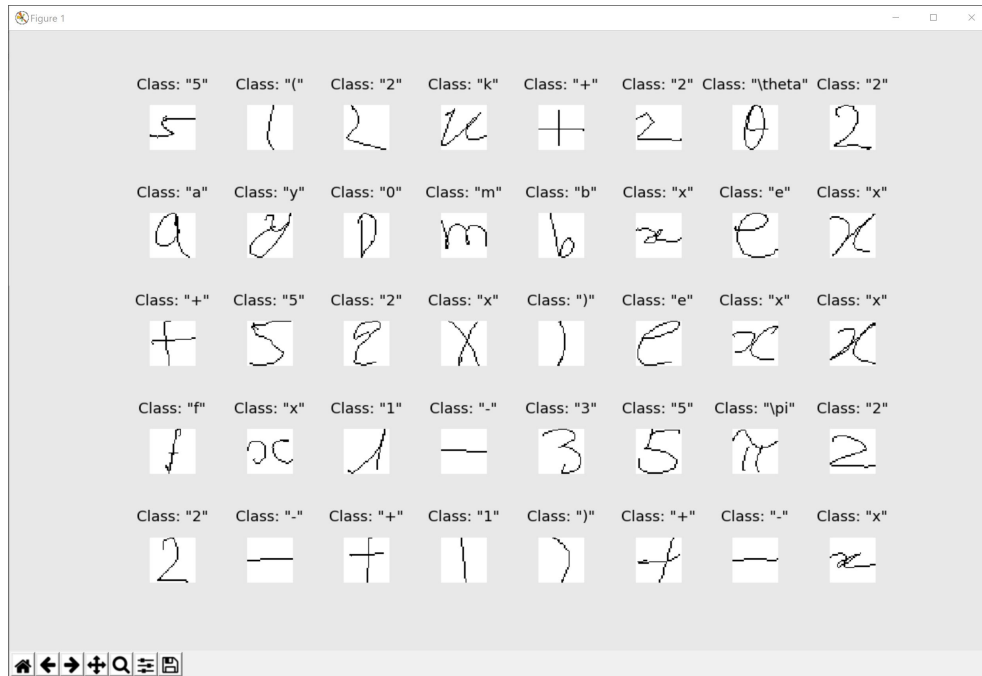
# CROHME data set [\[data\]](#)

CROHME datasets originally exhibit features designed for Online-handwriting recognition task.

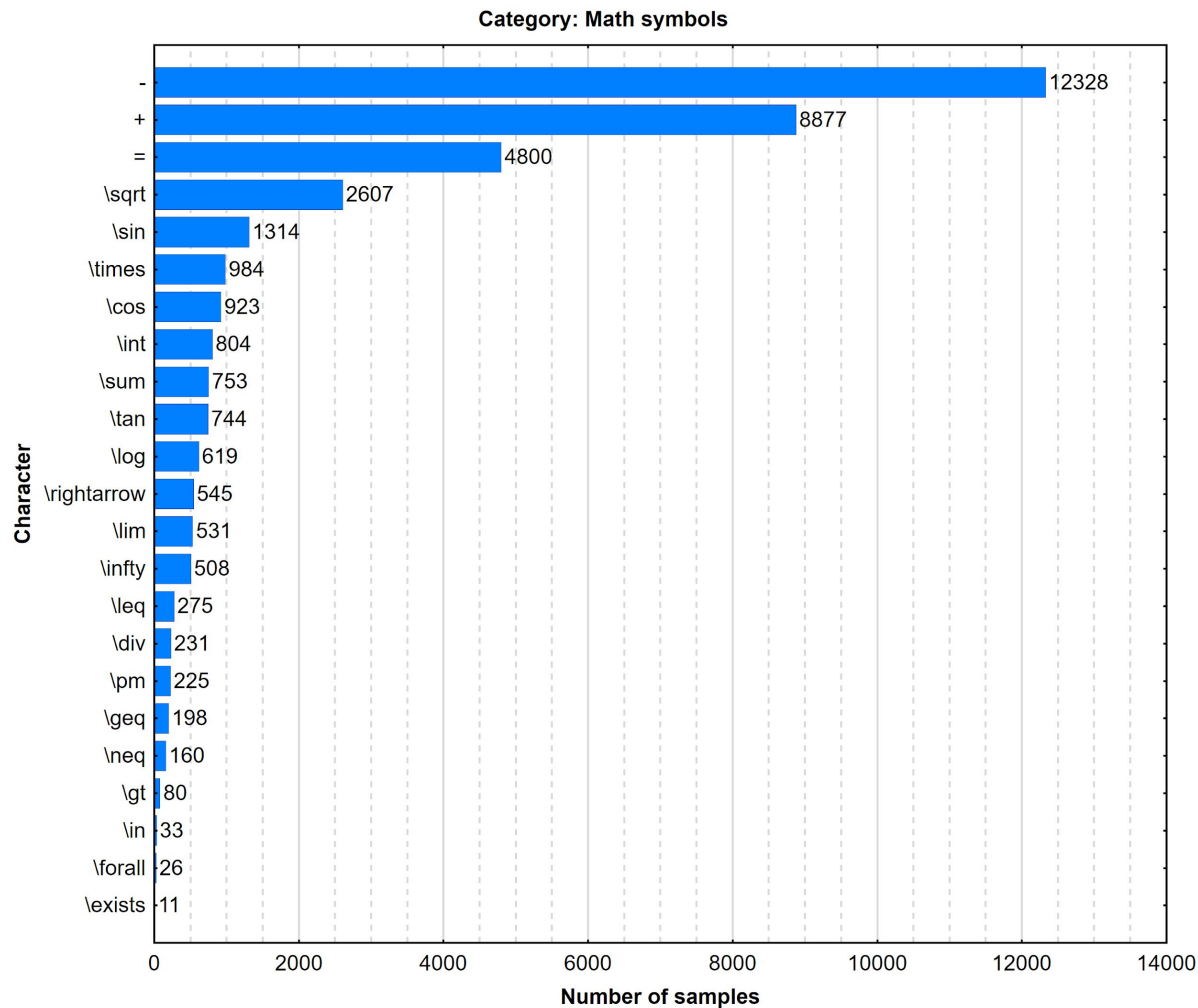
Apart from drawn traces being encoded, inkml files also contain trace drawing time captured.

