

# Supplemental Document for Combining Task Predictors via Enhancing Joint Predictability

Kwang In Kim<sup>1</sup>[0000-0002-6470-4571], Christian  
Richardt<sup>2</sup>[0000-0001-6716-9845], and Hyung Jin Chang<sup>3</sup>[0000-0001-7495-9677]

<sup>1</sup> UNIST, Korea    <sup>2</sup> University of Bath, UK    <sup>3</sup> University of Birmingham, UK

In this supplemental document, we present:

1. Details of the marginal likelihood calculation used in the automatic determination of relevant predictors  $\Sigma_L$  (Sec. 1.1);
2. A summary of our predictor combination algorithm (Sec. 1.2);
3. A detailed discussion of baseline algorithms, including:
  - (a) our adaptation of Mejjati et al.’s multi-task learning (*MTL*) algorithm [6] (Sec. 2.1),
  - (b) a derivation of Kim et al.’s original predictor combination (*OPC*) algorithm [5,4] (Sec. 2.2), and
  - (c) Evgeniou et al.’s graph Laplacian (*GL*)-based MTL algorithm and its adaptation to predictor combination (Sec. 2.3).

In the main paper, we only presented the ranking results for the first 10 attributes in each dataset. In Section 3, we provide the complete experimental results, including additional results of *GL* and tests of statistical significance of the accuracy improvements made by different algorithms. We reproduce some content from the main paper to make this document self-contained.

## 1 Details of the main algorithm

### 1.1 Calculating the marginal likelihood for linear Gaussian process prediction

Suppose that we have the following linear and nonlinear anisotropic covariance functions:

$$k_L(\mathbf{a}, \mathbf{b}) = \mathbf{a}^\top \Sigma_L \mathbf{b}, \quad (1)$$

$$k_A(\mathbf{a}, \mathbf{b}) = \exp(-(\mathbf{a} - \mathbf{b})^\top \Sigma_A (\mathbf{a} - \mathbf{b})), \quad (2)$$

where  $\Sigma_L = \text{diag}[\boldsymbol{\sigma}]$ , the diagonal matrix with elements  $\boldsymbol{\sigma} = [\sigma^1, \dots, \sigma^n]^\top$ .  $\Sigma_A$  is defined similarly. Our goal is to maximize the marginal likelihood  $p(\mathbf{f} | G, \Sigma_L)$  of the sampled predictor  $\mathbf{f}$  given the reference matrix  $G$  with respect to  $\boldsymbol{\sigma}$ . The log

marginal likelihood  $\log(p(\mathbf{f}|G, \Sigma_L))$  of linear Bayesian regression with Gaussian prior and i.i.d. Gaussian noise model is given as [7]:

$$\begin{aligned} \log(p(\mathbf{f}|G)) &= -\frac{1}{2} \log|G \cdot \text{diag}[\boldsymbol{\sigma}] \cdot G^\top + \lambda I| - \frac{N}{2} \log(2\pi) \\ &\quad - \frac{1}{2} \mathbf{f}^\top (G \cdot \text{diag}[\boldsymbol{\sigma}] \cdot G^\top + \lambda I)^{-1} \mathbf{f}. \end{aligned} \quad (3)$$

As maximizing  $p(\mathbf{f}|G, \Sigma_L)$  is equivalent to minimizing  $-\log(p(\mathbf{f}|G, \Sigma_L))$ , and the second term in  $\log(p(\mathbf{f}|G, \Sigma_L))$  is independent of  $\boldsymbol{\sigma}$ , we can find the optimal parameter vector  $\boldsymbol{\sigma}^*$  by minimizing the following energy:

$$\begin{aligned} \mathcal{E}(\boldsymbol{\sigma}) &= \log|G \cdot \text{diag}[\boldsymbol{\sigma}] \cdot G^\top + \lambda I| + \mathbf{f}^\top (G \cdot \text{diag}[\boldsymbol{\sigma}] \cdot G^\top + \lambda I)^{-1} \mathbf{f} \\ &= N \log|\lambda| + \sum_{i=1}^n \log(\sigma^i) + \log|\text{diag}[1./\boldsymbol{\sigma}] + G^\top G / \lambda| \\ &\quad + \mathbf{f}^\top \left( \frac{1}{\lambda} I - \frac{1}{\lambda^2} G \left( \text{diag}[1./\boldsymbol{\sigma}] + \frac{1}{\lambda} G^\top G \right)^{-1} G^\top \right) \mathbf{f}, \end{aligned} \quad (4)$$

where the second equation is obtained by applying the Sherman–Morrison–Woodbury formula [9] to both summands of  $\mathcal{E}$ , and ‘ $1./\boldsymbol{\sigma}$ ’ is the element-wise reciprocal of  $\boldsymbol{\sigma}$ . Since  $\frac{1}{\lambda} \mathbf{f}^\top \mathbf{f}$  and  $N \log|\lambda|$  are also independent of  $\boldsymbol{\sigma}$ , minimizing  $\mathcal{E}$  is equivalent to minimizing

$$\begin{aligned} \mathcal{E}'(\boldsymbol{\sigma}) &= \sum_{i=1}^n \log(\sigma^i) + \log|\text{diag}[1./\boldsymbol{\sigma}] + G^\top G / \lambda| \\ &\quad - \frac{1}{\lambda^2} \mathbf{f}^\top \left( G \left( \text{diag}[1./\boldsymbol{\sigma}] + \frac{1}{\lambda} G^\top G \right)^{-1} G^\top \right) \mathbf{f}. \end{aligned} \quad (5)$$

