

Towards style-based dating of historical documents

Sheng He*, Petros Samara†, Jan Burgers‡, Lambert Schomaker*

* University of Groningen, the Netherlands

Email: s.he@rug.nl, L.Schomaker@ai.rug.nl

†Instituut voor Nederlandse geschiedenis, Den Haag, the Netherlands

Email:petros.samara@huygens.knaw.nl

‡University of Amsterdam, the Netherlands

Email:J.W.J.Burgers@uva.nl

Abstract—Estimating the date of undated medieval manuscripts by evaluating the script they contain, using document image analysis, is helpful for scholars of various disciplines studying the Middle Ages. However, there are, as yet, no systems to automatically and effectively infer the age of historical scripts using machine learning methods. To build a system to date medieval documents is a challenging problem in several aspects: 1) As yet, no suitable reference dataset of medieval handwriting exists; 2) relatively little is known about the evolution of writing styles in the Middle Ages, and especially in the later Middle Ages. Our Medieval Paleographic Scale (MPS) project aims at solving these problems. We have collected a corpus of charters from the Medieval Dutch language area, dating from the period 1300 to 1550. A global and local regression method is proposed for learning and estimating the year in which these documents were written, using several features which have been successfully used in writer identification. The proposed system can serve as a basic tool for the medievalist or paleographer. The experimental results of the proposed method demonstrate its effectiveness.

Keywords—Medieval Paleographic Scale, historical document dating, age estimation, global and local regression

I. INTRODUCTION

The task of dating Medieval manuscripts is of the utmost importance to scholars of various disciplines studying the Middle Ages. Manuscripts that don't carry a date make it hard to assess their reliability as a historical source. However, this task is often regarded as the prerogative of a mere handful of specialists capable of correctly evaluating certain handwriting characteristics, but nevertheless sometimes conflicting conclusions. Usually, the dating of an instance of medieval script is based on the individual non-verbal intuition of the expert rather than on objective criteria. This state of affairs is not surprising, because there is a notorious lack of suitable collections of dated manuscripts that can be used as reference corpus. As the archaeologist has the ^{14}C technique to date organic materials, so the medievalist needs a method of dating manuscripts. The reliability of ^{14}C method is limited, however, when applied to medieval documents or manuscripts, and is, moreover, destructive because it requires physical samples.

In other domains, algorithms for determining the 'age' on the basis of image have been proposed [1], [2] from an human faces, while [3] addresses the age of color photographs. Both of those methods are based on extracting features from images and use regression [1], classification [3], or other learning methods [2]. However, there exists virtually no literature about dating historical documents or books. The emergence of



Fig. 1: Images of charters in our dataset from different cities.

historical handwriting collections, such as the Monk system [4] which contains more than 50k page images, has enabled progress on novel historical document dating methodology using pattern recognition and machine learning.

Given an undated manuscript, one possible way to estimate its year of origin is to search for similar writing styles in a large reference database consisting of dated documents, or to extract the general trend of writing styles in a certain period from the same database. A dating system such as this should, in other words, contain several steps: 1), a reference database which contains Medieval manuscripts or documents with year label is assembled; 2), several features are used to measure the similarity of writing styles in those documents. 3), machine learning methods are applied to perform the fine-tuned estimation of the year of origin of a given undated piece of writing.

From a paleographic perspective, it is not sensible to aim at accuracies higher than about a decade. Therefore, a collection of documents, and more specifically: charters, was composed dating from eleven so called "key periods" from the era 1300 to 1550, i.e. dating from the year marking each quarter century within this era and the decade surrounding



Fig. 2: An illustration of the development of the character ‘p’ from the ages 1300 to ages 1525. (Note: top-left is from ages 1300, bottom right is from ages 1525, in reading order.)

TABLE I: Current collection of reference documents over the key year.

Key year	1300	1325	1350	1375	1400	1425	1450	1475	1500	1525	1550
number of charters	95	141	89	168	196	198	245	205	141	115	113

it: 1300 ± 5 , 1325 ± 5 , 1350 ± 5 , ..., 1550 ± 5 . This design of the dataset allows both for regression and classification-based methods. Based on the collected data, we design a combined global and local regression method, for learning and estimating the age of late Medieval documents.

