

Report

Computer lab. OCR learned by Structured Output Perceptron

Morhunenko Mykola.

	Testing error in %	
	R_seq	R_char
independent multi-class classifier	70.6	26.39
structured, pair-wise dependency	11.8	4.643
structured, fixed number of sequences	1	0.609

The difference in performance caused by the inner structure of predictors: if remembering every single letter, it's the worse, not all of the information is utilized. Considering the pairwise dependencies between words is already much better because it reduces the probability of some obvious mistakes as two letters can not be near each other at all.

And remembering all vocabulary is the best for a small dataset because even during the first iteration with enough training data it can learn a lot.

The main advantage of an "independent multi-class classifier" is its simplicity.

The main advantage of an "structured, fixed number of sequences" is its performance if there is a small number of sequences to recognize.

If there are a lot of them, maybe "structured, pair-wise dependency" will be a much better choice, because it can remember subsequences (already better than just letterwise recognition) and is not influenced by the vocabulary size as the last one.