

## Learning to rank based on subsequences

Basura Fernando\*

ACRV, The Australian National University

Damien Muselet

Jean Monnet University, St-Etienne, France

Efstratios Gavves\*

QUVA Lab, University of Amsterdam, Netherlands

Tinne Tuytelaars

KU Leuven, ESAT-PSI, iMinds, Belgium

### Abstract

We present a supervised learning to rank algorithm that effectively orders images by exploiting the structure in image sequences. Most often in the supervised learning to rank literature, ranking is approached either by analyzing pairs of images or by optimizing a list-wise surrogate loss function on full sequences. In this work we propose MidRank, which learns from moderately sized sub-sequences instead. These sub-sequences contain useful structural ranking information that leads to better learnability during training and better generalization during testing. By exploiting sub-sequences, the proposed MidRank improves ranking accuracy considerably on an extensive array of image ranking applications and datasets.

### 1. Introduction

The objective of supervised learning-to-rank is to learn from training sequences a method that correctly orders unknown test sequences. This topic has been widely studied over the last years. Some applications include video analysis [28], person re-identification [31], zero-shot recognition [21], active learning [17], dimensionality reduction [8], 3D feature analysis [30], binary code learning [18], learning from privileged information [26] and interestingness prediction [10]. In particular, we focus on image re-ranking. Image re-ranking is useful in modern image search tools to facilitate user specific interests by re-ordering images based on some user specific criteria. In this context, we re-order the top-k retrieved images from an image search engine based on some criteria such as interestingness. In this paper, we propose a new, efficient and accurate supervised learning-to-rank method that uses the information in image subsequences effectively for image re-ranking.

Most learning-to-rank methods rely on pair-wise cues and constraints. However, pairs of ranked images provide rather weak, and often ambiguous, constraints, as illus-

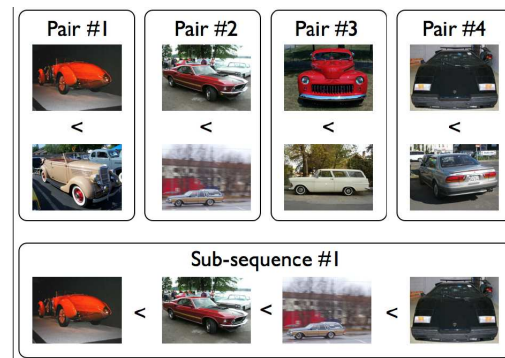


Figure 1: Imagine you want to learn a ranking function from the example pairs in the top row. What could be the ranking criterion implied by the data? Ranking more to less sporty cars could be a criterion. Ranking more to less colorful cars could be another. Given an example sequence, however, as in the bottom row, it becomes clear that the *chronological* order is the most plausible ranking criterion. We advocate the use of such (sub-) sequences, instead of just pairs of images or long sequences, for supervised learning of rankers.

trated in Fig. 1. Especially when complex and structured data is involved, considering only pairs during training may confuse the learning of rankers as shown in [6, 19, 35]. Learning-to-rank is in fact a prediction task on lists of data/images. Treatment of pairs of images as independent and identically distributed random variables during training is not ideal [6]. It is, therefore, better, to consider longer subsequences within a sequence, which contain more information than pair of elements, see Fig. 1.

To exploit the structure in long sequences, list-wise methods [6, 29, 32, 33, 34, 35] optimize for ranking losses defined over sequences. Although such an approach can exploit the structure in sequences, working on long sequences introduces a new problem. More specifically, as the number of wrong permutations grows exponentially with respect to the sequence length, list-wise methods often end up with a more difficult learning problem [19]. Hence, list-wise meth-

\*The work was conducted at KU Leuven, ESAT-PSI.

ods often lead to over-fitting, as we also observe in our experimental evaluations.

To overcome the above limitations, we propose to use *subsequences* to learn-to-rank. On one hand, the increased length of subsequences brings more information and less uncertainty than pairs. On the other hand, compared to full sequences learning on subsequences allows more regularity and better learnability.

Note that the training subsequences can be generated without any additional labeling cost: if the training set consists of sequences, we sub-sample; if the training set consists of pairs, we build subsequences by exploiting the transitivity property. We argue that every subsequence of any length sampled from a correctly ordered sequence also has the correct ordering *i.e. all subsequences of a correctly ordered sequence are also correctly ordered*. We exploit this property as follows. Given the training subsequences, we learn rankers that minimize the *zero-one loss* per subsequence length (or scale). During testing, using the above property, we evaluate all ranking functions of different lengths over the full test sequences using convolution. Then, to obtain the final ordering, we fuse the ranking results of different rankers.

Our major contributions are threefold. First, we propose a method, *MidRank*, that exploits subsequences to improve ranking both quantitatively and qualitatively. Second, we present a novel difference based vector representation that exploits the total ordering of subsequences. This representation, which we will refer to as *stacked difference vectors*, is discriminative and results in learning accurate rankers. Third, we introduce an accurate and efficient polynomial time testing algorithm for the NP-hard [23] linear ordering problem to re-rank images in moderately sized sequences.

We evaluate our method on three different applications: ranking images of famous people according to relative visual attributes [21], ranking images according to how interesting they are [10] and ranking car images according to the chronology [15]. Given an image search result obtained from an image search engine, we can use our method to re-rank images in a page to satisfy user specific criteria such as interestingness or chronology. Results show a consistent and significant accuracy improvement for an extensive palette of ranking criteria.