To extract CROHME data use this github repo

[https://github.com/ThomasLech/CROHME\\_extractor](https://github.com/ThomasLech/CROHME_extractor)

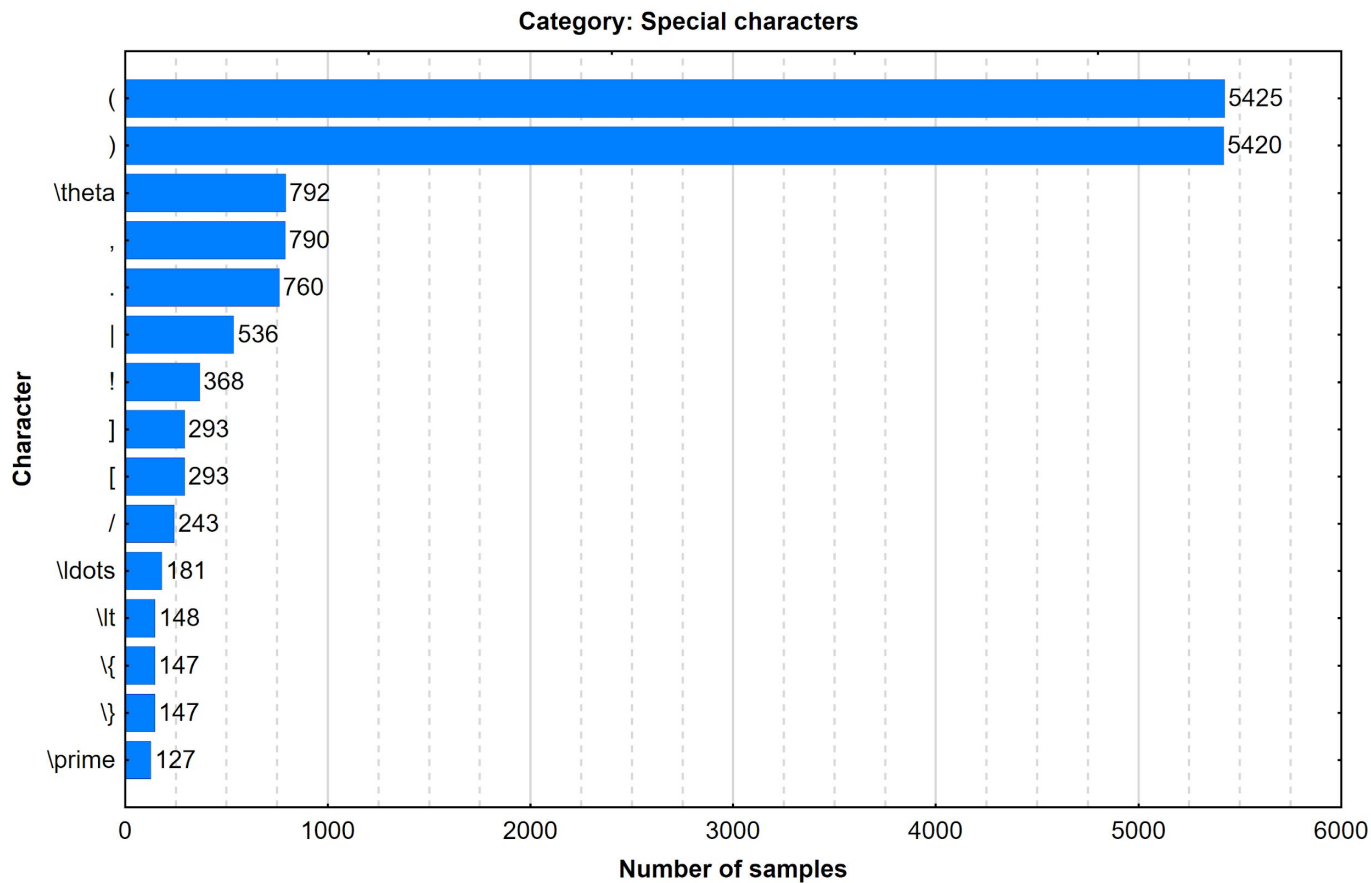


*Sample from CHROHME extractor*



This distribution shows the occurrence of Math symbols in the CROHME data set.