We organize the remainder of the paper as follows. In section 2, we will discuss the construction of the data used for the MPS proposed here. In section 3, we will present our approach for data-driven modeling of writing style processes. In section 4, we will then use a number of features and several machine learning methods to perform the dating. Section 5 presents the experimental results. Finally, in section 6, we will give the conclusion of the proposed method.

II. DATA

As discussed in Section 1, we aim to build a system to date medieval documents and manuscripts. Therefore, we construct an MPS dataset which contains images of charters from four cities representing the four corners of the Medieval Dutch language area: Arnhem, Leuven, Leiden and Groningen. Each charter is dated and was written within five years before or after one of the quarter century years, from the period under consideration here (1300-1550), such as 1300, 1325, 1350, ..., 1550. We designate each of these quarter century years as a ‘key year’.

A charter is a formal, public declaration of a certain legal or financial transaction or action, often functioning as proof of that transaction having taken place. It consists of one single piece of parchment (in almost all cases) or paper (very rarely). Usually, with very few exceptions, the chartered declaration is written on one side of the parchment only.

The most important reason charters are used in the MPS is that they are the only type of dated sources extant from the later middle ages available in large enough numbers to

assemble a statistically viable corpus of different writing hands without any large chronological gaps. Furthermore there is a certain consistency in the level of execution and the script type used in charters, so they are to a certain extent comparable. There are also script types which are very inconsistent with the script usually used in charters, notably a book script type called “textualis” [5]. The idea is to first establish a Paleographic scale using charters and charter script, usually written in a script type called “cursiva” [5], and then extrapolate and compare the results to and with the development of various book scripts, later. The overwhelming majority of the charters in MPS are instances of formal ‘cursive’ script according to paleographic criteria.

There are 1706 images of charters from four cities divided into 11 key years. The number of samples in each key year is shown in Table I. Ideally, the collections should be completely balanced. However, there is a natural and inherent imbalance resulting from the arbitrary way history bequeaths its records to us, so initially we use all the available data.

III. MODELING HANDWRITTEN STYLE PROCESSES

Having constructed the MPS database of charters, we now propose a data-driven technique for modeling the writing styles in the period of 1300-1550. The collected database and modeling are not only relevant to the problem of dating, but also help to develop character and word models for classification and retrieval [4], [6].

Our proposal is to apply a technique to model the writing styles in each key years and use it to estimate the year of origin of undated documents. Currently, there are not many studies in the area of style-based document dating. The ‘Graphem’ project [7] in France uses a combination of traditional paleography and pattern analysis. Similarly, allograph-based analysis is proposed in the DigiPal project [8]. However, these approaches assume a manual segmentation

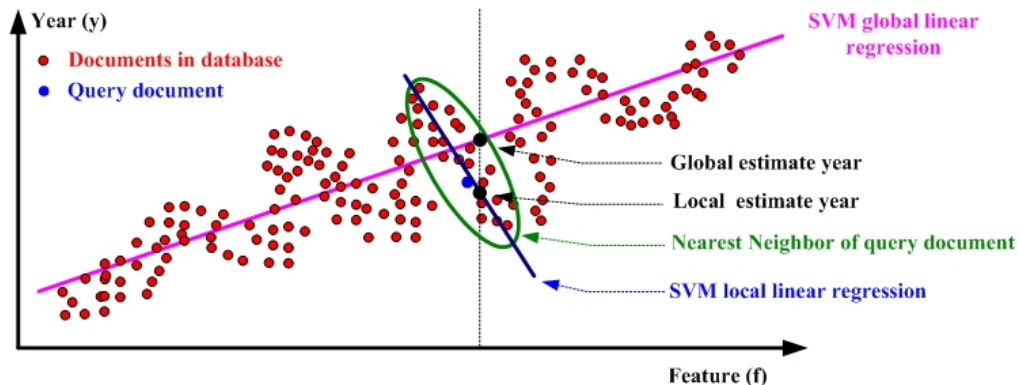


Fig. 3: Graphical depiction of the rationale for global and local regression in a hypothetical style evolution.

of characters. In [7], preliminary and inconclusive results on document dating were shown. The current study is the first attempt at digital dating using a systematically constructed scale of reference documents from medieval Dutch/Flemish manuscript collections. The assumption is that style variation can be traced over time using the evolution of ‘textural’ features of a fairly large region of interest in the document image, similar to [9], as opposed to using costly and detailed manual character segmentation.