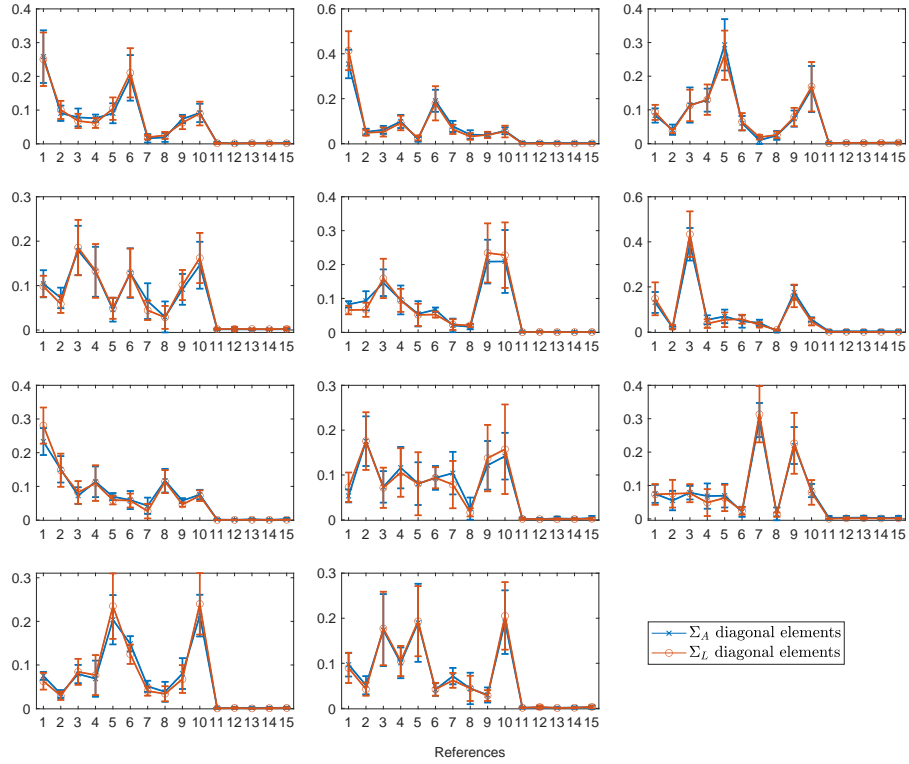
Since this energy  $\mathcal{E}'$  is a continuously differentiable function of  $\boldsymbol{\sigma}$ , it can be minimized by standard gradient descent. Figure 1 shows example parameters  $\boldsymbol{\sigma}^*$  optimized for the *Pubfig* dataset. For each of the 11 attributes in *Pubfig* as a target, we optimized the corresponding parameters  $\boldsymbol{\sigma}$  using the remaining attributes as references, plus 5 additional randomly generated references. As indicated by small magnitudes and the corresponding standard deviations of the  $\boldsymbol{\sigma}^*$  entries, our algorithm successfully disregards these irrelevant references. For comparison, we also show the corresponding parameters optimized for the anisotropic Gaussian kernel  $k_A$  (Eq. 2), demonstrating that once normalized, their relative scaling behaviors are similar, i.e.

$$\Sigma_A^* \approx \frac{\Sigma_L^*}{\sigma_k^2} \quad (6)$$

for a global scaling parameter  $\sigma_k^2$ . Our final algorithm uses  $\Sigma_L^* / \sigma_k^2$  as a surrogate to  $\Sigma_A^*$ , using  $\sigma_k^2$  as hyperparameter.

## 1.2 Algorithm summary

Given the reference matrix  $G$ , the initial predictor  $\mathbf{f}^0$ , and hyperparameters (noise variance  $\sigma^2$  in Eq. 8; global kernel scaling  $\sigma_k^2$  in Eq. 6; regularization parameter  $\lambda_J$



**Fig. 1.** The average diagonal values of  $\Sigma_A$  and  $\Sigma_L$  optimized for each attribute in the *Pubfig* dataset as the target with the remaining 10 attributes in the same dataset (references 1 to 10), plus 5 additional randomly generated attributes (references 11 to 15) as references. The values of  $\Sigma_A$  and  $\Sigma_L$  are normalized such that the respective sums total to one. Note that the unrelated references (11–15) are correctly detected as irrelevant (small magnitudes and standard deviations) and hence ignored. In addition, the linear kernel  $\Sigma_L$  (orange) is highly correlated to the anisotropic kernel  $\Sigma_A$  (blue), so we use a scaled version of  $\Sigma_L$  as a surrogate to the optimal  $\Sigma_A^*$  (Eq. 6).

in Eq. 7), our algorithm constructs a denoised predictor by iteratively maximizing the objective

$$\mathcal{O}_N(\mathbf{f}) = \frac{\mathbf{f}^\top \mathbf{A} \mathbf{f}}{\mathbf{f}^\top \mathbf{C}_N \mathbf{f}} \quad \text{with} \quad (7)$$

$$\mathbf{A} = (\mathbf{C}_N \mathbf{f}^t)(\mathbf{C}_N \mathbf{f}^t)^\top + \lambda_J \mathbf{Q}' \quad \text{and}$$

$$\mathbf{Q}' = \mathbf{C}_N (2\mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} - (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{K} \mathbf{K} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1}) \mathbf{C}_N. \quad (8)$$

Algorithm 1 summarizes this process.

---

**Algorithm 1** Nonlinear predictor combination

---

**Input:** Initial predictor  $\mathbf{f}^0$ , references  $\{\mathbf{g}_i\}_{i=1}^R$ , noise variance  $\sigma^2$  (Eq. 8), global kernel scaling  $\sigma_k^2$ , regularization parameter  $\lambda_J$  (Eq. 7), and iteration number  $S$ .

1:  $\mathcal{H}^0 = \{\mathbf{h}_0^0, \mathbf{h}_1^0, \dots, \mathbf{h}_R^0\} \leftarrow \{\mathbf{f}^0, \mathbf{g}_1, \dots, \mathbf{g}_R\}$ ;

2: Calculate the kernel parameter matrix  $(\Sigma_L^*)_i$  for each  $\mathbf{h}_i^0 \in \mathcal{H}^0$ ;  $(\Sigma_A)_i = (\Sigma_L^*)_i / \sigma_k^2$ ;

3: **for** step  $t \in \{0, \dots, S-1\}$  **do**

4:   **for** reference  $i \in \{0, \dots, R\}$  **do**

5:     Calculate  $\mathbf{h}_i^{t+1} \in \mathcal{H}^{t+1}$  by maximizing  $\mathcal{O}_N$  based on  $\mathcal{H}^t$  (Eq. 7);

6:   **end for**

7: **end for**

**Output:** Denoised target predictor  $\mathbf{f}^* = \mathbf{h}_0^S$ .

---

## 2 (Adapting) Existing algorithms

### 2.1 Mejjati et al.’s *MTL* algorithm

Mejjati et al.’s *MTL* algorithm considers each task-specific predictor as a random variable. Then, the relationships between tasks are modeled based on the statistical dependence estimated by evaluating these predictor random variables on a dataset  $X$  [6]. Adopting a nonparametric measure of statistical dependence, the finite set independence criterion (FSIC) [3], *MTL* enables training multiple predictors independently of their parametric forms and, therefore, it can be applied to predictor combination problems.