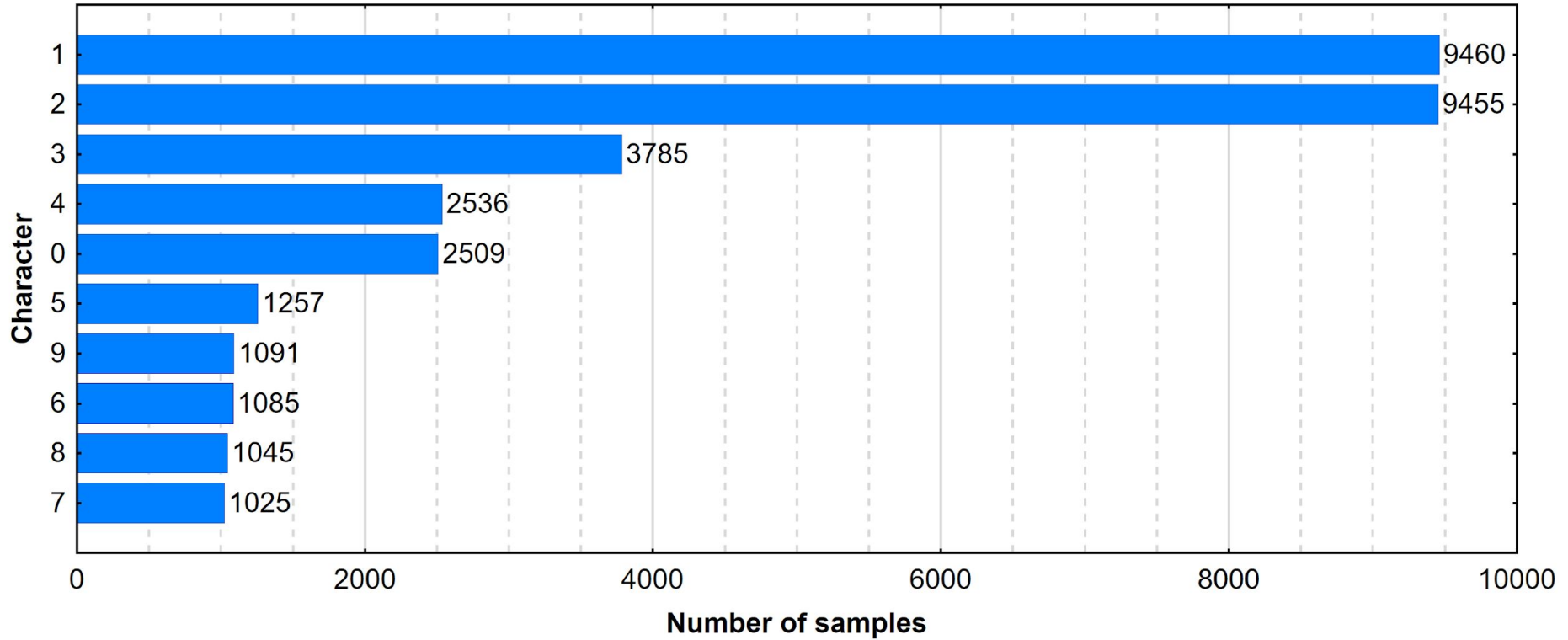
[https://raw.githubusercontent.com/ThomasLech/CROHME\\_extractor/master/histograms/math\\_symbols\\_distribution.png](https://raw.githubusercontent.com/ThomasLech/CROHME_extractor/master/histograms/math_symbols_distribution.png)



*This distribution shows the occurrence of Special characters in the CROHME data set.*

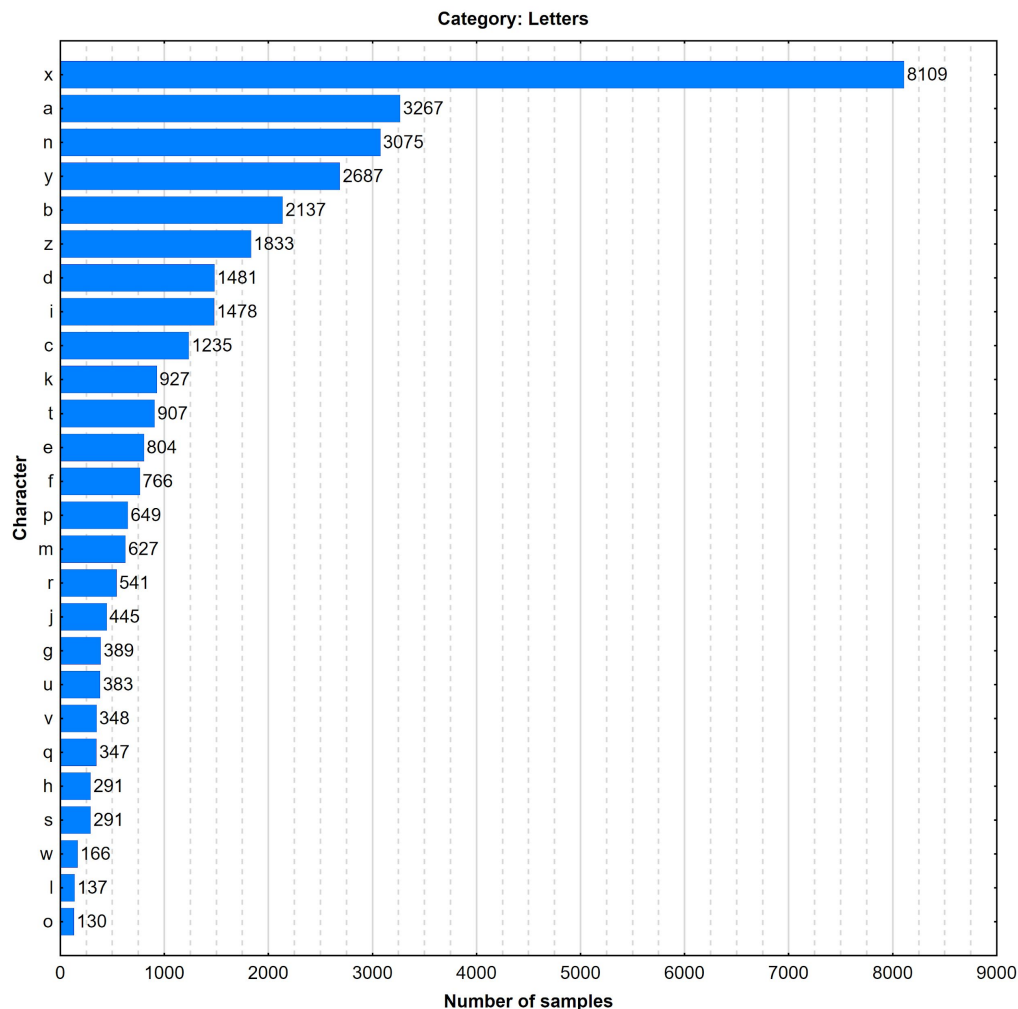
[https://raw.githubusercontent.com/ThomasLech/CROHME\\_extractor/master/histograms/special\\_characters\\_distribution.png](https://raw.githubusercontent.com/ThomasLech/CROHME_extractor/master/histograms/special_characters_distribution.png)