The underlying assumption for our proposed model is that writing styles changed gradually, continuously and in general within a relatively limited time frame (within 25 years) in the period under consideration 1300-1550. The rationale behind the assumption of a gradual style evolution comes from the observation that scribes were strictly and formally trained by experienced, older teachers. As an example, Fig. 2 shows the writing styles of character ‘p’, as it was written in different ways in the period from 1300 to 1525. If one wants to avoid the individual character segmentation and recognition, the question is whether the style evolution in the individual allograph is also reflected in overall page texture features, such as the Hinge or Fraglets [10]. In this paper, we try to give an answer of this question. Preliminary experiments with regressive and discriminative methods indicated that an advanced approach was necessary.

Fig. 3 shows a graphical depiction about a hypothetical evolution of writing styles in the period of 1300-1550 for some feature ‘f’. The red dots represent charters in the MPS, ordered in a 2D plane of their year (y) and the hypothetical-feature value (f). As discussed in [1], the global regression (the pink line) will yield a suboptimal performance because it attempts to find a straight curve to approximate the data. Therefore, local adjustments are necessary to obtain a higher performance. The use of all samples in a key year for local regression, as in [1], is risky, because multiple styles may be active during a key year. Therefore, we propose to use k-nearest neighbor (kNN), to first determine a subset of relevant samples as the basis for the local regression. Evidently, the value of k must be large enough to allow the determination of a good support-set. In this manner it is possible to zoom in on the relevant (sub)style, with limited computational cost (typically $k = 100$). There are two advantages to use kNN: 1) it can find the documents in MPS with similar writing styles; 2) it is efficient to train a

local linear regression (blue line) based on the subset of k neighbors.

IV. A THREE-STAGE GLOBAL/LOCAL ESTIMATION.

A. features

The **Hinge feature** is a texture-level feature, which captures the curvature of the ink trace of the character. It is very discriminatory between different writers and also a powerful feature to measure writing styles in historical documents [11]. The center idea is to consider two contour fragments attached at a common end pixel and compute the joint probability distribution of the orientations of the two legs of the obtained contour-hinge. The hinge feature is agglomerated in an interpolated $p \times p$ histogram, where p is the number of angle bins. For more details refer to [10]. In this paper, we set the $p = 40$, and remove the redundant information, thus the dimension of Hinge feature is 696.

The **Fraglets feature** is an allograph-level approach to writer identification. The underlying assumption of the Fraglets is that the writer acts as a stochastic generator of ink-blob shapes, or graphemes in [10]. The Fraglets feature, which is a probability distribution of grapheme usage, is characteristic of each document and computed using a common codebook of shapes obtained by Kohonen SOM or k-means methods. In our model, we collect all the character contours in MPS, and train the codebook with a size of 70×70 , the dimension of Fraglets feature being 4900.

For each document, the Hinge and Fraglets features are extracted from the full text portion of the whole document, and concatenated together to form the feature representation. Finally, the dimensionality of feature is 5596, and the scale is in probabilities. Although this dating problem may also indicate a solution in the direction of principal component analysis or a similar dimensionality reduction, we decided on the basis of pilot experiments, to exploit the stability provided by the key years as solid anchors along the time line, by using a mixed nearest-neighbor and regression scheme.