Applying this algorithm to the predictor combination setting, we construct the initial predictor matrix  $H^0$  by stacking column-wise, the initial target predictor  $\mathbf{f}^0$  and the references  $\{\mathbf{g}_1, \dots, \mathbf{g}_R\}$

$$H^0 = [\mathbf{f}^0, \mathbf{g}_1, \dots, \mathbf{g}_R]. \quad (9)$$

*MTL* then refines the initial predictor matrix  $H^0$  by minimizing the energy

$$\mathcal{E}_M(H) = \|H - H^0\|_F^2 - \lambda_1 \|\text{vec}(\Phi(H))\|_2^2 + \lambda_2 \|\text{vec}(\Phi(H))\|_1, \quad (10)$$

where  $\text{vec}(A)$  constructs a vector by concatenating columns of matrix  $A$ ,  $\Phi(H)$  is an  $(R+1) \times (R+1)$ -sized matrix consisting of pairwise FSIC evaluations:  $\Phi(H)_{[i,j]}$  takes a large positive value when  $H_{[:,i]}$  and  $H_{[:,j]}$  exhibit strong statistical dependence and it takes 0 when  $H_{[:,i]}$  and  $H_{[:,j]}$  are independent as realizations of random variables. Minimizing  $\mathcal{E}_M$  strengthens overall task dependence via (negation of) the  $L^2$  norm of  $\text{vec}(\Phi(H))$  and, at the same time, introduces sparsity in the task dependence via the  $L^1$  norm of  $\Phi(H)$ . Combing these two terms, *MTL* selectively enforces task dependence while suppressing the dependence of weakly related tasks as outliers. As  $\mathcal{E}_M$  is not differentiable, standard gradient-descent type algorithms are not applicable. Instead, it is minimized based on the alternating direction method of multipliers (ADMM) approach. This involves iteratively solving ADMM sub-problems [1], with the number of total iterations  $S$  as a hyperparameter. Once the optimal predictor matrix  $H^*$  is constructed, the final denoised predictor is obtained by extracting the first column of  $H^*$ :  $\mathbf{f}^* = H^*_{[:,1]}$ . Similarly to our algorithm,  $S$  is determined by setting the maximum number of iterations at 50

and selecting the iteration number achieving the highest validation accuracy. The other two hyperparameters,  $\lambda_1$  and  $\lambda_2$ , are tuned based on validation accuracy.

## 2.2 Derivation of Kim et al.’s algorithm (*OPC*).

Kim et al.’s original predictor combination (*OPC*) approach iteratively minimizes the following energy (Eq. 1 in the main paper):

$$\mathcal{E}_O(f) = D_{\text{KL}}(f | f^t)^2 + \lambda_O \sum_{i=1}^R w_i D_{\text{KL}}(f | g_i)^2, \quad (11)$$

$$w_i = \exp\left(-\frac{D_{\text{KL}}(f^t | g_i)^2}{\sigma_O^2}\right), \quad (12)$$

with  $\lambda_O, \sigma_O^2 > 0$  being hyperparameters. We present how this algorithm is obtained as an instance of Hein and Maier’s *Manifold Denoising* algorithm [2] by discretizing a diffusion process on a predictor manifold  $\mathcal{M}$ .

*Manifold denoising* [2]. Suppose that we have a set of data points  $\mathcal{H}^0 = \{\mathbf{h}_i^0\}_{i=1}^n$  presented as a sample from a Euclidean space  $\mathbb{R}^d$  and, further, that the points in  $\mathcal{H}^0$  are sampled from an underlying data-generating manifold  $\mathcal{M}$  embedded in  $\mathbb{R}^d$  ( $\iota(\mathcal{M}) \subset \mathbb{R}^d$  with  $\iota$  being the embedding), and they are observed as a subset of  $\mathbb{R}^d$  contaminated with i.i.d. Gaussian noise  $\epsilon$  in  $\mathbb{R}^d$ :

$$\mathbf{h}_i^0 = \iota(\tilde{\mathbf{h}}_i^0) + \epsilon \in \mathbb{R}^d \quad \text{for } \tilde{\mathbf{h}}_i^0 \in \mathcal{M}. \quad (13)$$

The manifold denoising algorithm denoises  $\mathcal{H}^0$  by simulating diffusion on a graph  $G$  that discretizes  $\mathcal{M}$  (each point  $\mathbf{h}_i^0 \in \mathcal{H}^0$  forms a vertex of  $G$ ):

$$\frac{\partial H}{\partial t} = -\delta L H, \quad (14)$$

where  $H = [\mathbf{h}_1, \dots, \mathbf{h}_n]^\top$  and  $L$  is the graph Laplacian:

$$L = I - D^{-1}W, \quad (15)$$

$$W_{[i,j]} = \exp\left(-\frac{\|\mathbf{h}_i - \mathbf{h}_j\|^2}{\sigma^2}\right), \quad (16)$$

and  $D$  is a diagonal matrix consisting of row sums of  $W$ , such that  $D_{ii} = \sum_{j=1} W_{[i,j]}$ . Now discretizing Eq. 14 using the implicit Euler method, we obtain

$$H^{t+1} - H^t = -\delta L H^{t+1}. \quad (17)$$

At each time step  $t$ , the solution  $H^{t+1}$  of Eq. 17 is obtained as the minimizer of the following energy:

$$\mathcal{E}_D(H) = \|H - H^t\|_F^2 + \delta \text{tr}[H^\top L H], \quad (18)$$

where  $\|A\|_F$  and  $\text{tr}[A]$  are the Frobenius norm and trace of matrix  $A$ , respectively. As the number of data points  $n$  grows to infinity,  $G$  becomes a precise representation of  $\mathcal{M}$  embedded in  $\mathbb{R}^d$ , and  $L$  converges to the Laplace-Beltrami operator  $\Delta_{\mathcal{M}}$  on  $\mathcal{M}$  casting Eq. 17 into a diffusion process on a continuous manifold

$\mathcal{M}$  [2]. It should be noted that the graph Laplacian was constructed based on the *ambient*  $L^2$  distance in  $\mathbb{R}^d$  rather than the intrinsic metric on  $\mathcal{M}$ . This facilitates building a practical, still consistent algorithm: Equation 18 only requires the ambient Euclidean distance (via  $L$ ) without having to access  $\mathcal{M}$  directly, but it guarantees the statistical consistency of  $L$  as proven by Hein and Maier [2].

