

Performance Prediction for Handwritten Word Recognizers and Its Application to Classifier Combination

Hanhong Xue and Venu Govindaraju
CEDAR, SUNY at Buffalo, Buffalo, NY 14260
{hxue,govind}@cedar.buffalo.edu

Abstract

This paper introduces a performance prediction model for handwritten word recognizers. This model considers the factors involved in word recognition, i.e. the recognizer, input images and lexicons, and presents a quantitative formula to associate performance with these factors. It produces a direct measure of recognition difficulty by predicted performance which can be utilized to improve the combination of multiple recognizers. We support the accuracy of our model by extensive experiments conducted on five word recognizers and its applications to multiple classifier systems.

1 Introduction

In handwritten word recognition, the performance of a recognizer is affected by many different factors including the nature of the recognizer, the quality of input images and the complexity of lexicons. It is helpful to know how a recognizer performs under different circumstances, so its performance can be improved by giving more considerations to its weakness. Or, when we have several recognizers, each recognizer can be assigned to the circumstances where it performs the best compared to others, so that higher overall performance can be achieved. Therefore, studying the influence of different factors on the performance of recognizers gives chances of improving them and methods of combining them.

A common expectation is that recognition with larger lexicons is more difficult because larger lexicons are more likely to contain similar words that confuse recognizers. Besides lexicon size, similarity between lexicon words is also considered as an important factor that poses difficulty to a recognition task [10, 4]. String edit distance, defined as the minimum number of insertions, deletions and substitutions to convert one string to another, is a recognizer-independent measure of word similarity.

Image quality is crucial to any image pattern recognition tasks including handwritten word recognition. One possible approach to characterizing image quality is the use of parameterized image defect models [1], where image size, resolution, skew, blur, binarization threshold, pixel sensitivity and other parameters are used to characterize image quality and to generate pseudo-images. The defect models have been applied to the evaluation of OCR accuracy on synthetic data [5]. However, to the best of our knowledge, there has been no application reported on the evaluation of handwritten word recognizers.

The common part of most previous researches on word recognition difficulty is the lack of a performance prediction model that combines the influences of recognizers, input images and lexicons. Data collected from extensive experiments is capable of showing the tendency of performance change when the circumstance changes. However, since there is no quantitative model attached to the data, it is not known whether the performance change is linear, quadratic or something else. So, the exact answer is actually hidden in the data and it will be inaccurate and inconvenient to use data plots to predict the performance of word recognizers especially when parameters are out of the plot range.

There are also extensive researches on combining multiple classifiers. Techniques reported for handwriting recognition include voting, Borda count [6], logistic regression [7], Bayesian [11] and Dempster-Shafer theory [11]. According to Xu and his colleagues [11], combination of multiple classifiers has three types according to the three levels of classifier's output information: a single class label, rankings and rankings with measures/scores. In our study, available recognizers are based on completely different methodologies and they output best word candidates according to incomparable measures. Therefore, focus is placed on the combination of rank-based decisions in this paper.

Realizing the importance of measuring recognition difficulty, we propose a model to consider the nature of the recognizer, the quality of input images and the complexity of lexicons in predicting performance which gives a direct measure of recognition difficulty.

2 The Performance Model

The performance function of word recognition is defined as $p(R, I, L)$ where R is the recognizer, I the image and L the lexicon. R , I and L can be also viewed as three sets of factors/parameters that characterize the recognizer, the image and the lexicon, respectively. When the parameters related to the recognizer and the image are obtained through a training procedure, the performance function is rewritten as $p_{R,I}(L)$, predicting performance given lexicon. From now on, the term *performance* refers to the top-choice accuracy rate when there is no rejection.

Consider the recognition process as a tournament where non-truths are matched against the truth. Let $t \in L$ be the truth and all other entries in L be the challengers against t . Denote the edit distance between a word w and the truth t as $d(w, t)$. Each of the $d(w, t)$ edit operations is considered as an evidence of t being the truth and w being the non-truth. If the recognizer is aware of at least one such evidence, t is ranked higher than w and thus t beats w . Let q be the probability of one edit operation being ignored by the recognizer and assume equal importance of all edit operations. Then, the probability that t beats w is $1 - q^{d(w,t)}$ when not all edit operations are ignored. In order for t to be the top choice, t must beat all $w \in L - \{t\}$.