Category: Digits



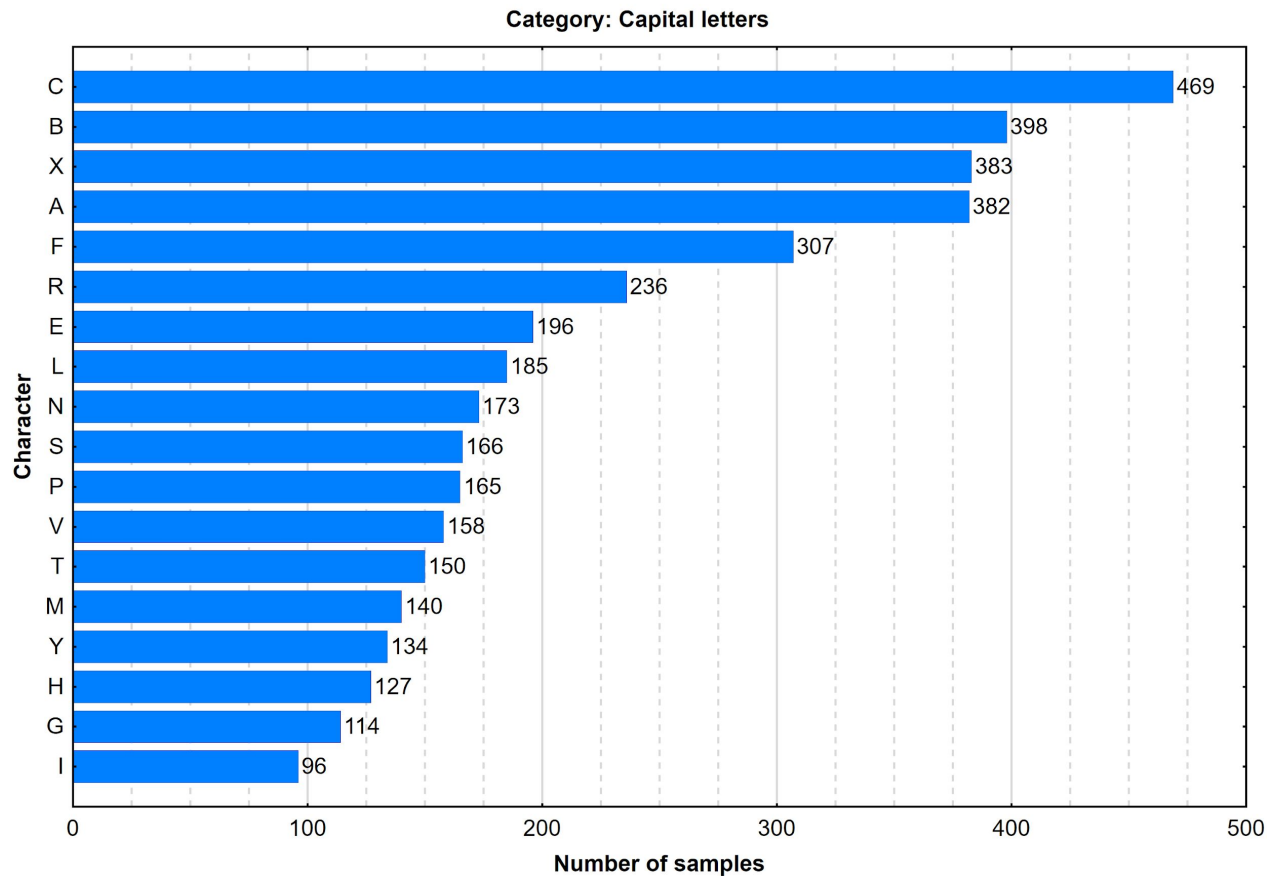
*Histogram showing the distribution of digits in the CROHME dataset*

[https://raw.githubusercontent.com/ThomasLech/CROHME\\_extractor/master/histograms/digits\\_distribution.png](https://raw.githubusercontent.com/ThomasLech/CROHME_extractor/master/histograms/digits_distribution.png)