Now applying this algorithm to the predictor combination setting and, therefore, assuming that only the first point  $\mathbf{h}_1^0 \in \mathcal{H}^0$  is noisy (i.e.  $\mathbf{h}_i = \iota(\tilde{\mathbf{h}}_i)$  for  $i = \{2, \dots, n\}$ ), we obtain an iterative update rule of  $\mathbf{h}_1^t$  given fixed *references*  $\{\mathbf{h}_i\}_{i=2}^n$ :

$$\mathcal{E}_D(\mathbf{h}) = \|\mathbf{h} - \mathbf{h}_1^t\|^2 + \delta \sum_{i=2}^n W_{[1,i]} \|\mathbf{h} - \mathbf{h}_i\|^2. \quad (19)$$

Finally,  $\mathcal{E}_O$  in Eq. 11 is obtained by replacing each point in  $\mathcal{H}^t$  and the corresponding  $L^2$  distances in  $\mathcal{E}_D$  with a Gaussian process predictor and Kullback–Leibler divergences, respectively:  $\mathbf{h}_1^t$  and  $\{\mathbf{h}_i\}_{i=2}^n$  (with  $n = R+1$ ) are considered as the target predictor  $\mathbf{f}$  and the corresponding references  $\{\mathbf{g}_i\}$ , respectively.

### 2.3 Evgeniou et al.’s graph Laplacian ( $GL$ )-based MTL algorithm.

Evgeniou et al.’s graph Laplacian ( $GL$ )-based algorithm learns predictors  $\mathcal{H} = \{h_i\}_{i=1}^n$  of multiple tasks by enforcing pairwise parameter similarities: Assuming that all predictors are linear, i.e.  $h_i(\mathbf{x}) = \mathbf{w}_i^\top \mathbf{x}$ , their algorithm estimates the predictor parameters  $W = \{\mathbf{w}_1, \dots, \mathbf{w}_n\}$  by minimizing the energy

$$\mathcal{E}_{GL}(W) = \sum_{i=1}^n l_i(h_i) + \lambda_1 \sum_{i=1}^n \|\mathbf{w}_i\|^2 + \lambda_2 \sum_{i=1}^n \sum_{j \neq i}^n U_{[i,j]} \|\mathbf{w}_i - \mathbf{w}_j\|^2, \quad (20)$$

where  $\{l_i(\cdot)\}_{i=1}^n$  are task-specific loss functions and  $U_{[i,j]} \geq 0$  represents the relationship between tasks  $i$  and  $j$ . Now adapting this algorithm to the predictor combination setting, we assume that the initial predictor  $f^0 = h_1$  and the corresponding references  $g_i = h_{i+1}$  for  $i \in \{1, \dots, R\}$  are given ( $n = R+1$ ). Then,  $f^0$  is refined by minimizing the energy

$$\mathcal{E}_{GL}(\mathbf{w}) = \|\mathbf{w} - \mathbf{w}_1^0\|^2 + \lambda_{GL} \sum_{j=2}^n U_{[1,j]} \|\mathbf{w} - \mathbf{w}_j\|^2. \quad (21)$$

In general, determining the task relationship parameters  $\{U_{[1,j]}\}$  is a challenging problem. Here, we determine them by adopting Kim et al.’s approach: We iteratively update  $\{U_{[1,j]}\}$  by minimizing  $\mathcal{E}_{GL}$  at each time step  $t$  with

$$U_{[1,j]} = \exp\left(-\frac{\|\mathbf{w}_i - \mathbf{w}_1^t\|^2}{\sigma_{GL}^2}\right). \quad (22)$$

Further, adopting Kim and Chang’s approach [4], we explicitly constrain all predictor parameter vectors to have unit norm:  $\|\mathbf{w}_i\| = 1$ , enabling the comparison of task predictor parameters independently of their scales. The two hyperparameters  $\sigma_{GL}^2$  and  $\lambda_{GL}$  are determined based on validation accuracy. As often, ranking problems are nonlinear, we extend this framework by adopting

the linear-in-parameter model:

$$f(\mathbf{x}) = \phi(\mathbf{x})^\top \mathbf{w}_f, \quad (23)$$

where  $\phi: \mathcal{X} \rightarrow \mathcal{F}_k$  with  $\mathcal{F}_k$  being the reproducing kernel Hilbert space (RKHS) corresponding to a Gaussian kernel with hyperparameter  $\sigma_k^2$  [8]:

$$k(\mathbf{a}, \mathbf{b}) = \exp\left(-\frac{\|\mathbf{a} - \mathbf{b}\|^2}{\sigma_k^2}\right). \quad (24)$$

In this case, the target predictor  $f$  is represented based the original parameter vector  $\mathbf{w}_f$  as well as its dual parameter vector  $\mathbf{a}_f = [a_f^1, \dots, a_f^{N'}]^\top$ :

$$f(\mathbf{x}) := \phi(\mathbf{x})^\top \mathbf{w}_f = \sum_{j=1}^{N'} a_f^j k(\mathbf{b}_j, \mathbf{x}), \quad (25)$$

with  $\{\mathbf{b}_i\}_{i=1}^{N'}$  being a set of *basis vectors*. The reference predictors  $\{g_i\}_{i=1}^R$  are represented similarly:

$$g_i(\mathbf{x}) := \phi(\mathbf{x})^\top \mathbf{w}_i = \sum_{j=1}^{N'} a_i^j k(\mathbf{b}_j, \mathbf{x}). \quad (26)$$