Now consider the tournament in an incremental way where random entries are added one by one and a new entry w is qualified to challenge the truth only when w is better than all entries that have been defeated by the truth. Suppose the list of entries defeated by the truth is F . The chance of w being the best in $F \cup \{w\}$ is $1/|F \cup \{w\}|$ since all entries in $F \cup \{w\}$ are random.

To simplify the analysis, we assume that all entries in $L - \{t\}$ have the same edit distance to the truth, i.e. the average edit distance $\bar{d}(t) = \frac{1}{n} \sum_{w \in L - \{t\}} d(w, t)$. Let $p(n)$ denote the recognizer's performance on a lexicon of size n and $r = q^{\bar{d}(t)}$. We have $p(1) = 1$ because a lexicon of size 1 contains only the truth. When $n > 1$, there is $\frac{1}{n-1}$ chance that the n -th entry challenges the truth and the probability that the truth wins is $1 - r$. So the probability that the truth is still the top choice after the addition of the n -th entry is $\frac{1}{n-1}(1 - r) + \frac{n-2}{n-1} = 1 - \frac{r}{n-1}$. Therefore, $p(n)$ can be computed by

$$p(n) = \begin{cases} 1 & , n = 1 \\ p(n-1)(1 - \frac{r}{n-1}) & , n > 1 \end{cases} \quad (1)$$

When $n > 1$, we have

$$p(n) = (1 - \frac{r}{1})(1 - \frac{r}{2}) \dots (1 - \frac{r}{n-1}) = \frac{\Gamma(n-r)}{\Gamma(n)\Gamma(1-r)}$$

The Γ function is the well-known extension of factorial to non-integer values and it has following properties, $\Gamma(x+1) = x\Gamma(x)$ and $\Gamma(n+1) = n!$, where x is a real number

and n is an integer. When $x \rightarrow \infty$, we apply the Stirling's asymptotic formula

$$\Gamma(x+1) = \sqrt{2\pi x} \left(\frac{x}{e}\right)^x \left(1 + \frac{1}{12x} + \dots\right) \simeq \sqrt{2\pi x} \left(\frac{x}{e}\right)^x$$

and get

$$p(n+1) \simeq n^{-r} / \Gamma(1-r) = e^{-r \ln n + c} \quad (2)$$

where $c = \ln \frac{1}{\Gamma(1-r)}$, for $n \rightarrow \infty$.

To compensate for the simplifying assumptions leading to Equation 2, we introduce more parameters into the performance function and conjecture that the performance function has the form

$$p_{q,k,a,t}(L) = (e^{-r})^{f(n)} \quad (3)$$

where $f(n) = k \ln^a n$ and c is ignored to meet the initial condition $p(1) = 1$.

For testing images whose truths are unknown, $\bar{d}(t)$ has to be approximated by the average edit distance between any two entries and the performance function is re-written as

$$p_{q,k,a}(L) = (e^{-q^D})^{f(n)} \quad (4)$$

where $D = \frac{1}{n(n-1)} \sum_{w,v \in L} d(w,v)$ and $f(n) = k \ln^a n$.

The advantage of such a model is that it can be converted into a multiple regression model.

$$p = (e^{-q^D})^{k \ln^a n} \Leftrightarrow \ln(-\ln p) = D \ln q + a \ln \ln n + \ln k$$

Suppose we have a set of observations $\{p_i, n_i, D_i\}$. Let $P_i = \ln(-\ln p_i)$ be the dependent variable, $N_i = \ln \ln n_i$ and D_i be the independent variables, and $\ln q$, a and $\ln k$ be the regression parameters. We get a multiple regression model

$$P_i = (\ln q)D_i + aN_i + \ln k + e_i = \hat{P}_i + e_i \quad (5)$$

where \hat{P}_i is the predicted performance and e_i is the residual.