B. regression

In this section, we use statistical methods to construct a model according to our assumption in section 3. The proposed

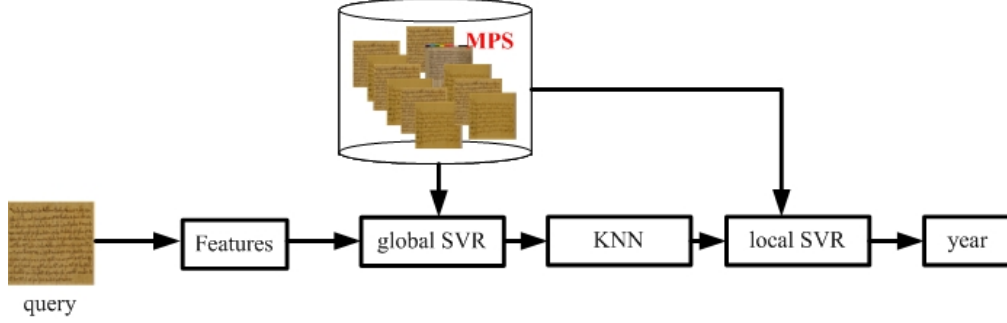


Fig. 4: The framework of proposed method for historical document dating.

algorithm is divided into three steps. 1) Global regression; 2) Local support-set selection; 3) Local regression.

Global regression: The global regression can capture the global trend of the writing styles in our database (see the pink line in Fig. 3). For the purpose of robust regression of the development of the writing style in historical documents in our database, we adopt the support vector regression (SVR) method [12], which is widely used for the problem of estimating the age of human faces [1] and historical color images [3].

As discussed in [1] the global regression method cannot show an optimal performance for age prediction from face images and this is also expected to be the case for historical documents. In fig. 3 it is illustrated that the global regression method may estimate the global year trend linearly, but cannot estimate the year precisely. In order to solve this problem, a locally adjusted robust regressor (LARR) method was proposed in [1], which used a linear Support Vector Machine (SVM) on a local search region around the global estimation age. However, there is an unsolved problem in [1] that concerns determining the local search region. In this paper, we apply the kNN method to select the local samples, and train a local SVR to estimate the year of documents.

Local support-set selection: Given the feature representation of a query document x , and its estimated year y from global SVR, and the Mean Absolute Error (MAE), denoted as MAE_{global} which is defined as the average of the absolute errors between the estimated years and the ground truth year [1], [13] (see Eq(1)) of global SVR on the all documents in our database, we set the search region λ times of MAE_{global} around the estimated year y . Mathematically, the search region is defined as $\mathcal{R} = [y - \lambda MAE_{global}, y + \lambda MAE_{global}]$. However, unlike LARR [1], not all the samples in the search region \mathcal{R} are used to locally adjust the estimated year. In our paper, we follow the work [14] that applies k-nearest neighbor (kNN) method to find the k nearest neighbors of the query document in \mathcal{R} . In our experiments, we set $k = 100$ and $\lambda = 2$.

Local regression: Given the local support document set from k nearest neighbors in the search region \mathcal{R} , another support vector regression (SVR) which we call “local SVR” is trained using the same method as global regression. The local SVR will fine tune the initially estimated age value according to the global SVR towards the true age as close as possible (see Fig. 3) [1].

V. EXPERIMENTS

A. Experiment Setup

In our experiments, we train a linear SVR and non-linear SVR method both for global and local regression. The ϵ -insensitive loss function is adopted in SVR, in which we disregard errors as long as they are less than ϵ , and we set $\epsilon = 0.1$ in our experiments. To compare with linear SVR, a nonlinear regression function is used. For our document dating, the Gaussian radial basis function kernel is adopted. A radial basis function (RBF) is $k(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$, where γ is a constant to adjust the width of the Gaussian function.

The hyperparameter in linear SVR is a cost parameter (C), and in non-linear SVR an additional gamma parameter (γ). Regular grid search method is commonly used to find the suitable values of parameters. In our experiments, we try C and γ values in logarithmic scale $\{2^{-15}, 2^{-14}, \dots, 2^{14}, 2^{15}\}$.

In the experiments, we randomly split the dataset into two parts: 90% for training and 10% for testing, and 10-fold cross-validation is used. Finally, the average results are reported in this paper. Care was taken in sample selection to exclude the risk of inappropriate good results due to writer identification. In the application domain, the paleographic users want to be fairly certain that the date estimate is based on general style characteristics not an individual writing style.