Under this setting, the parameter similarity  $\|\mathbf{w}_f - \mathbf{w}_i\|$  can be calculated using the standard kernel trick [8] as

$$\|\mathbf{w}_f - \mathbf{w}_i\| = \mathbf{a}_f^\top K \mathbf{a}_f^\top + \mathbf{a}_i^\top K \mathbf{a}_i^\top - 2\mathbf{a}_f^\top K \mathbf{a}_i^\top, \quad (27)$$

with  $K_{[i,j]} = k(\mathbf{b}_i, \mathbf{b}_j)$ . It should be noted that efficient<sup>1</sup> calculation of  $\|\mathbf{w}_f - \mathbf{w}_i\|$  based on Eq. 27 requires that all predictors should share the same RKHS determined by the kernel parameter  $\sigma_k$ . To facilitate this, in our experiments, we first determine  $f^0$  as nonlinear rank support vector machine that minimizes the regularized energy

$$\mathcal{E}_S(f) = \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in U} \mathcal{L}(f, (\mathbf{x}_i, \mathbf{x}_j)) + C_f \|\mathbf{w}\|^2, \quad (28)$$

$$\mathcal{L}(f, (\mathbf{a}, \mathbf{b})) = \max(1 - (f(\mathbf{a}) - f(\mathbf{b})), 0)^2 \quad (29)$$

for the rank loss  $\mathcal{L}$  defined on ground-truth ranked pairs  $U \subset X \times X$  and tune the hyperparameters  $\sigma_k^2$  and  $C_f$  based on validation accuracy. Once  $f^0$  is fixed in this way, the reference predictors  $\{g_i\}_{i=1}^R$  are determined by minimizing  $\mathcal{E}_S$  for the respective rank labels. However, for these references, only the respective regularization hyperparameters  $\{C_i\}_{i=1}^R$  are tuned while the corresponding kernel parameters are all fixed as  $\sigma_k^2$  (optimized for  $f^0$ ), to facilitate the computation of  $\|\mathbf{w} - \mathbf{w}_i\|$  (Eq. 27). We fixed  $N'$  at 500 and selected the basis vectors  $\{\mathbf{b}_i\}_{i=1}^{N'}$  as the cluster centers of input data points  $X$ , estimated using  $k$ -means clustering.

Note that this setting violates the application conditions of predictor combination: It requires access to the forms of all predictors  $\{f, g_1, \dots, g_R\}$  and, further,

<sup>1</sup> The RKHS  $\mathcal{F}_k$  corresponding to a Gaussian kernel  $k$  is infinite-dimensional. Therefore, each parameter vector  $\mathbf{w} \in \mathcal{F}_k$  is an infinite-dimensional object, making the direct evaluation of  $\|\mathbf{w}_f - \mathbf{w}_j\|$  infeasible.

it assumes that all predictors share the same form (Eqs. 25 and 26). We show in Sec. 3 that the latter homogeneity requirement poses a severe limitation on predictor combination performance. Even when *GL* took advantage of known predictor forms, except for a few cases, the other predictor combination algorithms significantly outperformed *GL*. Often, the results of *GL* are even worse than the initial predictors  $\mathbf{f}^0$  that are obtained by selecting the best predictors (via validation) from the heterogeneous predictor pools.

### 3 Complete ranking results

Table 1 summarizes the results for the relative attributes ranking experiments (see Tables 2–6 for complete results). Our algorithm *NPC* performs best for 87% (162/186) of attributes. In particular, it showed statistically significant improvement on 74 out of 80 *AWA2* attributes, while the baselines *OPC* and *MTL* achieved significant performance gains only on 8 and 33 attributes, respectively.<sup>2</sup>

Overall, our algorithm *NPC* is often statistically significantly better than these methods and – apart from only one attribute (for *CUB*) out of 186 – ours is not statistically significantly worse than the other methods. This demonstrates that the baselines *OPC* and *MTL* are limited in that they can only capture pairwise dependence between the target predictor and each reference. Taking into account the dependence present among the references, and thereby *jointly* exploiting them in improving the target predictor, our algorithm *NPC (Ours)* significantly improves the performance.

## References

1. Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J.: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning* **3**(1), 1–122 (2010). <https://doi.org/10.1561/22000000016> 4
2. Hein, M., Maier, M.: Manifold denoising. In: NIPS. pp. 561–568 (2007) 5, 6
3. Jitkrittum, W., Szabó, Z., Gretton, A.: An adaptive test of independence with analytic kernel embeddings. In: PMLR (Proc. ICML). pp. 1742–1751 (2017) 4
4. Kim, K.I., Chang, H.J.: Joint manifold diffusion for combining predictions on decoupled observations. In: CVPR. pp. 7549–7557 (2019) 1, 6
5. Kim, K.I., Tompkin, J., Richardt, C.: Predictor combination at test time. In: ICCV. pp. 3553–3561 (2017) 1
6. Mejjati, Y.A., Cosker, D., Kim, K.I.: Multi-task learning by maximizing statistical dependence. In: CVPR. pp. 3465–3473 (2018) 1, 4

<sup>2</sup> We used a t-test with  $\alpha = 0.95$ . Note that statistical significance tests do not necessarily evaluate how significant the improvements are in an absolute scale: Even when the improvements are marginal, if they are consistent, the result of statistical significant tests can be positive. For instance, for *AWA2* attribute 4, *MTL* achieved rather moderate improvements (with mean 0.05) but the test of statistical significance is positive as the results consistently improved the performance from the baseline, as indicated by the small standard deviation.



**Table 1.** A summary of the results of statistical significance tests of our method NPC compared to baseline  $\mathbf{f}^0$ , *GL*, *OPC* and *MTL*, based on a t-test with  $\alpha = 0.95$ . For each method, we show #attributes where our *NPC* is statistically significantly better (first column), on par with (second column), and statistically significantly worse (third column).