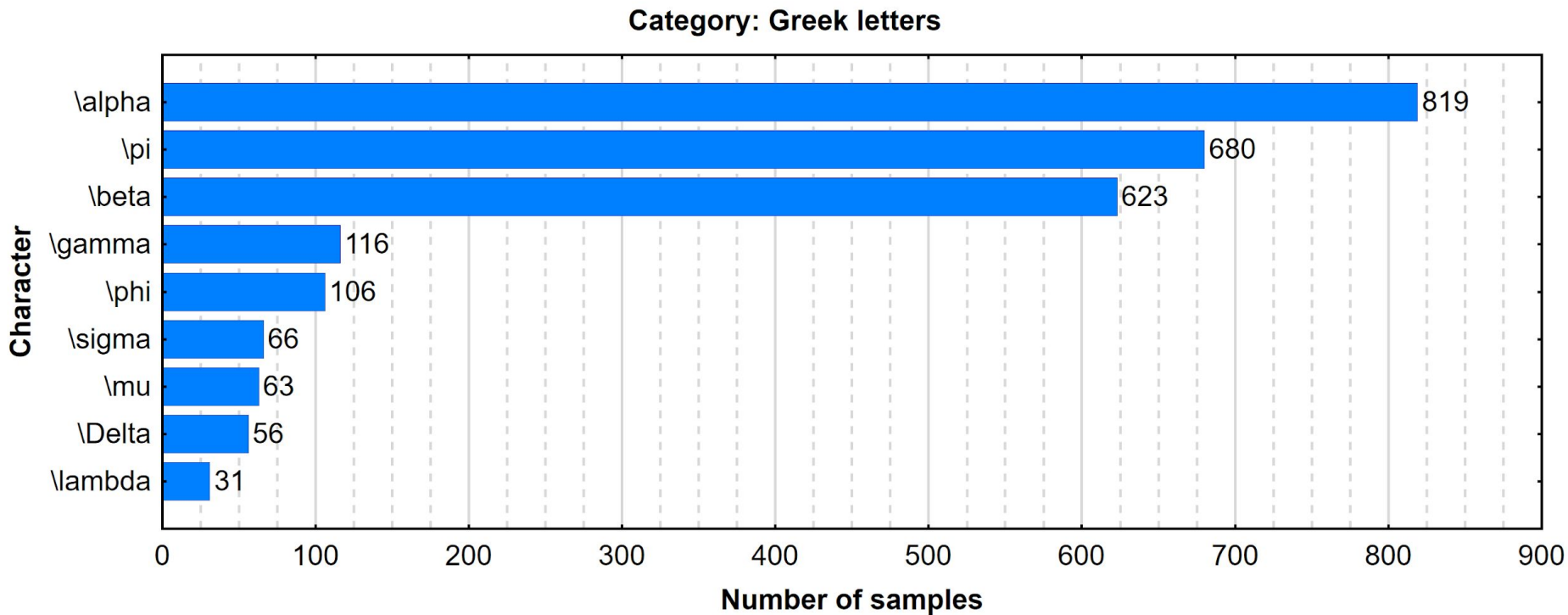
This plots shows the distribution of lower case letters in the CROHME dataset

[https://raw.githubusercontent.com/ThomasLech/CROHME\\_extractor/master/histograms/lowercase\\_letters\\_distribution.png](https://raw.githubusercontent.com/ThomasLech/CROHME_extractor/master/histograms/lowercase_letters_distribution.png)



*Distribution of capital letters in the CHROME dataset*

[https://raw.githubusercontent.com/ThomasLech/CROHME\\_extractor/master/histograms/capital\\_letters\\_distribution.png](https://raw.githubusercontent.com/ThomasLech/CROHME_extractor/master/histograms/capital_letters_distribution.png)

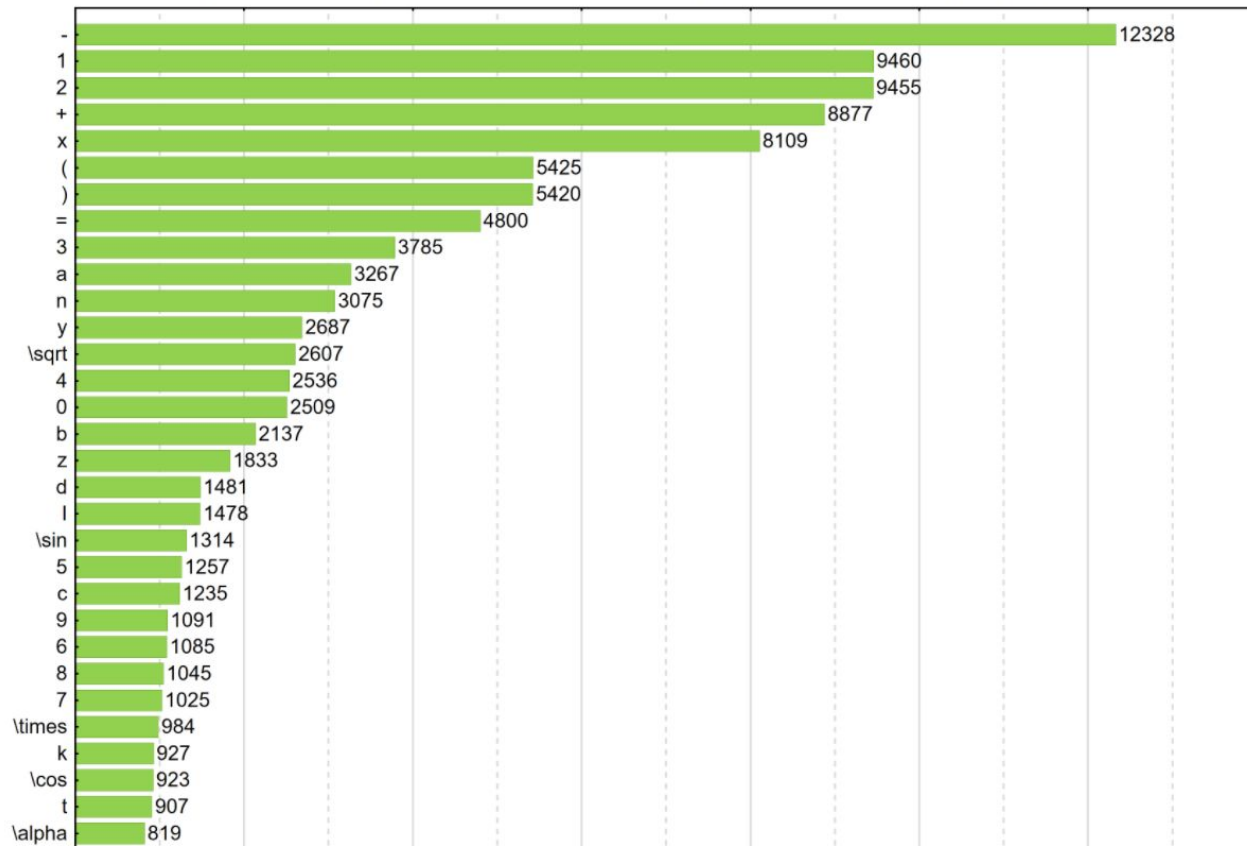


*Distribution of greek letters found in the CROHME dataset*

[https://raw.githubusercontent.com/ThomasLech/CROHME\\_extractor/master/histograms/greek\\_letters\\_distribution.png](https://raw.githubusercontent.com/ThomasLech/CROHME_extractor/master/histograms/greek_letters_distribution.png)