Now we examine the arguments and parameters in the model. Firstly, the model takes the lexicon as its input argument, considering both lexicon size n and the average edit distance D between lexicon entries. Secondly, q , a and k are model parameters where q is defined as the probability that an edit operation between the truth and the non-truth is ignored by the recognizer and $f(n) = k \ln^a n$ indicates the recognizer's sensitivity to lexicon size. Thirdly, the parameter q also reflects the quality of input images. The better the image quality, the less likely that an edit operation will be ignored. $1 - q$ can be considered as the recognizer's ability to distinguish characters, so the dependence of word recognition on image quality is actually decomposed into two parts: first word recognition on character recognition and then character recognition on image quality.

Therefore, this model takes into account the nature of the recognizer, the complexity of lexicons and the quality of input images and aims in revealing statistically their relation to word recognition performance.

3 Experimental Resources

Recognizers: We use 5 different word recognizers in our experiments. WR1 adopts an over-segmentation methodology along with word model based recognition using dynamic programming [8]. WR2 is similar to WR1 except for its segmentation and preprocessing algorithms [9]. WR3 is grapheme based and involves no explicit segmentation [12]. It uses word model based recognition with dynamic programming. WR4 adopts an over-segmentation methodology along with character model based recognition using dynamic programming [3]. WR5 uses over-segmentation and character model based recognition with continuous density and variable duration hidden Markov models [2].

Images: All experiments are conducted on a set of 3000 US postal word images (digitized at 212 dpi) of unconstrained writing styles. The images are divided into equal halves, one for training and the other for testing.

Lexicons: To test the dependence of performance on lexicon size, we generate lexicons of size 5,10,20 and 40 for each image. For each lexicon size, 10 lexicons are generated and ordered in ascending order of average edit distance. Therefore, there are 40 lexicons in total associated with each image. In order to allow wide variation of average edit distance, these 40 lexicons actually contain meaningless entries that are random combinations of characters.

4 Deciding Model Parameters

We run each word recognizer on the training set, which contains 1,500 images and 40 lexicons for each image. Then, the 60,000 image-lexicon instances are partitioned by lexicon size and the average edit distance between lexicon entries, resulting in 44 partitions with those of size under 100 ignored. Table 1 gives the performance data of the recognizers on all partitions. It should be noticed that D_i is actually the average of average edit distance for all lexicons in a partition and the number of instances will be used as weights in multiple regression.

We apply the regression model (Equation 5) to Table 1 and obtain model parameters in Table 2. t -ratio, defined as the ratio of a value to its standard error, measures the significance of a parameter. The larger the absolute value of t -ratio, the less likely that a parameter has an actual value of zero and the less likely it is redundant. Therefore, Table 2 shows high significance of model parameters for all five word recognizers. R^2 , namely the *Coefficient of Multiple Determination*, indicates the proportion of variation in the data which is explained by the regression model and thus measures the overall significance of the regression model. Table 2 shows that about 98% of data variation for WR1-4 and about 96% for WR5 have been “explained” by the regression model.