Dataset	vs. baseline $\mathbf{f}^0$			vs. <i>GL</i>			vs. <i>OPC</i>			vs. <i>MTL</i>			# total attr.
<i>Shoes</i>	9	1	0	10	0	0	8	2	0	9	1	0	10
<i>Pubfig</i>	11	0	0	11	0	0	11	0	0	11	0	0	11
<i>OSR</i>	6	0	0	2	4	0	4	2	0	5	1	0	6
<i>OSR (ResNet)</i>	6	0	0	6	0	0	6	0	0	6	0	0	6
<i>aPascal</i>	25	4	0	16	13	0	9	20	0	19	10	0	29
<i>CUB</i>	31	9	0	23	17	0	26	14	0	12	27	1	40
<i>AWA2</i>	74	6	0	73	7	0	71	9	0	72	8	0	80
<i>Zap50K</i>	0	4	0	0	4	0	0	4	0	0	4	0	4
# total attr.	162	24	0	141	45	0	135	51	0	134	51	1	186

7. Rasmussen, C.E., Williams, C.K.I.: Gaussian Processes for Machine Learning. MIT Press, Cambridge, MA (2006) 2
8. Schölkopf, B., Smola, A.J.: Learning with Kernels. MIT Press, Cambridge, MA (2002) 7
9. Sherman, J., Morrison, W.J.: Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. The Annals of Mathematical Statistics 21(1), 124–127 (1950) 2

**Table 2.** Ranking accuracies of different predictor combination algorithms on the *Shoes*, *Pubfig*, *OSR*, and *OSR (ResNet)* datasets. For each dataset, we repeated experiments 10 times with different training, validation, and test set splits. For baseline  $\mathbf{f}^0$  (second column), Kendall’s Tau correlations $\times 100$  (standard deviations in parentheses) are presented. For the remaining algorithms (third to sixth columns), the accuracy offsets from  $\mathbf{f}^0$  are presented. The best and second best results are highlighted with **bold** and *italic* fonts, respectively. The results of statistical significance test based on a t-test with  $\alpha=0.95$  are highlighted in **green** (significantly positive) and **orange** (significantly negative). The last three columns show the results of statistical significance test of our algorithm with *GL*, *OPC*, and *MTL*, respectively (+/-: significantly positive/negative).

<i>Shoes</i>									
Attr.	Baseline $\mathbf{f}^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>	
1	72.09 (1.71)	-0.36 (1.44)	<i>2.36 (0.79)</i>	2.03 (0.49)	<b>3.21 (0.86)</b>	+	+	+	
2	63.84 (1.87)	-2.04 (2.94)	<i>1.57 (1.25)</i>	0.70 (0.39)	<b>2.26 (1.38)</b>	+	0	+	
3	38.07 (2.11)	-1.38 (2.43)	-0.24 (0.65)	<i>0.16 (0.70)</i>	<b>4.58 (2.45)</b>	+	+	+	
4	<i>50.10 (2.75)</i>	<b>-2.45 (2.59)</b>	-0.88 (3.29)	-0.08 (0.51)	<b>1.63 (2.25)</b>	+	+	+	
5	65.76 (1.20)	-1.11 (1.84)	-0.05 (0.16)	<i>0.13 (0.34)</i>	<b>0.92 (2.46)</b>	+	0	0	
6	65.02 (1.83)	<b>-0.86 (1.13)</b>	0.68 (0.87)	<i>0.81 (0.86)</i>	<b>4.18 (1.54)</b>	+	+	+	
7	59.38 (2.06)	<b>-3.31 (1.59)</b>	<i>0.78 (1.16)</i>	0.45 (0.42)	<b>4.14 (3.20)</b>	+	+	+	
8	56.85 (2.04)	<b>-2.57 (1.46)</b>	0.19 (0.46)	<i>0.40 (0.51)</i>	<b>2.62 (0.87)</b>	+	+	+	
9	65.15 (1.94)	0.35 (1.94)	<i>2.49 (1.27)</i>	1.40 (0.72)	<b>4.58 (1.72)</b>	+	+	+	
10	72.10 (1.24)	<b>-1.26 (1.55)</b>	<i>1.71 (0.77)</i>	1.47 (0.77)	<b>2.75 (1.05)</b>	+	+	+	
<i>Pubfig</i>									
Attr.	Baseline $\mathbf{f}^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>	
1	67.13 (2.75)	<b>4.89 (2.48)</b>	8.37 (3.84)	<i>9.37 (2.94)</i>	<b>15.45 (2.59)</b>	+	+	+	
2	62.49 (2.41)	-0.96 (2.62)	-0.31 (0.91)	<i>2.24 (1.55)</i>	<b>13.78 (3.23)</b>	+	+	+	
3	68.31 (2.33)	2.27 (3.27)	<b>3.06 (2.52)</b>	<i>6.25 (3.04)</i>	<b>11.33 (3.10)</b>	+	+	+	
4	63.98 (3.46)	<i>8.44 (4.26)</i>	4.88 (3.14)	7.84 (3.49)	<b>17.80 (4.26)</b>	+	+	+	
5	61.27 (2.96)	<i>6.15 (2.92)</i>	3.33 (4.23)	3.26 (3.70)	<b>16.62 (6.02)</b>	+	+	+	
6	81.60 (1.26)	-1.09 (2.67)	-0.03 (1.25)	<i>0.44 (1.58)</i>	<b>6.17 (2.03)</b>	+	+	+	
7	64.23 (2.88)	<b>2.87 (3.13)</b>	1.68 (2.67)	<i>3.14 (4.61)</i>	<b>15.66 (3.14)</b>	+	+	+	
8	66.10 (3.53)	<i>0.38 (2.66)</i>	0.19 (0.21)	0.10 (0.49)	<b>12.16 (3.17)</b>	+	+	+	
9	59.73 (4.79)	<i>3.58 (3.78)</i>	-0.11 (1.89)	1.96 (3.20)	<b>17.74 (4.68)</b>	+	+	+	
10	63.58 (3.48)	<i>5.79 (3.55)</i>	3.24 (1.92)	4.06 (2.23)	<b>14.16 (2.58)</b>	+	+	+	
11	69.12 (2.87)	<b>7.76 (2.73)</b>	9.30 (2.48)	<i>9.49 (2.73)</i>	<b>15.20 (2.87)</b>	+	+	+	
<i>OSR</i>									
Attr.	Baseline $\mathbf{f}^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>	
1	88.57 (0.93)	<b>3.16 (1.07)</b>	2.06 (0.83)	2.19 (1.03)	<i>2.72 (1.39)</i>	0	0	0	
2	87.52 (0.89)	<b>-1.17 (1.06)</b>	-0.00 (0.19)	<i>0.03 (0.12)</i>	<b>0.93 (0.69)</b>	+	+	+	
3	76.12 (0.95)	0.50 (1.32)	1.31 (0.92)	<i>1.99 (1.26)</i>	<b>3.25 (1.49)</b>	+	+	+	
4	77.67 (0.92)	<i>1.29 (1.11)</i>	0.69 (0.70)	0.90 (0.68)	<b>2.10 (1.29)</b>	0	+	+	
5	79.58 (0.65)	<i>2.50 (0.72)</i>	2.26 (0.68)	1.43 (0.86)	<b>2.89 (1.04)</b>	0	0	+	
6	80.49 (1.22)	<i>0.70 (1.00)</i>	0.09 (0.52)	0.03 (0.43)	<b>1.46 (0.84)</b>	0	+	+	
<i>OSR (ResNet)</i>									
Attr.	Baseline $\mathbf{f}^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>	
1	96.12 (0.59)	-0.07 (0.14)	<i>0.39 (0.33)</i>	0.26 (0.37)	<b>1.33 (0.48)</b>	+	+	+	
2	84.73 (0.87)	0.00 (0.27)	-0.11 (0.25)	<i>0.10 (0.25)</i>	<b>2.51 (1.14)</b>	+	+	+	
3	84.46 (1.08)	-0.01 (0.03)	<b>0.43 (0.30)</b>	<i>0.89 (0.66)</i>	<b>2.56 (1.21)</b>	+	+	+	
4	85.14 (1.27)	-0.09 (0.44)	-0.03 (0.10)	<i>0.57 (0.89)</i>	<b>2.45 (0.78)</b>	+	+	+	
5	88.00 (0.78)	0.03 (0.24)	0.55 (0.82)	<i>0.62 (0.93)</i>	<b>3.52 (1.66)</b>	+	+	+	
6	90.88 (0.88)	-0.08 (0.17)	0.06 (0.25)	<i>0.71 (0.57)</i>	<b>1.56 (1.08)</b>	+	+	+	