Categories: Math symbols, Greek letters, Capital letters, Special characters, Letters and Digits



Most of the dataset contains:

- Numbers
- Variable  $x$
- $- + x =$
- Brackets  $()$
- $\sqrt$
- Sin, Cos

This plot shows the distribution of symbols in the CROHME data set.

<https://cloud.githubusercontent.com/assets/22115481/26694312/413fb646-4707-11e7-943c-b8eceb0c986.png>

# HASY data set

- HASY is derived from the HWRT dataset.
- HWRT is an on-line recognition dataset, meaning it does not contain the handwritten symbols as images, but as point-sequences.
- HWRT contains strictly more information than HASY.
- The smaller dimension of each recordings bounding box was scaled to be 32 px.
- The image was then centered within the 32 px × 32 px bounding box

- 32px x 32px images
- 369 symbol classes
- over 150,000 instances of handwritten symbols.

Link to a paper on the HASY data set [\[paper\]](#)

“A weakness of HASYv2 is the amount of available data per class.  
For some classes, there are only 51 elements in the test set”

**MNIST has grayscale images while HASY has black and white images.**

## hasy-data (168k files)



*Sample of HASY data found here*

<https://www.kaggle.com/kerneler/starter-hasyv2-df54c04f-b?>

There are 168233 items stored as png files.



Figure 1: 100 recordings of the HASYv2 data set.

The ten classes with most recordings are:

$$\int, \sum, \infty, \alpha, \xi, \equiv, \partial, \mathbb{R}, \in, \square$$

*Same algorithms applied to MNIST and HASY datasets.*

Classifier	Test Accuracy		
	MNIST	HASY	min – max
TF-CNN	99.20 %	81.0 %	80.6 % – 81.5 %
Random Forest	96.41 %	62.4 %	62.1 % – 62.8 %
MLP (1 Layer)	89.09 %	62.2 %	61.7 % – 62.9 %
LDA	86.42 %	46.8 %	46.3 % – 47.7 %
QDA	55.61 %	25.4 %	24.9 % – 26.2 %
Decision Tree	65.40 %	11.0 %	10.4 % – 11.6 %
Naive Bayes	56.15 %	8.3 %	7.9 % – 8.7 %
AdaBoost	73.67 %	3.3 %	2.1 % – 3.9 %

The following observations are noteworthy:

- All algorithms achieve much higher accuracy on MNIST than on HASYv2.
- While a single Decision Tree performs much better on MNIST than QDA, it is exactly the other way around for HASY. One possible explanation is that MNIST has grayscale images while HASY has black and white images.

LaTeX	Rendered	Total	Confused with
$\backslash mid$		34	
$\backslash triangle$	$\Delta$	32	$\backslash Delta$
$\backslash mathds{1}$	1	32	$\backslash mathbb{1}$
$\backslash checked$	✓	28	$\backslash checkmark$
$\backslash shortrightarrow$	→	28	$\backslash rightarrow$
$\backslash Longrightarrow$	⇒	27	$\backslash Rightarrow$
$\backslash backslash$	\	26	$\backslash setminus$
$\backslash O$	Ø	24	$\backslash emptyset$
$\backslash with$	&	21	$\backslash \&$
$\backslash diameter$	∅	20	$\backslash emptyset$
$\backslash triangledown$	∇	20	$\backslash nabla$
$\backslash longmapsto$	↦	19	$\backslash mapsto$
$\backslash dotsc$	...	15	$\backslash dots$
$\backslash fullmoon$	○	15	$\backslash circ$
$\backslash varpropto$	α	14	$\backslash propto$
$\backslash mathsection$	§	13	$\backslash S$
$\backslash vartriangle$	Λ	12	$\backslash Delta$
$O$	<i>O</i>	9	$\backslash circ$
$o$	<i>o</i>	7	$\backslash circ$
$c$	<i>c</i>	7	$\backslash subset$
$v$	<i>v</i>	7	$\backslash vee$
$x$	<i>x</i>	7	$\backslash times$
$\backslash mathbb{Z}$	<b>Z</b>	7	$\backslash mathds{Z}$
$T$	<i>T</i>	6	$\backslash top$
$V$	<i>V</i>	6	$\backslash vee$
$g$	<i>g</i>	6	9
$l$	<i>l</i>	6	
$s$	<i>s</i>	6	$\backslash mathcal{S}$
$z$	<i>z</i>	6	$\backslash mathcal{Z}$
$\backslash mathbb{R}$	<b>R</b>	6	$\backslash mathds{R}$
$\backslash mathbb{Q}$	<b>Q</b>	6	$\backslash mathds{Q}$
$\backslash mathbb{N}$	<b>N</b>	6	$\backslash mathds{N}$