Lex size n_i	Average edit dist D_i	num. inst.	Performance				
			WR1	WR2	WR3	WR4	WR5
5	1.1814	2169	3400	8368	6611	7865	8220
5	1.6136	3869	8736	8555	7085	7914	8431
5	2.1644	2368	9155	8830	7390	8269	8796
5	2.6169	1580	9310	9171	7911	8475	8994
5	3.1167	1294	9498	9359	8176	8655	9057
5	3.6191	1100	9645	9464	8518	8973	9291
5	4.1174	850	9718	9400	8706	8859	9471
5	4.6125	629	9825	9650	9269	9173	9396
5	5.1145	465	9871	9656	9054	9484	9441
5	5.6038	277	9856	9856	9422	9386	9567
5	6.1213	171	9942	9883	9532	9474	9591
10	1.3344	2844	7546	7254	5030	6600	7416
10	1.6960	3824	7790	7577	5221	6556	7892
10	2.1958	2301	8483	8270	5958	7149	8192
10	2.7178	1410	8723	8624	6780	7681	8482
10	3.2030	1423	9213	8904	7182	7920	8566
10	3.7166	960	9302	8958	7385	8073	8906
10	4.2110	967	9607	9400	8004	8449	9060
10	4.7110	554	9603	9314	7960	8755	8971
10	5.1923	371	9730	9569	8329	8841	9030
10	5.7205	196	9643	9796	8827	9082	9031
10	6.1845	168	9881	9524	8631	8929	9226
20	1.3870	2402	6428	6220	3739	5462	6873
20	1.6985	4202	6671	6323	3665	5305	7175
20	2.2062	2064	7413	7093	4409	5891	7398
20	2.7374	1625	7920	7674	5114	6437	7692
20	3.2136	1311	8345	8168	5919	6819	8108
20	3.7390	1023	8719	8399	6197	7067	8289
20	4.2147	944	9121	8761	6716	7288	8305
20	4.7181	551	9292	8875	7132	8004	8475
20	5.2251	389	9563	9383	7609	8123	8638
20	5.7209	202	9505	9307	8119	8515	8762
20	6.2043	127	9685	9528	7638	8268	8031
40	1.4187	1881	5253	5114	2621	4460	6321
40	1.6956	4735	5333	4813	2161	3880	6583
40	2.2127	1978	6360	6087	2968	4560	6790
40	2.7469	1707	6942	6731	3855	5231	7077
40	3.2290	1312	7589	7081	3925	5716	7188
40	3.7507	1015	7970	7330	4581	5882	7507
40	4.2350	1044	8625	7998	5817	6922	7799
40	4.7133	530	8774	8509	6151	7019	7849
40	5.2332	375	8907	8613	5653	7333	8000
40	5.7281	190	9105	8789	6632	7632	7684
40	6.2192	120	9167	8583	6833	7167	8000

Table 1. Performance data on the training set

5 Classifier Combination

Previously, logistic regression [7] proves to be successful in providing fixed weights for recognizers to be combined. In dealing with the situation where a recognizer performs well on some inputs but not well on some others, the inputs are divided into partitions according to the *state of agreement*, i.e. how the top choices returned by recognizers agree with each other, and logistic regression is performed on each partition. However, situations also exist when the relative performance of recognizers changes and fixed weights are not sufficient. For example, Figure 1 gives two performance curves when lexicon size is 40 for WR1 and WR5. When the average distance is more than 4.2, WR1 is the best; otherwise, WR5 is the best. Under this situation, partitioning inputs by states of agreement will not help because it does not separate the cases when WR1 is better and the others when WR5 is better.

After parameters have been decided, the performance model allows performance prediction given lexicons. The predictions can be used as weights in combining multiple recognizers. Moreover, since there is strong dependence of performance on lexicons, partitioning inputs by lexicon size and the average edit distance between lexicon entries instead of states of agreement can be another solution.

In weighted recognizer combination, the combined score

	ln q			α		
	value	std err	t-ratio	value	std err	t-ratio
WR1	-.5773	.0147	-32.75	2.0500	.0562	36.51
WR2	-.4722	.0157	-30.04	1.8964	.0602	31.50
WR3	-.3636	.0117	-31.04	1.8044	.0449	40.23
WR4	-.3089	.0106	-28.94	1.7162	.0409	41.99
WR5	-.2688	.0146	-18.39	1.3716	.0560	24.51

	ln k			R^2
	value	std err	t-ratio	
WR1	-2.1620	.0660	-39.36	.9850
WR2	-2.0797	.0708	-29.38	.9785
WR3	-1.3683	.0527	-25.94	.9841
WR4	-1.8395	.0481	-38.27	.9842
WR5	-2.1209	.0658	-32.24	.9575

Table 2. Regression model parameters.

m	Lexicon size 20			Lexicon size 40		
	1	2	10	1	2	10
LR	.7514	.8009	.8318	.6437	.7100	.7676
LR-PL	.7710	.8057	.8402	.6889	.7322	.7777
PP	.7737	.8095	.8403	.7001	.7358	.7801

Table 3. Combining WR1 and WR5.