## 2. Related work

Supervised learning-to-rank algorithms are categorized as point-wise, pair-wise and list-wise methods. Point-wise methods [7], which process each element in the sequence individually, are easy to train but prone to over-fitting. Pair-wise methods [12, 13, 24] compute the differences between two input elements at a time and learn a binary decision function that outputs whether one element precedes the other or vice-versa. These methods are

restricted to pair-wise loss functions. Naturally, pair-wise methods do not explicitly exploit any structural information beyond what a pair of elements can yield. List-wise methods [6, 32, 34, 35, 29, 33], on the other hand formulate a loss on whole lists, thus being able to optimize more relevant ranking measures like the NDCG or the Kendall-Tau.

We present MidRank, which belongs to a fourth family of learning to rank methods, that is positioned between pair-wise and list-wise methods. Similar to pair-wise methods, MidRank uses pairwise relations but extends to more informative subsequences, by considering multiple pairs within a subsequence simultaneously. Similar to list-wise methods, MidRank optimizes a list-wise ranking loss, but unlike most list-wise methods we use zero-one sequence loss. This is done at sub-sequences learning thus allowing to exploit the *regularity* in them during learning.

In [9] Dokania *et al.* propose to optimize average precision *information retrieval loss* using point-wise and pair-wise feature representations. However, this method only focuses on information retrieval.

MidRank is also different from existing methods that use multiple weak rankers, such as LambdaMART [32] and AdaRank [34], which propose a linear combination of weak rankers, with iterative re-weighting of the training samples and rankers during training. In contrast, MidRank learns multiple ranking functions, one for each subsequence length, therefore focusing more on the latent structure inside the subsequences.

## 3. MidRank

We start from a training set of ordered image sequences. Each sequence orders the images according to a predetermined criterion, *e.g.* images ranging from the most to the least happy face or from the oldest to the most modern car. Our goal is to learn from data in a supervised manner a ranker, such that we can order *a new* list of *unseen* images according to the same criterion.

**Basic notations.** Our training set is composed of  $N$  ordered image sequences,  $D = \{\mathbf{X}^i, \mathbf{Y}^i, \ell^i\}, i = 1, \dots, N$ .  $\mathbf{X}^i$  stands for an image sequence  $[\mathbf{x}_1^i, \mathbf{x}_2^i, \dots, \mathbf{x}_{\ell^i}^i]$  containing  $\ell^i$  images, where  $\ell^i$  can vary for different sequences  $\mathbf{X}^i$ .  $\mathbf{Y}^i$  is a permutation vector  $\mathbf{Y}^i = [\pi(1), \dots, \pi(\ell^i)]$ , and represents that the correct order of the images in the sequence is  $\mathbf{x}_{\pi(1)}^i \succ \mathbf{x}_{\pi(2)}^i \succ \dots \succ \mathbf{x}_{\pi(\ell^i)}^i$ . Henceforth, whenever we speak of a sequence  $\mathbf{X}^i$ , we imply that it is *unordered*, and when we speak of an *ordered sequence*, we imply a tuple  $\{\mathbf{X}^i, \mathbf{Y}^i, \ell^i\}$ . To reduce notation clutter, whenever it is clear from the context we drop the superscript  $i$  referring to the  $i$ -th sequence.

### 3.1. Ranking sequences

Assume a new list of previously unseen images,  $\mathbf{X}' = [\mathbf{x}'_1, \dots, \mathbf{x}'_{\ell'}]$ . We define a ranking score function  $z(\mathbf{X}', \mathbf{Y}')$ , which should return the highest score for the correct order of images  $\mathbf{Y}'^*$ . Given an appropriate loss function  $\delta(\cdot, \cdot)$  our learning objective is

$$\arg \min_{\vartheta} \delta(\mathbf{Y}'^*, \widehat{\mathbf{Y}}'), \quad (1)$$

$$\widehat{\mathbf{Y}}' = \arg \max_{\mathbf{Y}'} z(\mathbf{X}', \mathbf{Y}'; \vartheta), \quad (2)$$

where  $\widehat{\mathbf{Y}}'$  is the highest scoring order for  $\mathbf{X}'$  and  $\vartheta$  are the parameters for our ranking score function.

The score function  $z(\mathbf{X}', \mathbf{Y}'; \vartheta)$  should be applicable for any length  $\ell'$  that a new sequence  $\mathbf{X}'$  might have. To this end we decompose the sequence  $\mathbf{X}'$  into a set of subsequences of a particular length,  $\lambda$ . We only consider consecutive subsequences, *i.e.* subsequences of the form  $\mathbf{X}'_j = [\mathbf{x}'_{j:j+\lambda}]$ . The following proposition holds for any  $\lambda \in [2 \dots \ell']$ :

**Proposition 1.** *A sequence of length  $\ell$  is correctly ordered if and only if all of its  $\ell - \lambda + 1$  consecutive subsequences of length  $\lambda$  are correctly ordered.*