B. Random Guess

In order to validate the performance of the proposed method, we also conduct a random guess for purposes of comparison. Without any prior information, we assume that documents pertain to each of the key years with equal probability. We repeat the random guess process 100 times to obtain a more stable estimation.

C. Measurements

The performance of year estimation can be measured by two different measures: the mean absolute error (MAE) and the cumulative score (CS) as is the case with age estimation using face images [1], [13]. The performance measurement of MAE is typically defined as:

$$MAE = \sum_{j=1}^N |\widetilde{K}(y_j) - K(y_j)|/N \quad (1)$$

TABLE II: Mean Absolute Error(MAE) and Standard Deviation(SD)(in years) of different methods

Methods	linear SVR	non-linear SVR	random guess
MAE	36.6	35.4	85.3
SD	34.7	32.9	58.8

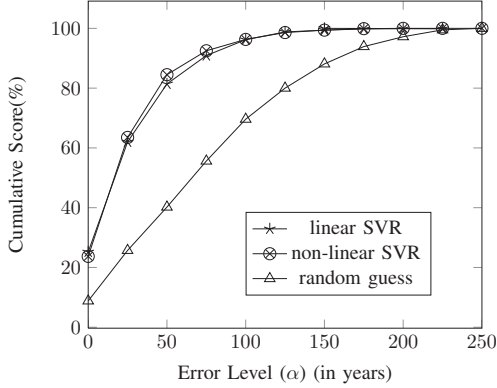


Fig. 5: Cumulative probability of $p(MAE \leq \alpha)$ on the error levels (α) from 0 to 250 years of linear SVR, non-linear SVR and random guess.

where $K(y_j)$ is the ground truth key year of the input document y_j , and $\hat{K}(y_j)$ is the estimated key year, N is the number of test documents.

The cumulative score (CS) [13] is defined as:

$$CS(\alpha) = N_{e \leq \alpha} / N \times 100\% \quad (2)$$

where $N_{e \leq \alpha}$ is the number of test images on which the year estimation makes an absolute error no higher than α years.

We also give a confusion matrix to describe the performance of each year. This matrix evaluates the probability ($\times 100\%$) which belong to year y_1 are estimated to belong to year y_2 . This matrix can show the transition probability between the ground-truth and estimated year. If the values on the diagonal are high, the performance of the system is better.

D. Dating results

Table II gives the MAEs and SDs of linear SVR, non-linear SVR and random guess. One can see that non-linear SVR works a little better than linear SVR. The errors of both linear and non-linear SVR are much lower than random guess. The best results in terms of MAE is 35.4 years using non-linear SVR.

The CS measures are shown in Fig. 5. One can observe that the score of non-linear SVR is slightly higher than linear SVR, especially in error levels less than 75 years. Both linear SVR and non-linear SVR improve the score significantly over the random guess.

Table III and IV show the performance matrices of linear SVR and non-linear SVR, respectively. From the two tables

TABLE III: confusion matrix of linear SVR

numbers	Estimated Year										
	1300	1325	1350	1375	1400	1425	1450	1475	1500	1525	1550
1300	14.7	34.7	22.1	16.8	8.4	2.1	1.1	0.0	0.0	0.0	0.0
1325	21.3	45.4	22.0	7.1	2.8	0.7	0.0	0.7	0.0	0.0	0.0
1350	1.1	12.4	22.5	22.5	30.3	10.1	1.1	0.0	0.0	0.0	0.0
1375	0.0	7.7	21.4	31.0	26.2	10.1	3.0	0.6	0.0	0.0	0.0
1400	2.0	1.5	3.6	20.9	40.8	18.9	7.7	2.0	1.0	1.0	0.5
1425	0.5	1.5	2.0	7.1	33.8	32.8	16.2	4.0	1.0	0.5	0.5
1450	0.0	0.4	2.4	3.3	18.4	35.9	31.8	6.9	0.4	0.0	0.4
1475	0.0	0.0	1.5	3.9	10.7	22.9	36.6	20.5	3.9	0.0	0.0
1500	0.0	0.7	0.0	7.8	10.6	23.4	34.0	19.1	3.5	0.7	0.0
1525	0.9	0.0	1.7	0.9	7.8	11.3	27.0	25.2	17.4	7.0	0.9
1550	0.9	0.0	0.9	4.4	5.3	11.5	23.9	19.5	23.9	8.0	1.8

