

UvA Applied Machine Learning Course

Group project

Project setup

- Participate in a Kaggle competition by solving a real-world ML task
- Submit your evaluation results directly to Kaggle and see the leaderboard updated in real time!
- Choose one of three project topics:
 - Project 1: Read between the lines
 - Project 2: Feathers in focus
 - Project 3: Unified tabular learning

Your tasks

1. Register a team of 3 on Canvas by Monday 17 Nov at 12:00pm
If you don't have a team, we will find you one. We will allocate any non-registered students to teams on Monday afternoon
2. Decide which project topic you want to do with your team.
3. [Start with a pretrained model from HuggingFace](#), this will serve as your baseline
Be sure to name it “baseline” when submitting results to Kaggle
4. Develop your own ML model and try to beat the baseline!
Note: it is unlikely you will actually beat the baseline, but it's good to have a sense of the upper bound on performance that's possible

Deliverables

- Final deliverable: poster with main outcomes, findings, novelties, etc.
No final report!
- Include analysis of computational complexity of your model vs your baseline
Baseline will probably have better results, but at what cost? Much more compute/flops/parameters/etc
- Each poster is graded by two people.
- Lab sessions vital for progress! Communicate often with many TAs.

Project 1: Read between the lines

Reading comprehension with logic

Most people who are skilled banjo players are also skilled guitar players.
But most people who are skilled guitar players are not skilled banjo players.

Q: If the statements above are true, which one of the following must also be true?

A1: There are more people who are skilled at playing the guitar than there are people who are skilled at playing the banjo.

A2: There are more people who are skilled at playing the banjo than there are people who are skilled at playing the guitar.

A3: A person trying to learn how to play the guitar is more likely to succeed in doing so than is a person trying to learn how to play the banjo.

A4: There are more people who are skilled at playing both the guitar and the banjo than there are people who are skilled at playing only one of the two instruments.

<https://www.kaggle.com/t/a2e3f70c477a48dba2627d2cf42e699b>

TH

White throated Sparrow



Green Jay



White breasted kingfisher



Yellow bellied flycatcher



Final project 2: Feather in Focus

Classifying images of bird species

<https://www.kaggle.com/t/0e9856f5cb5f40af8739be017cc75b9b>

Project 3: Unified tabular learning

Learning to classify multiple tabular datasets



Forest Cover Dataset



Credit Card Fraud Dataset



Bank Marketing Dataset

<https://www.kaggle.com/t/acbc4bb2ee8149e6a74e808c9795794d>

Things to include on your posters

- Main research question: what is the problem we're trying to solve?
 - Hint: RQs are questions that can be answered concretely i.e., "yes", "no", or a number
- Figure that explains your model architecture
 - Include explanation of how an individual sample moves through the architecture
- Comparison to baselines, both simple and complex
 - Be able to explain the design choices you made
- Error analysis: where does your model get it wrong and why?
 - Show 1-2 examples

Spam Filter

CIVILIA DONKER- C.I.E.DONKER@STUDENT.VU.NL

ROOS SLINGERLAND- ROOS.SLINGERLAND@STUDENT.UVA.NL



The Project

PREVIOUS RESEARCH

- Spam: sell product or services to customers available on the internet via email, also bulk-email [7]
- Because of the increase of email use, bulk-email increased as well [4]
- Research is often done, but spam keeps developing [4] and labelled data is often an issue [7]
- Length could be an indicator of spam [5]
- Metadata such as
- Mail is often forged
- Decision trees provide a good classification field [7, 8]

ABOUT THE DATA

4021 training examples

24% spam
76% not spam

Identifying Quora Question Pair Duplicates

Imaz Binsbok, Tom Dap (Quora-The-Explorer)

1. The Problem



Our Ensemble Approach

Utilizing a combination of LSTMs (Long Short-Term Memory) and MLP features, feeding two various classification algorithms to predict duplicate questions.

2. Data



2.1 Preprocessing



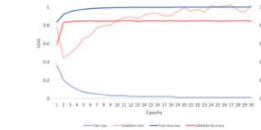
2.2 Embeddings



3. Feature Engineering



4.2 Model Performance



5. Classification Results

After combining outputs from LSTM and engineered features, we tried several different algorithms for classifying duplicate questions and compared results (Accuracy).

Neural Network	3 Dense Layers 200 nodes each, 20% dropout	0.84515
Logistic Regression		0.84585
Random Forest	max depth 2	0.82639
SVM	rd kernel, gamma=1	0.81001

6. Summary

- Although we achieved a reasonable final results (84.6% accuracy), our LSTM performance graph, suggests the model was overfitting.
- While accuracy was used to judge the final result, precision and recall may have been a better measure of performance due to imbalanced classes.
- Although we attempted to tune our model hyperparameters, we were limited by time and computational power and so were only able to test each model for 5 epochs.

PLANKTON IMAGE CLASSIFICATION

AUTOMISATION OF THE PLANKTON IMAGE IDENTIFICATION PROCESS BY MAKING USE OF MACHINE LEARNING TECHNIQUES

DATASET

24304 training images
6132 test images
~30px x 30px smallest
~400px x 400px biggest
Classes not uniformly distributed

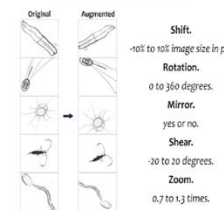
PREPROCESSING



SOFTWARE



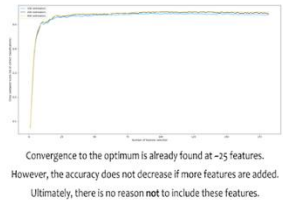
DATA AUGMENTATION



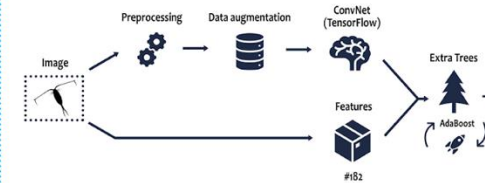
FEATURE EXTRACTION

182 FEATURES, AMONG WHICH:
Centroid
Aspect ratio
Local Binary Patterns
Hu and Zernike moments
Parameter Free Threshold Adjacency Statistics
Mean distance (and σ) to center of image
Number of filled pixels
Haralick's Features
Orientation
Solidity

OPTIMAL NUMBER OF FEATURES (EXTRA TREES)



MACHINE LEARNING MODEL



PROCESS

NETS
experimented with:
Random Forest
AdaBoost
Extra Trees
Logistic Regression
Multi Layer Perceptron
ConvNet (based on VGGNet)

DECISION TREE
BASED NETS
experimented with:
Number of estimators, more is better.
Number of features, 25 or more.

DETAILS CONVNET

LAYER TYPE	SIZE
Convolution	2D 3x3 filter
Convolution	4x4 filter
Max pooling	2x2 with stride of 2
Convolution	6x6 filter
Convolution	2D 3x3 filter
Max pooling	2x2 with stride of 2
Convolution	10x10 filter
Convolution	10x10 filter
Convolution	10x10 filter
Max pooling	2x2 with stride of 2
Flattening	8192x100
Dense + dropout	512
Dense + dropout	256
Length	1

RESULTS

82.7%
CORRECT
PREDICTIONS
Kaggle rank #2

CONVNET
experimented with:
Size of images, bigger is better.
Less or more layers, more is better.
Number of filters, more is better.
Size of filters, 3x3 is best.
Type of pooling, does not really matter.
Activation functions, Leaky ReLU.
Learning rates, decreasing over time.
Different optimizers, Ada Gradient.