This proposition follows easily from the transitivity property of inequalities. We have, therefore, transformed our goal from ranking an unconstrained sequence of images, to ranking images inside each of the constrained, length-specific subsequences. To get to the final ranking of the original sequence we need to combine the rankings from the subsequences. Based on proposition 1, we define the ranking inference as

$$\widehat{\mathbf{Y}}' = \arg \max_{\mathbf{Y}'} \sum_{j=1}^{\ell-\lambda+1} z(\mathbf{X}'_j, \mathbf{Y}'_j; \vartheta). \quad (3)$$

with  $\mathbf{Y}'_j = [\pi'(j) \dots \pi'(j + \lambda - 1)]$  and  $z(\cdot)$  the ranking score function for fixed-length subsequences. The simplest choice for  $z(\cdot)$  would be a linear classifier, *i.e.*  $z(\mathbf{X}', \mathbf{Y}'; \vartheta) = \vartheta^T \psi(\mathbf{X}', \mathbf{Y}')$ , where  $\psi(\mathbf{X}^i, \mathbf{Y}^i)$  is a feature function that we will revisit in the next subsection. However, since eq. (3) sums the ranking scores for all subsequences together, a non-linear feature mapping is to be preferred, otherwise the effect of different subsequences will be cancelled out. In practice, we use  $z(\mathbf{X}', \mathbf{Y}'; \vartheta) = \text{sign}(\vartheta^T \psi(\mathbf{X}', \mathbf{Y}')) \cdot |\vartheta^T \psi(\mathbf{X}', \mathbf{Y}')|^{\frac{1}{2}}$ . It is worth noting that the ranking inference of eq. (3) is a generalization of the inference in pairwise ranking methods, like RankSVM [13], where  $\lambda = 2$ . Moreover, eq. (3) is similar in spirit to convolutional neural network models [14], which decompose the input of unconstrained size to a series of overlapping operations.

### 3.2. Training $\lambda$ -subsequence rankers

Having decomposed an unconstrained ranking problem into a combination of constrained, length-specific subsequence ranking problems, we need a learning algorithm for optimizing  $\vartheta$ . Considering prop. (1) and eq. (3) we train the parameters  $\vartheta$  for a given subsequence length, as follows

$$\arg \min_{\vartheta} \frac{\mu}{2} \|\vartheta\|^2 + \sum_{i,j} \mathcal{L}_{(\lambda)_j}^i \quad (4)$$

$$\mathcal{L}_{(\lambda)_j}^i = \max\{0, 1 - \delta(\mathbf{Y}_j^i, \mathbf{Y}_j^{i*}) \vartheta^T \cdot \psi(\mathbf{X}_j^i, \mathbf{Y}_j^i)\}.$$

The loss function  $\mathcal{L}_{(\lambda)_j}^i$  measures the loss of the  $j$ -th subsequence of length  $\lambda$  of the  $i$ -th original sequence. For discriminative learning in eq. (4) we need both positive and negative subsequences. During training we sample subsequences of standardized lengths. Although we can mine positive subsequences of all possible lengths, in practice we focus on subsequences up to length 7-10. To generate negative subsequences during training, we scramble the correct order randomly sampled subsequences. For each positive subsequence of length  $\lambda$ , we can generate theoretically up to  $\lambda! - 1$  negatives. However, this would create a heavily imbalanced dataset, which might influence the generalization of the learnt rankers [1]. Furthermore, keeping all possible negative subsequences would have a severe impact on memory. Instead, we generate as many negative subsequences as our positives. For the optimization of eq. (4) we use the stochastic dual coordinate ascent (SDCA) method [25], which can handle imbalanced or very large training problems [25]. In practice, when adding more negatives we did not witness any significant differences.

The  $\delta(\cdot)$  function operates as a weight coefficient. The optimization uses indirectly the ranking disagreements to emphasize the wrong classifications proportionally to the magnitude of the  $\delta(\cdot)$  disagreement value. At the same time, the optimization is expressed as a zero-one loss based classification using hinge loss. This allows to maximize the margin between the correct and the incorrect subsequences of specific length.

As we are interested in obtaining the optimal ranking, we could implement  $\delta(\cdot)$  using any ranking metric, such as sequence accuracy, Kendall-Tau or the NDCG. From the above metrics the sequence accuracy,  $\delta(\mathbf{Y}, \mathbf{Y}^*) = 2[\mathbf{Y} = \mathbf{Y}^*] - 1$  is the strictest one, where  $[\cdot]$  is the Iverson's bracket. List-wise methods usually employ relatively relaxed ranking metrics, *e.g.* based on the NDCG or Kendall-Tau measure. This is mostly because for longer, unconstrained sequences on which they operate directly, the zero-one loss is too restrictive. In our case, the advantage is we have standardized, relatively small, and length-specific subsequences, on which the zero-one loss can be easily applied. With zero-one loss we observe better discrimination of correct subsequences from incorrect ones and, therefore, better

generalization. To this end in our implementation we use the zero-one loss, although other measures can also be easily introduced.

### 3.3. Ranking-friendly feature maps

To learn accurate rankers we need discriminative feature representations  $\psi(\mathbf{X}_j^i, \mathbf{Y}_j^i)$ . We discuss three different representations for  $\mathbf{X}_{(\lambda)}$ .

*Mean pairwise difference representation.* Herbrich *et al.* [12] eloquently showed that *the difference of vector* representation yields accurate results for learning a pairwise ranking function.

For this representation we have  $\psi(\mathbf{X}, \mathbf{Y}) = \frac{1}{|\{(i,j)|\mathbf{x}_{\pi(i)} \succ \mathbf{x}_{\pi(j)}\}|} \sum_{\forall \{(i,j)|\mathbf{x}_{\pi(i)} \succ \mathbf{x}_{\pi(j)}\}} (\mathbf{x}_{\pi(i)} - \mathbf{x}_{\pi(j)})$ . The mean pairwise difference representation is perhaps the most popular choice in the learning-to-rank literature [12, 27, 13, 21, 18, 9]. In the specific case of  $\lambda = 2$  we end up with the standard rank SVM [13] and the SVR [27] learning objectives.

*Stacked representation.* Standard pairwise max margin rankers make the simplifying assumption that a sequence is equivalent to a collection of pairwise inequalities. To exploit the structural information beyond pairs, we propose to use the stacked representation as  $\psi(\mathbf{X}, \mathbf{Y}) = [\mathbf{x}_{\pi(1)}^T, \dots, \mathbf{x}_{\pi(\ell)}^T]^T$  for subsequences.