TABLE IV: confusion matrix of non-linear SVR

numbers	Estimated Year										
	1300	1325	1350	1375	1400	1425	1450	1475	1500	1525	1550
1300	17.9	53.7	13.7	10.5	3.2	0.0	1.1	0.0	0.0	0.0	0.0
1325	17.0	50.4	23.4	7.1	1.4	0.0	0.0	0.7	0.0	0.0	0.0
1350	2.2	21.3	16.9	29.2	20.2	5.6	3.4	0.0	0.0	1.1	0.0
1375	0.0	13.1	23.8	26.2	20.8	14.3	1.8	0.0	0.0	0.0	0.0
1400	1.0	1.0	7.7	23.5	35.2	18.4	7.7	4.1	0.5	0.5	0.5
1425	0.0	0.0	2.5	13.1	28.8	29.3	17.7	5.1	1.0	2.5	0.0
1450	0.0	0.8	1.6	7.3	17.6	33.1	28.2	9.8	1.6	0.0	0.0
1475	0.0	0.5	2.9	2.9	7.8	25.9	30.7	21.0	8.3	0.0	0.0
1500	0.0	1.4	0.7	8.5	7.8	14.9	34.0	24.8	7.8	0.0	0.0
1525	0.0	0.9	0.0	2.6	7.8	10.4	23.5	25.2	22.6	6.1	0.9
1550	0.0	0.0	0.9	4.4	5.3	8.0	15.9	15.0	23.9	26.5	0.0

TABLE V: Mean Absolute Error (in years) in three periods

	1300-1375	1400-1475	1500-1550
linear SVR	27.9	27.1	70.1
non-linear SVR	25.8	28.8	63.3
random guess	93.8	73.5	103.0
Number of documents in MPS	493	844	369

we can see that the proposed method has a large bias to its nearby key years.

The documents in MPS can be divided into three periods: 1300-1375, 1400-1475 and 1500-1550. Table V shows the MAEs of the three periods. The MAEs of periods 1300-1375 and 1400-1475 are almost same. This trend is consistent with the number of documents in the three periods of the MPS.

In order to evaluate the geographical differences between writings produced at the same dates, we conduct an experiment using documents from three cities to train and the documents from the fourth city to test. Table VI and VII show the performances of linear SVR and non-linear SVR method, respectively. However, there is no paleographic assumption that the development of scripts in one city is identical to its development in other cities.

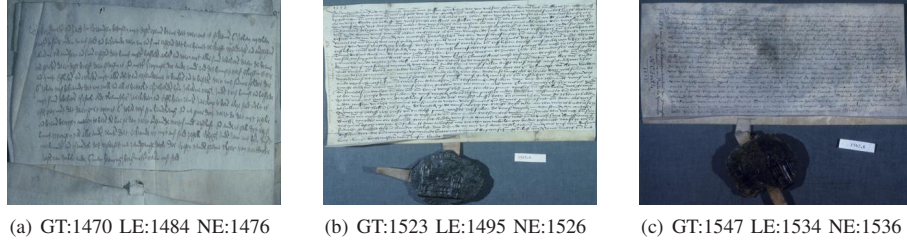


Fig. 6: Several examples of document with ground truth (GT), linear SVR estimated year (LE) and non-linear SVR estimated year (NE). (Note: the images are resized to same size for visualization.)

TABLE VI: Mean Absolute Error(MAE) and Standard Deviation(SD)(in years) of geographical testing using **linear-SVR**

city	Arnhem	Leiden	Leuven	Groningen
MAE	51.0	38.4	45.5	64.1
SD	36.6	32.4	29.7	51.8

TABLE VII: Mean Absolute Error(MAE) and Standard Deviation(SD)(in years) of geographical testing using **non-linear-SVR**

city	Arnhem	Leiden	Leuven	Groningen
MAE	51.3	40.8	44.3	64.9
SD	37.8	33.9	32.0	51.7

VI. CONCLUSION

In this paper, we present the preliminary, earlier results of the Medieval Paleographic Scale (MPS) project. This project aims at computer-based dating of charters collected from the period between 1300 and 1550 from four cities of the Medieval Dutch language area. Based on this historical documentary data, we introduced the novel task of automatically estimating the age of undated Medieval documents. We showed that given successful features, both linear SVR and non-linear SVR achieve promising results significantly better than randomly guessing.