**Table 3.** Ranking accuracies of different predictor combination algorithms on the *aPascal* and *Zap50K* datasets. For each dataset, we repeated experiments 10 times with different training, validation, and test set splits. For baseline  $f^0$  (second column), Kendall’s Tau correlations $\times 100$  (standard deviations in parentheses) are presented. For the remaining algorithms (third to sixth columns), the accuracy offsets from  $f^0$  are presented. The best and second best results are highlighted with **bold** and *italic* fonts, respectively. The results of statistical significance test based on a t-test with  $\alpha=0.95$  are highlighted in **green** (significantly positive) and **orange** (significantly negative). The last three columns show the results of statistical significance test of our algorithm with *GL*, *OPC*, and *MTL*, respectively (+/-: significantly positive/negative).

<i>aPascal</i>									
Attr.	Baseline $f^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>	
1	59.44 (4.19)	3.11 (2.88)	<i>3.38 (2.31)</i>	0.87 (0.64)	<b>6.04 (3.43)</b>	+	+	+	
2	68.21 (4.50)	-0.02 (0.12)	<i>0.42 (0.87)</i>	0.27 (0.13)	<b>1.19 (2.17)</b>	0	0	0	
3	14.45 (4.62)	1.15 (2.40)	<i>2.09 (5.79)</i>	0.69 (0.33)	<b>4.24 (4.67)</b>	0	0	+	
4	65.20 (3.20)	0.23 (0.57)	<b>1.82 (1.38)</b>	0.07 (0.04)	<i>1.19 (2.67)</i>	0	0	0	
5	57.87 (4.38)	2.08 (1.69)	<i>4.10 (2.03)</i>	1.19 (0.67)	<b>4.10 (2.26)</b>	+	0	+	
6	57.37 (5.76)	3.21 (1.49)	<b>4.18 (1.23)</b>	1.36 (0.70)	<i>3.90 (2.28)</i>	0	0	+	
7	71.34 (2.69)	0.05 (0.36)	<i>2.29 (2.27)</i>	1.05 (0.53)	<b>2.91 (1.85)</b>	+	0	+	
8	67.08 (4.56)	2.25 (1.79)	<i>3.78 (3.30)</i>	1.65 (0.51)	<b>3.79 (3.11)</b>	+	0	+	
9	62.36 (4.83)	2.54 (1.51)	<i>3.91 (2.27)</i>	2.03 (0.55)	<b>4.69 (1.34)</b>	+	0	+	
10	57.09 (4.36)	4.79 (2.25)	<i>5.38 (2.86)</i>	2.92 (1.55)	<b>7.00 (2.74)</b>	+	+	+	
11	62.25 (3.59)	2.05 (1.54)	<i>2.61 (1.39)</i>	1.76 (0.75)	<b>4.12 (3.01)</b>	0	0	+	
12	60.58 (4.76)	<i>3.48 (2.13)</i>	<b>4.25 (2.38)</b>	1.54 (0.67)	3.43 (2.63)	0	0	+	
13	46.71 (4.95)	3.41 (3.54)	<i>4.20 (4.16)</i>	1.47 (0.60)	<b>4.86 (3.99)</b>	0	0	+	
14	52.31 (3.52)	3.52 (2.89)	<i>4.37 (1.68)</i>	2.29 (1.07)	<b>6.75 (2.27)</b>	+	+	+	
15	52.07 (6.42)	2.14 (2.77)	<i>2.34 (2.77)</i>	1.46 (0.77)	<b>5.65 (3.71)</b>	+	+	+	
16	47.33 (4.15)	0.46 (0.72)	<i>1.55 (1.22)</i>	1.16 (0.41)	<b>2.76 (2.55)</b>	+	0	0	
17	49.53 (4.57)	-0.18 (0.96)	0.42 (1.70)	<i>0.95 (0.45)</i>	<b>2.49 (2.44)</b>	+	0	0	
18	82.22 (3.20)	0.60 (0.81)	<i>1.37 (0.70)</i>	1.10 (0.36)	<b>1.53 (0.88)</b>	+	0	0	
19	70.22 (3.21)	1.04 (1.05)	1.01 (1.37)	<i>1.42 (0.61)</i>	<b>2.40 (2.19)</b>	0	0	0	
20	79.62 (2.60)	1.87 (0.88)	<b>2.10 (1.19)</b>	1.44 (0.58)	<i>2.09 (2.31)</i>	0	0	0	
21	53.01 (5.18)	<i>0.69 (1.34)</i>	0.02 (0.37)	0.37 (0.25)	<b>3.70 (4.07)</b>	+	+	+	
22	55.52 (3.98)	2.83 (2.30)	<i>3.04 (1.78)</i>	2.22 (0.71)	<b>6.41 (3.92)</b>	+	+	+	
23	70.53 (5.23)	2.31 (2.23)	2.03 (1.14)	<i>2.34 (1.11)</i>	<b>3.82 (2.24)</b>	+	+	+	
24	38.48 (5.51)	<i>1.12 (2.63)</i>	0.95 (3.22)	0.91 (0.43)	<b>1.81 (4.71)</b>	0	0	0	
25	48.45 (3.52)	<i>0.89 (2.29)</i>	0.76 (1.65)	0.58 (0.51)	<b>2.48 (2.88)</b>	0	+	+	
26	49.93 (3.91)	2.53 (2.56)	<i>3.60 (3.00)</i>	2.31 (0.55)	<b>6.00 (2.57)</b>	+	+	+	
27	72.87 (3.05)	0.04 (0.31)	<i>1.07 (1.49)</i>	0.91 (0.29)	<b>1.65 (1.55)</b>	+	0	0	
28	64.35 (2.27)	0.30 (0.66)	<b>1.39 (1.33)</b>	0.42 (0.42)	<i>0.43 (3.09)</i>	0	0	0	
29	53.84 (3.47)	3.97 (2.36)	<i>4.16 (2.19)</i>	2.53 (0.78)	<b>5.17 (2.20)</b>	0	0	+	