An interesting property of stacked representations comes from the field of combinatorial geometry [3]. From combinatorial geometry we know that all permutations of a vector are vertices of a *Birkhoff polytope* [3]. As the Birkhoff polytope is provably convex, there exists at least one hyperplane that separates each vertex/permutation from all others. Hence, there exists also at least one hyperplane to separate the optimally permuted sequence  $\mathbf{X}_i^*$  from all others, *i.e.*, we have a linearly separable problem. Of course this linear separability applies only for the different permutations of a particular list of elements  $\mathbf{X}_i^i$ . Thus is not guaranteed that all correctly permuted  $\mathbf{X}_i^{i*}, \forall i$  training sequences will be linearly separable from all incorrectly permuted ones. However, this property ensures that from all possible permutations of a sequence, the correct one can always be linearly separated from the incorrect ones.

The same advantage of better separability between different orderings of the same sequence could be obtained by nonlinear kernels. Such kernels, however, are too expensive to apply on many realistic scenarios, when thousands of sequences are considered at a time. Furthermore, the design of such kernels is application dependent, thus making them less general.

*Stacked difference representation.* Inspired from [12]

and the stacked representations, we can also represent a sequence of images as  $\psi(\mathbf{X}, \mathbf{Y}) = [(\mathbf{x}_{\pi(1)} - \mathbf{x}_{\pi(2)})^T, \dots, (\mathbf{x}_{\pi(\lambda-1)} - \mathbf{x}_{\pi(\lambda)})^T]^T$ . Similar to mean pairwise difference representations, they model only the difference between neighboring elements in a rank, thus being invariant to the absolute magnitudes of the elements in  $\mathbf{X}_\ell$ . Furthermore it is easy to show that stacked difference representations maintain total order structure (proof in supplementary material<sup>1</sup>). As a result, if there is some latent structure in the subsequence explaining why a particular order is correct, the stacked difference representation will capture it, to the extent of the feature’s capacity.

### 3.4. Multi-length MidRank

So far we focused on subsequences of fixed length  $\lambda$ . A natural extension is to consider multiple subsequence lengths, as different lengths are likely to capture different aspects of the example sequences and subsequences. To train a multi-length ranker we simply consider each length in eq. (4) separately, namely  $\lambda = 2, 3, 4, \dots, L$ . To infer the final ranking we need to fuse the output from the different length rankers. To this end we propose a weighted majority voting scheme.

For each test instance  $\mathbf{X}'$  we obtain a ranking per  $\lambda$  and the respective ranking score from  $z(\cdot)$ . As a result we have rankings for each of the  $L - 1$  rankers. Then, each image in the test sequence gets a vote for its particular position as returned from each ranker. Also, each image gets a voting score that is proportional to the ranking score from  $z(\cdot)$ , and, therefore, the confidence of the ranker for placing that image to the particular position. We weight the image position with the voting scores and compute the weighted votes for all images for all positions. Then, we decide the final position of each image starting rank 1, selecting the image with the highest weighted vote at rank 1. Then we eliminate this image from the subsequence comparisons and continue iteratively, until there are no images to be ranked.

## 4. Efficient inference

Solving eq. (3) requires an explicit search over the space of all possible permutations in  $\mathbf{Y}$ , which amounts to  $l!$ . Hence, for a successful as well as practical inference we need to resolve this combinatorial issue. Inspired by random forests [5] and the best-bin-first search strategies [2], we propose the following greedy search algorithm. For a visual illustration of the algorithm we refer to Figure 2.

We start from an initial ranking solution  $\hat{\mathbf{Y}}^{(0)}$  obtained from a less accurate, but straightforward ranking algorithm (*i.e.* RankSVM). Given  $\hat{\mathbf{Y}}^{(0)}$ , we generate a set of permutations denoted by  $\{\hat{\mathbf{Y}}^{(1)}\}$ , such that the new permutations are obtained by only swapping a pair of elements

<sup>1</sup>users.cecs.anu.edu.au/~basura/ICCV15\_sup.pdf

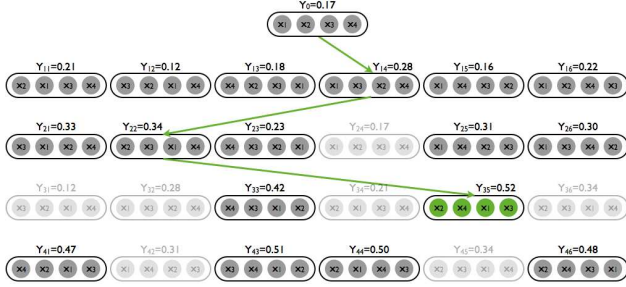


Figure 2: Illustration of the proposed greedy inference algorithm. Starting from the root, we visit all child permutations where only a single pair is switched. Then we select the child permutation with the maximum score. If this maximum score is greater than the parent score, then we expand the maximum child, otherwise we stop. In the process we do not revisit the permutations that have already been visited, see faded nodes. We limit the search procedure to a maximum of  $l$  depths of the tree.

of  $\hat{\mathbf{Y}}^{(0)}$ . From all permutations of  $\{\hat{\mathbf{Y}}^{(1)}\}$ , we compute the ranking scores using Eq. 3 and pick the permutation with the maximum score denoted by  $\hat{\mathbf{Y}}^{(1)}$ . If this score is larger than the score of the parent permutation (i.e.  $\sum z(\mathbf{X}'_j, \hat{\mathbf{Y}}^{(1)}_j; \vartheta, \lambda) > \sum z(\mathbf{X}'_j, \hat{\mathbf{Y}}^{(0)}_j; \vartheta, \lambda)$ ), we set the permutation  $\hat{\mathbf{Y}}^{(1)}$  as the new parent and traverse the solution space recursively using the same strategy. Permutations that have already been visited are removed from any future searches. The procedure of traversing through the above strategy forms a tree of permutations (see Figure 2). We stop when *i*) no other solution with a higher score can be found (score criterion) or *ii*) when we reach the  $l$ -th level of the permutation tree (depth criterion).