However, there are still challenges for dating Medieval documents in several aspects: (1) More documents should be added into MPS, especially in the period of 1300-1400. From our experiments, we can observe that the larger the number of documents given, the better the performance obtained (in the period of 1400-1475). (2) Although our experiments have shown an improving performance, the MAE of the proposed method (35.4) is still higher than 25 years which is the interval between key years. (3) Due to the poor image quality of historical documents, several pre-processing methods are needed, such as image de-noising, image enhancement, and robust binarization. (4) Apart from Hinge and Fraglets, other features will need to be developed, tuned to ? and the task of dating Medieval documents, such as the Quill and Quill-Hinge features [11] and $\Delta^n Hinge$ features [15].

ACKNOWLEDGMENT

This research was granted by the Dutch Organization for Scientific Research NWO (project No. 380-50-006) in 2012.

REFERENCES

- [1] G. Guo, Y. Fu, C. R. Dyer, and T. S. Huang, "Image-based human age estimation by manifold learning and locally adjusted robust regression," *Image Processing, IEEE Transactions on*, vol. 17, no. 7, pp. 1178–1188, 2008.
- [2] K. Chen, S. Gong, T. Xiang, Q. Mary, and C. C. Loy, "Cumulative attribute space for age and crowd density estimation," *Computer Vision and Pattern Recognition*, 2013.
- [3] F. Palermo, J. Hays, and A. A. Efros, "Dating historical color images," in *Computer Vision–ECCV 2012*. Springer, 2012, pp. 499–512.
- [4] T. Van der Zant, L. Schomaker, and K. Haak, "Handwritten-word spotting using biologically inspired features," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 30, no. 11, pp. 1945–1957, 2008.
- [5] J. P. Gumbert, "A proposal for a cartesian nomenclature," *Essays presented to G. I. Liefinck, IV: miniatures, scripts, collections. (Litterae Textuales)*, Ed., J.P. Gumbert and M.J.M. de Haan., pp. 45–52, Amsterdam 1976.
- [6] J.-P. Van Oosten and L. Schomaker, "Separability versus prototypicality in handwritten word-image retrieval," *Pattern Recognition*, vol. 47, no. 3, pp. 1031–1038, 2014.
- [7] "<http://liris.cnrs.fr/graphem/>."
- [8] P. Stokes, "Modeling medieval handwriting: A new approach to digital palaeography," in *DH2012 Book of Abstracts*, ed. by Jan Christoph Meister et al., pp. 382–385, 2012.
- [9] L. Schomaker, K. Franke, and M. Bulacu, "Using codebooks of fragmented connected-component contours in forensic and historic writer identification," *Pattern Recognition Letters*, vol. 28, no. 6, pp. 719–727, 2007.
- [10] M. Bulacu and L. Schomaker, "Text-independent writer identification and verification using textural and allographic features," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 29, no. 4, pp. 701–717, 2007.
- [11] A. Brink, J. Smit, M. Bulacu, and L. Schomaker, "Writer identification using directional ink-trace width measurements," *Pattern Recognition*, vol. 45, no. 1, pp. 162–171, 2012.
- [12] V. N. Vapnik, "Statistical learning theory," 1998.
- [13] X. Geng, Z.-H. Zhou, and K. Smith-Miles, "Automatic age estimation based on facial aging patterns," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 29, no. 12, pp. 2234–2240, 2007.
- [14] H. Zhang, A. C. Berg, M. Maire, and J. Malik, "Svm-knn: Discriminative nearest neighbor classification for visual category recognition," in *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, vol. 2. IEEE, 2006, pp. 2126–2136.
- [15] S. He and L. Schomaker, "Delta-n hinge: rotation-invariant features for writer identification," in *International Conference on Pattern Recognition*, 2014.