of class c is defined as $\sum_R w_R r_R(c)$ where w_R is the weight of recognizer R and $r_R(c)$ is the rank score of class c given by recognizer R . We use Borda count as rank score, i.e. $m + 1 - n$ for a class ranked as the n -th choice among total m choices participating in the combination. For logistic regression(LR), weights are decided by training data and remain fixed for all testing data. For logistic regression with partitions of lexicons(LR-PL), weights are decided for each partition and remain fixed for that partition. However, a performance prediction(PP) method calculates weights for every input, thus is more dynamic. Table 3 gives the results of combining WR1 and WR5 on the testing set using these three methods. Generally, PP is better than LR- PL and LR-PL better than LR. However, it should be noticed that performance prediction does not help when the relative performance of recognizers remains unchanged across all lexicons, such as in combining WR1 and WR5 for lexicon size 5 and 10 or in combining WR1-4 for all lexicon sizes.

Performance prediction can be also applied to dynamic classifier selection where only one recognizer with the highest predicted performance gets running for an input. Figure 1 gives the result of selecting between WR1 and WR5 when lexicon size is 40. It can be seen that dynamic classifier selection using predicted performance results in higher performance than any of WR1 and WR5 individually.

6 Conclusions

We present a performance prediction model for handwritten word recognizers. This model considers factors involved in word recognition, including the nature of a recognizer, the complexity of lexicons and the quality of input images. Multiple regression is used to decide model parameters and statistical analysis has shown the significance of

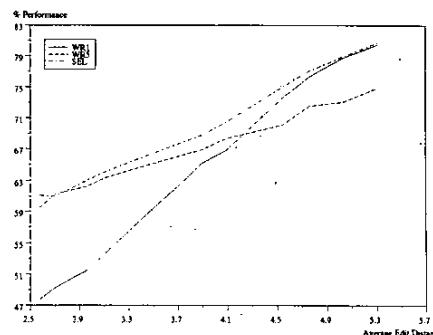


Figure 1. Dynamic classifier selection.

this model in revealing the dependence of word recognition performance on these factors. We also apply this model to combining multiple recognizers for the cases when the relative performance of recognizers changes with lexicons.

References

- [1] H. S. Baird. State of the art of document image degradation modeling. In *IAPR Workshop on DAS*, 2000.
- [2] M. Chen, A. Kundu, and S. Srihari. Variable duration hidden Markov model and morphological segmentation for handwritten word recognition. *IEEE Trans. on Image Processing*, 4(12):1675–1688, 1995.
- [3] J. Favata. Character model word recognition. In *Proc. of 5th IWFHR*, pages 437–440, 1996.
- [4] V. Govindaraju, P. Slavik, and H. Xue. Use of lexicon density in evaluating word recognizers. *IEEE Trans. on PAMI*, To appear.
- [5] T. K. Ho and H. S. Baird. Large-scale simulation studies in image pattern recognition. *IEEE Trans. on PAMI*, 19(10):1067–1079, 1997.
- [6] T. K. Ho, J. J. Hull, and S. N. Srihari. On multiple classifier systems for pattern recognition. *Proc. of 11th ICPR*, II:84–87, 1992.
- [7] T. K. Ho, J. J. Hull, and S. N. Srihari. Decision combination in multiple classifier systems. *IEEE Trans. on PAMI*, 16(1):66–75, 1994.
- [8] G. Kim and V. Govindaraju. A lexicon driven approach to handwritten word recognition for real-time applications. *IEEE Trans. on PAMI*, 19(4):366–379, 1997.
- [9] P. Slavik and V. Govindaraju. An overview of run-length encoding of handwritten word images. Technical Report 09, SUNY at Buffalo, 2000.
- [10] P. Slavik and V. Govindaraju. Use of lexicon density in evaluating word recognizers. In *Multiple Classifier Systems*, pages 310–319, 2000.
- [11] L. Xu, A. Krzyzak, and C. Y. Suen. Methods for combining multiple classifiers and their applications to handwriting recognition. *IEEE trans. on SMC*, 23(3):418–435, 1992.
- [12] H. Xue and V. Govindaraju. Building skeletal graphs for structural feature extraction on handwriting images. In *Proc. of 6th ICDAR*, pages 96–100, 2001.