At each level of the tree, we traverse a maximum of  $\ell^i \cdot (\ell^i - 1)/2$  nodes. Experimentally, we observe that after only a few levels we satisfy the score criterion, thus having an average complexity of  $O(\ell^2)$ . When the depth criterion is satisfied, we have the worst case complexity of  $O(\ell^3)$ .

To further ensure that the search algorithm has not converged to a poor local optimum, we repeat the above procedure starting from different initial solutions. For maximum coverage of the solution space, we carefully select the new  $\hat{\mathbf{Y}}^{(0)}$ , such that  $\hat{\mathbf{Y}}^{(0)}$  was not seen in the previous trees. As we also show in the experiments, the proposed search strategy allows for obtaining good results using very few trees. Inarguably, our efficient inference enables ranking with subsequences, as the exhaustive search is impractical for moderately sized sequences (8-15 elements long), and intractable for longer sequences.

## 5. Experiments

We select three supervised *image re-ranking applications* to compare MidRank with other supervised learning-to-rank algorithms, namely, ranking public figures (section 5.2), ordering images based on interestingness (section 5.3) and chronological ordering of car images (section 5.4). Next, we analyze the properties and the efficiency of MidRank under different parameterizations in section 5.5.

### 5.1. Evaluation criteria & implementation details

We evaluate all methods with respect to the following ranking metrics. First, we use the *normalized discounted cumulative gain* NDCG, commonly used to evaluate ranking algorithms [19]. The discounted cumulative gain at position  $k$  is defined as  $DCG@k = \sum_{i=1}^k \frac{2^{rel_i} - 1}{\log_2(i+1)}$ , where  $rel_i$  is the relevance of the image at position  $i$ . To obtain the normalized DCG, the  $DCG@k$  score is divided by the ideal DCG score. NDCG, whose range is  $[0, 1]$ , is strongly *non-linear*. For example going from 0.940 to 0.950 indicates a significant improvement.

We also use the *Kendall-Tau*, which captures better how close we are to the perfect ranking. The *Kendall-Tau* accuracy is defined as  $KT = \frac{l^+ - l^-}{0.5l(l-1)}$ , where  $l^+, l^-$  stand for the number of all pairs in the sequence that are correctly and incorrectly ordered respectively, and  $l = l^+ + l^-$ . *Kendall-Tau* varies between  $-1$  and  $+1$  where a perfect ranker will have a score of  $+1$ . For completeness we also use pairwise accuracy as an evaluation criterion in which we count the percentage of correctly ranked pairs of elements in all sequences.

We compare our MidRank with point-wise methods such as SVR [27], McRank [16], pair-wise methods such as RankSVM [13], Relative Attributes [21] and CRR [24]. We also compare with list-wise methods such as AdaRank [34], LambdaMART [32], ListNET [6] and ListMLE [33]. For all methods we use the publically available code as provided by the authors, the same features and recommended settings for fine-tuning. All these baselines and MidRank take the same set of training sequences as input. There is no overlap between elements of train and test sequences. All training sequences are sampled from the training set of each dataset and testing sequences are sampled from the testing set. We make these train and test sequences along with the data publicly available. We evaluate all methods on a large number of 20,000 test sequences. We experimentally found that *the standard deviations are quite small for all methods*.

For MidRank we pick the values for any hyper-parameters (e.g. the cost parameter in eq. (4)) after cross-validation. We investigate MidRank rankers of length 3 – 8 and merge the results with the majority weighted voting, as described in Sect. 3.4. For the efficient inference, we ini-

Method	Public figures			Scene interestingness			Car chronology		
	NDCG	KT	Pair Acc.	NDCG	KT	Pair Acc.	NDCG	KT	Pair Acc.
SVR	.882	.349	65.7	.870	.317	65.8	.910	.399	69.9
McRank	.921	.540	76.9	.859	.295	64.8	.921	.477	70.6
RankSVM	.947	.617	80.8	.870	.317	65.8	.928	.482	74.1
Rel. Attributes	.947	.616	80.8	.869	.315	65.7	.927	.479	73.9
CRR	.945	.612	80.6	.846	.273	63.6	.912	.394	69.7
AdaRank	.836	.154	57.7	.745	-.077	46.1	.827	.118	55.9
LambdaMART	.855	.207	60.4	.860	.315	64.3	.935	.409	70.6
ListNET	.866	.314	65.7	.821	.118	55.9	.872	.291	64.5
ListMLE	.851	.262	63.1	.862	.282	64.1	.854	.278	63.9
MidRank	<b>.954</b>	<b>.722</b>	<b>84.7</b>	<b>.887</b>	<b>.347</b>	<b>67.4</b>	<b>.949</b>	<b>.553</b>	<b>76.9</b>

Table 1: Evaluating ranking methods on three datasets and applications: ranking public figures using relative attributes of [21], ranking scenes according to how interesting they look [11] and ranking cars according to their manufacturing date [15]. For public figures we have sequences of 8 images because of the size of the dataset. For the scenes and the cars datasets we have sequences of 20 images. Similar trends were observed with sequences of 80 images. For all baselines we use the publicly available code and the recommended settings.