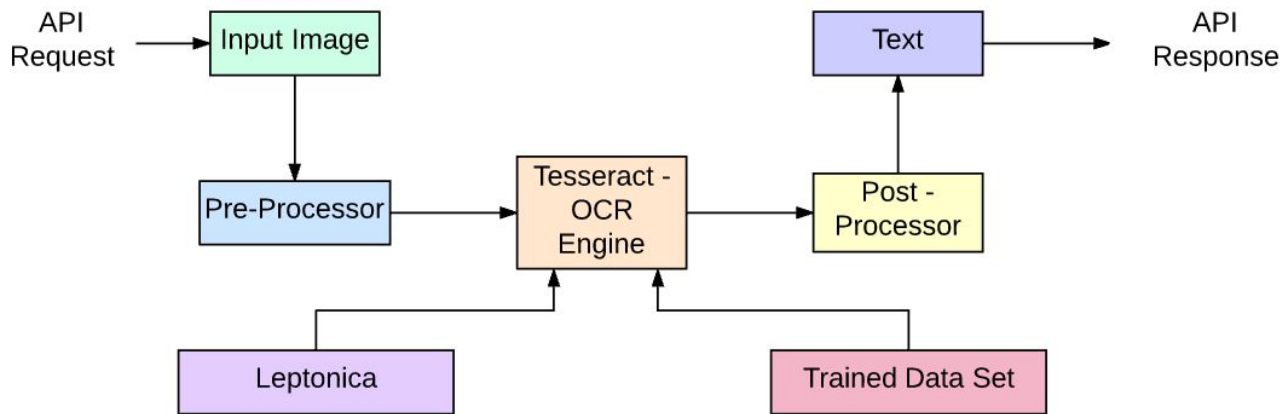
Table IV: 32 classes which were not a single time classified correctly by the best CNN.

# Tesseract OCR:

- Tesseract is an open source text recognition (OCR) Engine.
- It can be used directly, or using an API.
- It supports multiple languages
- Web Page: <https://nanonets.com/blog/ocr-with-tesseract/>
- Installation with Python: `pip install pytesseract`

# Tesseract OCR Process Flow

OCR Process Flow



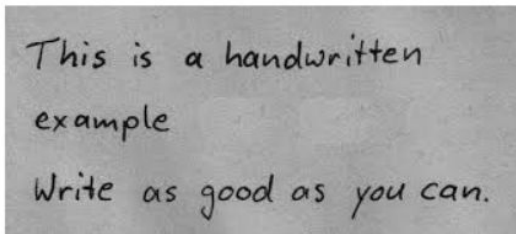
# PyTesseract

This is simple python optical character recognition.

- Uses PyTesseract Library
- Webpage: <https://stackabuse.com/pytesseract-simple-python-optical-character-recognition/>
- Webpage contains python code that can be implemented.
- More so focused on Strings than math, but the concept is still there and could be modified.
- Webpage uses flask to deploy there code.

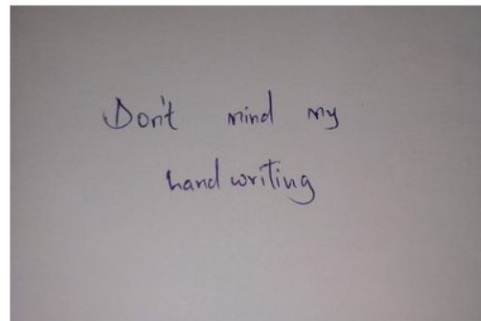
# Example Outputs:

## Result:



The extracted text from the image above is: **This is a handwritten example Write as geol as you can.**

## Result:



The extracted text from the image above is: **ad oviling**

The Difference between the two outputs is really quite something, showing that handwriting does in fact play a crucial role, we could try to improve these by training our model on the datasets we discussed earlier.