Figure 3: Example of ordering images with MidRank according to how interesting they look, or how old is the car they depict. Although both tasks seem rather difficult with the naked eye, MidRank returns rankings very close to the ground truth. We include more visual results in the supplementary material.

tialize the parent node with the solution obtained from the pair-wise RankSVM [13]. We also tried various ranking score normalizations between the different length rankers, but we found experimentally that results did not improve significantly. Consequently, we opted for directly using the unnormalized ranking scores from different length rankers. In all cases the feature vectors are L2-normalized.

## 5.2. Ranking Public Figures

First we evaluate MidRank on ranking public figure images with respect to relative visual attributes [21], using the features and the train/test splits provided by [21]. The dataset consists of images from eight public figures and eleven visual attributes of theirs, such as *big lips*, *white* and *chubby*. Our goal is to learn rankers for these attributes. Since there 8 public figures, we report results on the longest possible test sequence size composed of 8 images. For each of the 11 attributes we sample 10,000 train sequences of

length 8 and 20,000 test sequences, totaling to 220,000 test sequences for all attributes. We use the standard GIST features provided with the dataset. The results are reported by taking the average over all eleven attributes and over all test sequences. See results in Table 1.

We observe that MidRank improves the accuracy of the ranking significantly for all the evaluation criteria. For Kendall-Tau, MidRank brings a +10.5% absolute improvement. It is worth mentioning that for this dataset the best individual MidRank function was of length 7, which in isolation scored a 0.684 Kendall-Tau accuracy.

## 5.3. Ranking Interestingness

Next, we evaluate MidRank on ranking images according to how interesting they are. We consider train and test sequences of size 20. It is not possible to consider much longer sequences as the annotation pool for interestingness is limited in this dataset. In practice even for humans rating

more than 20 images based on interestingness would be difficult. We use the scene categories dataset from [20], whose images were later annotated with respect to interestingness by [11]. We extract GIST [20] features and construct 10,000 train sequences and 20,000 test sequences. Note that this is a difficult task, as interestingness is a subjective criterion which can be attributed to many different factors within an image. See results in Table 1.

We observe that also for this dataset MidRank has a significantly better accuracy than the competitor methods for all the evaluation criteria. For visual results obtained from our method, see Fig. 3 (random example). As we can see, MidRank returns a rank which visually is very close to the ground truth.

#### 5.4. Chronological ordering of cars

As a final application, we consider the task of re-ranking images chronologically. We use the car dataset of [15]. The chronology of the cars is in the range 1920, 1921, ..., 1999. As image representation we use 64-Gaussian Fisher vectors [22] computed on dense SIFT features, after being reduced to 64 dimensions with PCA. To control the dimensionality we also reduce the Fisher vectors to 1000 dimensions using PCA again. Similar to the previous experiment, we generate 10,000 training sequences and 20,000 test sequences of length 20. See results in Table 1.

Again, MidRank obtains significantly better results than the competitor methods for all the evaluation criteria and especially for the Kendall-Tau accuracy (+7.1%). We show some visual results in Fig. 3. We also experimented with training sequence lengths of 5, 10, 15, 20 and testing sequence lengths of 5, 10, 20, 80. Due to brevity and space, we report on test sequences of length 20 only (which seems a more practical scenario in image search applications). However, in all these cases MidRank outperforms all other methods. Note that despite the uncanny resemblance between the cars, especially the older ones, MidRank returns a ranking very close to the true one.

#### 5.5. Detailed analysis of MidRank

**Effect of subsequence length on MidRank.** Next, we evaluate the relation between the MidRank accuracy and the training subsequence sizes. We use sequences of size 20 for training and testing generated from the car dataset. We evaluate different MidRank ranking functions of size 3 up to 8. We plot the results in Fig. 4(a).

For all ranking measures the best MidRank is of size 7. Interestingly, the ranking performance gradually increases up to a point as the training subsequence size increases. This indicates that MidRank ranking models trained on moderately sized subsequences perform better than very small or very large ones. Small MidRank models are easy to train (small training errors), but solve a relatively easy

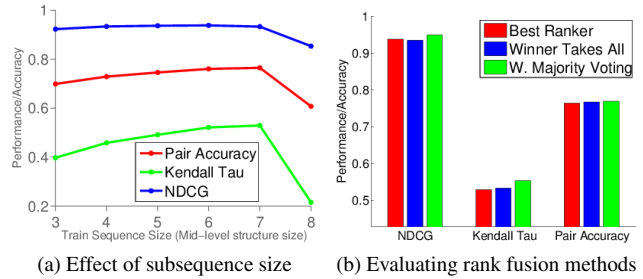


Figure 4: (a) Evaluation of different individual MidRank ranking functions of lengths 3 to 8 on the chronological ordering of cars task. We show how pair accuracy, Kendall Tau and NDCG vary over test sequences of size 20. (b) Comparison of several rank fusion strategies. Weighted majority voting method is the most effective strategy for combining rankers.

ranking sub-problem. In contrast, large MidRank models are more difficult to train (larger training errors). Our experiments suggest that moderately sized subsequences are the most suitable for MidRank. It is also interesting to see that all three ranking measures used are consistent (—see Fig. 4(a)). However, the sensitivity of Kendall Tau seems to be larger than the other two ranking criteria (NDCG and pair-accuracy).

**Evaluating subsequence representations.** In this experiment we evaluate the effectiveness of the stacked difference representation introduced in section 3.3 compared to other alternatives see Fig. 5 (left). The max-pooling or the mean pooling of difference vectors of a subsequence hinders useful ranking information, such as subtle variations between neighboring elements. Full stacked difference vectors (option (c)) does not perform as well as other stacked versions (d) and (e), probably due to the curse of dimensionality. Note that we also evaluated the mean representations on longer subsequences (20 images per sequence) and obtained 0.05 points lower in KT than the proposed stacked representations. This shows the best results are obtained with our stacked difference representation.

**Evaluating efficient inference vs exhaustive inference.** In this section we compare our efficient inference strategy with the exhaustive inference. As explained earlier, the exhaustive strategy has a complexity of  $O(\ell^i!)$ , whereas the proposed efficient inference strategy has an average complexity of  $O(\ell^i^2)$  and a worst case complexity of  $O(\ell^i^3)$ .

First, we plot how the execution time varies during inference for different test sequence sizes. We compare our inference method with the exhaustive search in Fig. 6 (a). For moderately long test sequences, e.g. up to size 8 in this experiment, our method is 50 times faster than exhaustive search. For longer sequences exhaustive inference is not

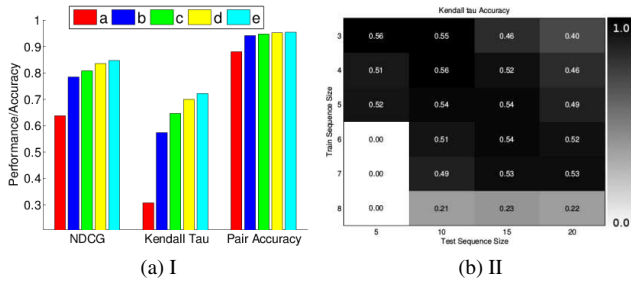


Figure 5: (I) Comparing different representations for mid-level sequences detailed in 3.3 on the public figures dataset, also considering two more choices from the literature. : (a) max pooling of difference vectors as in [4], (b) Mean pairwise difference representations, (c) the full stacked difference vector representation between all elements  $i, j$  in a sequence, (d) the stacked representation, and (e) the stacked difference vector representation. The stacked difference vectors outperform all other alternatives. (II) How KT-accuracy changes when varying the train and test sequence lengths (cars dataset).

even an option, as the number of possible combinations that need to be explored becomes very impractical, or even intractable for longer sequences.

Our inference method discovers the optimal order for a sequence of length 20 in  $0.75 \pm 0.1$  seconds, and in practice sequences of up to 500 elements can be easily processed. Hence, MidRank may easily be employed in the standard supervised image ranking and re-ranking scenarios, *e.g.* improving the image search results based on user preferences.

In Fig. 6(b) we show significant improvement in execution time does not hurt the accuracy with respect to the exhaustive search. In this experiment we vary the number of trees used in our efficient algorithm and report the Kendall-Tau score and the percentage of sequences that agrees with the solution obtained with exhaustive search (blue line in Fig. 6(b)). Interestingly, using a single tree, we obtain a better Kendall-Tau score than the one obtained with the exhaustive search. We attribute this to some degree of over-fitting that might occur during learning. With 3 trees we obtain the same solution as exhaustive search for 97% of the times, whereas with five trees we obtain exactly the same results as the exhaustive search.

**Evaluating Majority Voting Scheme.** Last, we evaluate the effectiveness of different rank fusion strategies for MidRank using the cars dataset. We compare the proposed weighted majority voting with *winner-takes-all* strategy in which the ranker with the highest ranking score is used to define the final ordering. We also compare with the best individual ranker, where we use cross-validation to find the

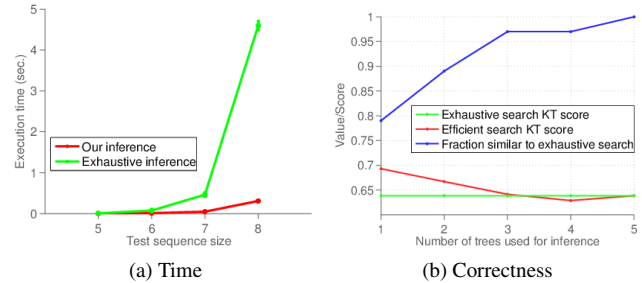


Figure 6: (a) Comparison of execution time between the proposed efficient inference vs exhaustive search on public figures dataset. (b) Our efficient inference algorithm uses multiple trees. This figure shows how ranking performance (Kendall Tau) varies with respect to the number of trees used. The blue plot shows the fraction of solutions (generated by the efficient algorithm) that agree with the solution obtained with exhaustive search.

best ranker given a test sequence length.

As can be seen from Fig. 4(b), the weighted majority voting scheme works best. The results indicate that each ranker from different mid-level structure sizes exploits different types of structural information. Similar conclusions were derived for the other datasets.

**Evaluating on sequences of different lengths.** In this experiment we evaluate how pair accuracy and KT vary for different train and test sequence sizes. We use the cars dataset for this experiment. From results reported in Fig. 5 (right) we see that for smaller test sequences of 5 and 10 the best results are obtained using train subsequences of size 3 or 4. However, for larger test sequences of size 15 and 20 the best results are obtained for train subsequences of size 5, 6 and 7. Interestingly, the largest train subsequence size of 8 reports the worst results. These observations are valid for both pair accuracy as well as *Kendall-Tau* performance.

## 6. Conclusion

In this paper we present a supervised learning to rank method, MidRank, that learns from sub-sequences. A novel *stacked difference vectors* representation and an effective ranking algorithm that uses sub-sequences during the learning is presented. The proposed method obtains significant improvements over state-of-the-art pair-wise and list-wise ranking methods. Moreover, we show that by exploiting the structural information and the regularity in sub-sequences, MidRank allows for a better learning of ranking functions on several image ordering tasks.