# Example Outputs:

## Result:

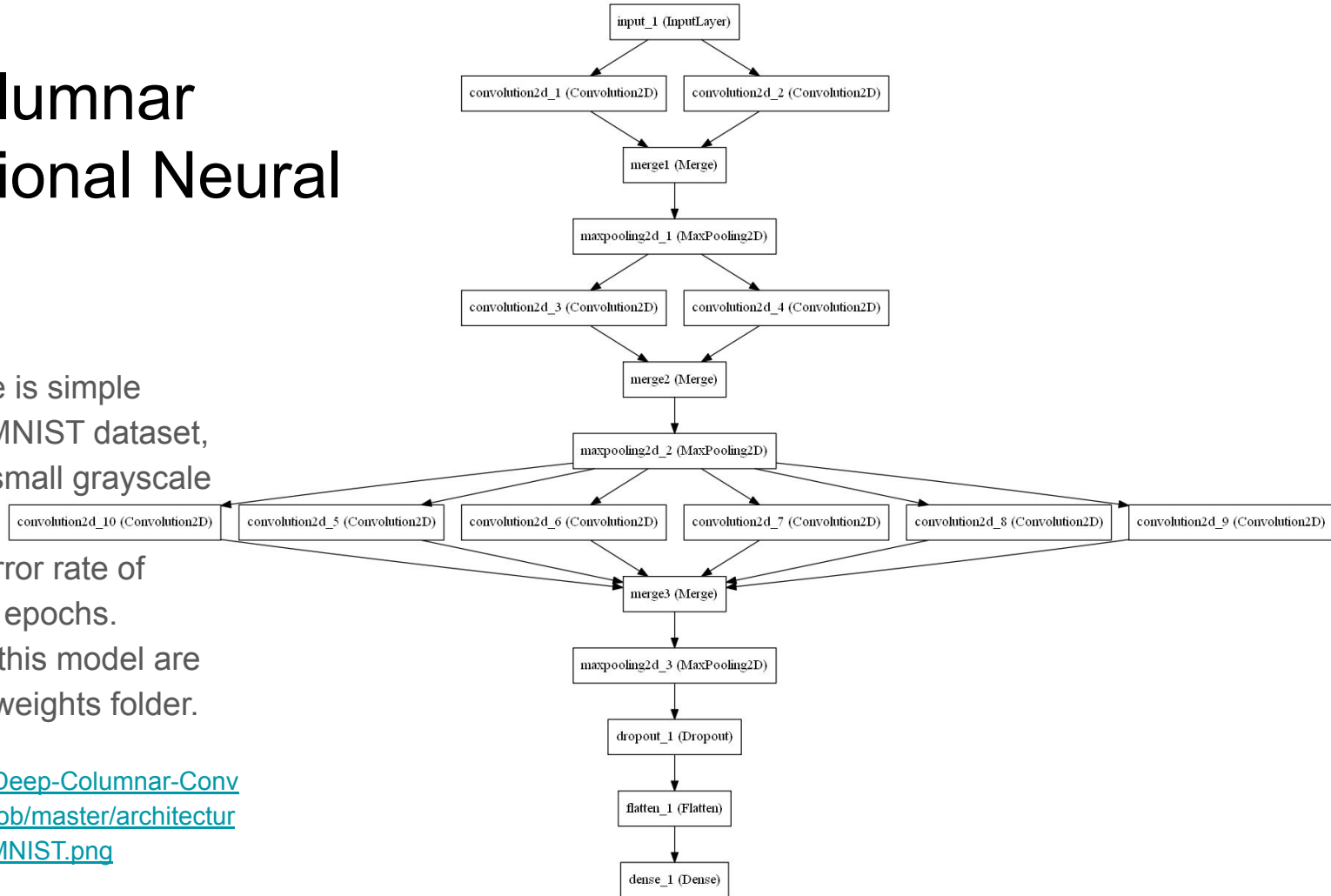


The extracted text from the image above is: **CETTE ALARME EST LOCALE SEULEMENT, EN CAS DE FEU SIGNALER 911**

Included this image to show how it did in fact extract a number being 911, however this is a lot different to a handwritten image. The difference being though that we would have trained our model on the previous datasets.

# Deep Columnar Convolutional Neural Network

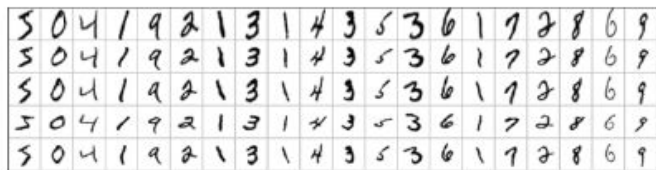
- This architecture is simple enough for the MNIST dataset, which contains small grayscale images.
- It achieves an error rate of 0.23% after 500 epochs.
- The weights for this model are available in the weights folder.



<https://github.com/titu1994/Deep-Columnar-Convolutional-Neural-Network/blob/master/architectures/DCCNN%20MNIST.png>

# How does it work?

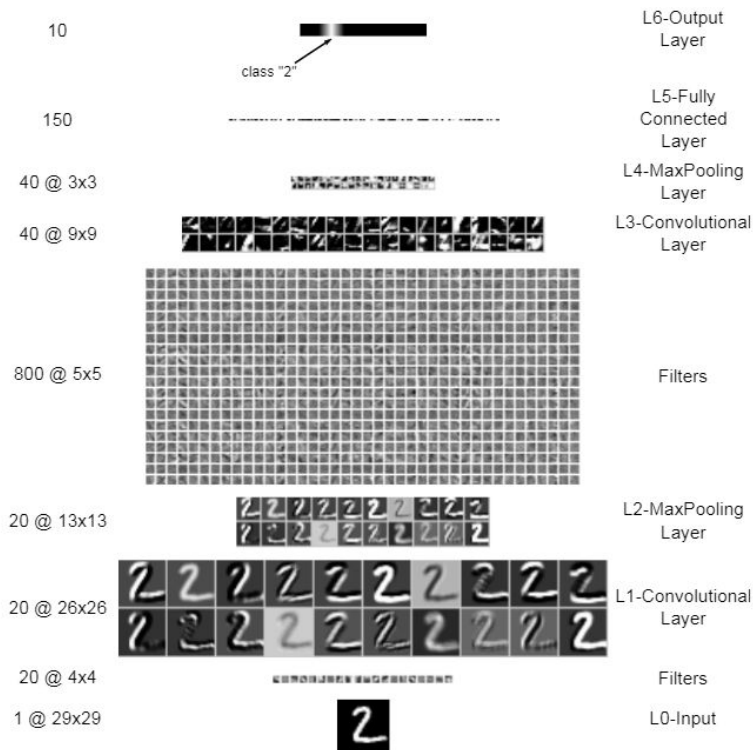
[\[notebook\]](#)



*Handwritten digits from the training set (top row) and their distorted versions after each epoch (second to fifth row)*



*The 23 errors of the MCDNN, with correct label (up right) and first and secondbest predictions (down left and right)*



*DNN architecture for MNIST. Output layer not drawn to scale; weights of fully connected layers not displayed.*

# Results

Table 4. Average error rates of MCDNN for all experiments, plus results from the literature. \* case insensitive

Data (task)	MCDNN error [%]	Published results Error[%] and paper	
all (62)	<b>11.63</b>		
digits (10)	<b>0.77</b>	3.71 [12]	1.88 [23]
letters (52)	<b>21.01</b>	30.91[16]	
letters* (26)	<b>7.37</b>	13.00 [4]	13.66[16]
merged (37)	<b>7.99</b>		
uppercase (26)	<b>1.83</b>	10.00 [4]	6.44 [9]
lowercase (26)	<b>7.47</b>	16.00 [4]	13.27 [16]

There seems to be a larger error when classifying letters compared to digits