**Acknowledgement** Authors acknowledge the support from the PARIS project (IWT-SBO-Nr.110067), iMinds project HiViz and ANR project SoLStiCe (ANR-13-BS02-0002-01).



## References

- [1] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid. Good practice in large-scale learning for image classification. *PAMI*, 36:507–520, 2014.
- [2] J. S. Beis and D. G. Lowe. Shape indexing using approximate nearest-neighbour search in high-dimensional spaces. In *CVPR*, 1997.
- [3] D. Birkhoff. Tres observaciones sobre el algebra lineal. *Universidad Nacional de Tucuman Revista , Serie A*, 5:147151, 1946.
- [4] Y. Boureau, J. Ponce, and Y. LeCun. A theoretical analysis of feature pooling in visual recognition. In *ICML*, 2010.
- [5] L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- [6] Z. Cao, T. Qin, T. Liu, M. Tsai, and H. Li. Learning to rank: from pairwise approach to listwise approach. In *ICML*, 2007.
- [7] W. S. Cooper, F. C. Gey, and D. P. Dabney. Probabilistic retrieval based on staged logistic regression. In *SIGIR*, 1992.
- [8] W. Deng, J. Hu, and J. Guo. Linear ranking analysis. In *CVPR*, 2014.
- [9] P. Dokania, A. Behl, C. Jawahar, and M. Kumar. Learning to rank using high-order information. In *ECCV*, 2014.
- [10] Y. Fu, T. Hospedales, T. Xiang, S. Gong, and Y. Yao. Interestingness prediction by robust learning to rank. In *ECCV*, 2014.
- [11] M. Gygli, H. Grabner, H. Riemenschneider, F. Nater, and L. Van Gool. The interestingness of images. In *ICCV*, 2013.
- [12] R. Herbrich, T. Graepel, and K. Obermayer. Large margin rank boundaries for ordinal regression. In *NIPS*, pages 115–132, 1999.
- [13] T. Joachims. Training linear svms in linear time. In *SIGKDD*, 2006.
- [14] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*. 2012.
- [15] Y. J. Lee, A. Efros, and M. Hebert. Style-aware mid-level representation for discovering visual connections in space and time. In *ICCV*, 2013.
- [16] P. Li, Q. Wu, and C. J. Burges. Mcrank: Learning to rank using multiple classification and gradient boosting. In *NIPS*, pages 897–904, 2007.
- [17] L. Liang and K. Grauman. Beyond comparing image pairs: Setwise active learning for relative attributes. In *CVPR*, 2014.
- [18] G. Lin, C. Shen, and J. Wu. Optimizing ranking measures for compact binary code learning. In *ECCV*, 2014.
- [19] T.-Y. Liu. Learning to rank for information retrieval. *Foundations and Trends in Information Retrieval*, 3(3):225–331, 2009.
- [20] A. Oliva and A. Torralba. Modeling the shape of the scene: A holistic representation of the spatial envelope. *IJCV*, 42(3):145–175, 2001.
- [21] D. Parikh and K. Grauman. Relative attributes. In *ICCV*, 2011.
- [22] F. Perronnin, J. Sánchez, and T. Mensink. Improving the fisher kernel for large-scale image classification. In *ECCV*, 2010.
- [23] T. Schiavinotto and T. Stützle. The linear ordering problem: Instances, search space analysis and algorithms. *Journal of Mathematical Modelling and Algorithms*, 3(4):367–402, 2004.
- [24] D. Sculley. Large scale learning to rank. In *NIPS*, 2009.
- [25] S. Shalev-Shwartz, Y. Singer, N. Srebro, and A. Cotter. Pegasos: Primal estimated sub-gradient solver for svm. *Mathematical programming*, 127(1):3–30, 2011.
- [26] V. Sharmanska, N. Quadrianto, and C. H. Lampert. Learning to rank using privileged information. In *ICCV*, 2013.
- [27] A. J. Smola and B. Schölkopf. A tutorial on support vector regression. *Statistics and computing*, 14(3):199–222, 2004.
- [28] M. Sun, A. Farhadi, and S. M. Seitz. Ranking domain-specific highlights by analyzing edited videos. In *ECCV*, 2014.
- [29] M. Taylor, J. Guiver, S. Robertson, and T. Minka. Soft-rank: optimizing non-smooth rank metrics. In *Proceedings of the 2008 International Conference on Web Search and Data Mining*, pages 77–86. ACM, 2008.
- [30] O. Tuzel, M. Liu, Y. Taguchi, and A. Raghunathan. Learning to rank 3d features. In *ECCV*, 2014.
- [31] T. Wang, S. Gong, X. Zhu, and S. Wang. Person re-identification by video ranking. In *ECCV*, 2014.
- [32] Q. Wu, C. J. Burges, K. M. Svore, and J. Gao. Adapting boosting for information retrieval measures. *Information Retrieval*, 13(3):254–270, 2010.
- [33] F. Xia, T.-Y. Liu, J. Wang, W. Zhang, and H. Li. Listwise approach to learning to rank: theory and algorithm. In *ICML*, pages 1192–1199. ACM, 2008.
- [34] J. Xu and H. Li. Adarank: a boosting algorithm for information retrieval. In *SIGIR*, 2007.
- [35] Y. Yue, T. Finley, F. Radlinski, and T. Joachims. A support vector method for optimizing average precision. In *ACM SIGIR*, pages 271–278. ACM, 2007.