<i>Zap50K</i>									
Attr.	Baseline $f^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>	
1	87.97 (0.99)	-0.00 (0.38)	<b>0.27 (0.60)</b>	-0.20 (0.67)	<i>0.27 (0.75)</i>	0	0	0	
2	89.43 (1.56)	<i>0.13 (0.82)</i>	-0.27 (1.10)	<b>0.40 (0.86)</b>	0.03 (0.95)	0	0	0	
3	90.67 (1.29)	0.43 (0.80)	<b>0.80 (0.83)</b>	<i>0.70 (0.84)</i>	0.67 (1.23)	0	0	0	
4	90.33 (1.56)	<i>0.23 (1.14)</i>	0.17 (0.98)	0.03 (0.82)	<b>0.37 (0.87)</b>	0	0	0	

**Table 4.** Ranking accuracies of different predictor combination algorithms on the *CUB* dataset. We repeated experiments 10 times with different training, validation, and test set splits. For baseline  $\mathbf{f}^0$  (second column), Kendall’s Tau correlations $\times 100$  (standard deviations in parentheses) are presented. For the remaining algorithms (third to sixth columns), the accuracy offsets from  $\mathbf{f}^0$  are presented. The best and second best results are highlighted with **bold** and *italic* fonts, respectively. The results of statistical significance test based on a t-test with  $\alpha=0.95$  are highlighted in **green** (significantly positive) and **orange** (significantly negative). The last three columns show the results of statistical significance test of our algorithm with *GL*, *OPC*, and *MTL*, respectively (+/-: significantly positive/negative).

Attr.	Baseline $\mathbf{f}^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>
1	68.80 (3.98)	-0.02 (0.07)	-0.07 (0.28)	<i>0.16 (0.24)</i>	<b>1.47 (1.05)</b>	+	+	+
2	74.83 (3.89)	<b>0.64 (0.66)</b>	<i>1.12 (0.64)</i>	0.87 (1.53)	<b>2.00 (1.10)</b>	+	+	+
3	78.59 (2.36)	<b>0.86 (0.66)</b>	<b>1.49 (1.10)</b>	<i>1.93 (1.23)</i>	<b>2.22 (1.10)</b>	+	+	0
4	73.92 (2.51)	-0.17 (0.39)	0.27 (0.42)	<b>1.82 (1.72)</b>	<i>1.38 (2.03)</i>	+	0	0
5	74.61 (3.37)	<b>1.30 (1.35)</b>	<b>0.98 (1.35)</b>	<i>2.36 (1.46)</i>	<b>2.73 (2.01)</b>	+	+	0
6	63.86 (5.24)	-0.00 (0.41)	0.58 (0.91)	<b>1.48 (1.99)</b>	<i>0.89 (1.34)</i>	0	0	0
7	76.97 (2.21)	<b>1.02 (0.70)</b>	<i>0.54 (0.34)</i>	<b>1.18 (0.66)</b>	<i>1.06 (0.79)</i>	0	+	0
8	62.97 (3.05)	-0.00 (0.04)	<i>0.26 (0.45)</i>	0.20 (0.49)	<b>0.76 (1.11)</b>	0	0	0
9	72.52 (2.57)	<b>1.05 (0.99)</b>	<b>1.06 (0.57)</b>	<i>1.53 (1.19)</i>	<b>3.08 (1.84)</b>	+	+	+
10	63.62 (2.99)	0.09 (0.35)	0.50 (1.69)	<i>0.75 (0.86)</i>	<b>2.30 (1.47)</b>	+	0	+
11	59.70 (3.69)	0.02 (0.30)	0.02 (0.29)	<b>0.66 (0.98)</b>	<i>0.54 (1.30)</i>	0	0	0
12	71.08 (2.09)	0.17 (0.40)	-0.04 (0.74)	<b>0.87 (0.95)</b>	<i>0.87 (0.72)</i>	+	+	0
13	78.10 (2.31)	<b>0.25 (0.31)</b>	0.11 (0.41)	<b>1.88 (1.09)</b>	<i>1.31 (1.27)</i>	+	+	-
14	74.13 (1.90)	<b>1.08 (0.89)</b>	<b>0.48 (0.39)</b>	<i>1.59 (1.32)</i>	<b>1.85 (1.93)</b>	0	+	0
15	72.23 (3.07)	0.02 (0.62)	0.04 (0.30)	<i>1.25 (0.92)</i>	<b>1.84 (1.28)</b>	+	+	+
16	73.32 (1.97)	<b>1.02 (1.13)</b>	<b>0.67 (0.67)</b>	<b>1.65 (1.81)</b>	<i>1.54 (2.30)</i>	0	0	0
17	58.11 (4.61)	0.09 (0.15)	0.18 (0.39)	<i>0.51 (0.61)</i>	<b>1.16 (0.95)</b>	+	+	+
18	57.35 (4.92)	-0.04 (0.23)	0.29 (0.65)	<i>0.55 (0.92)</i>	<b>0.71 (1.85)</b>	0	0	0
19	76.67 (3.06)	<b>1.28 (1.15)</b>	0.43 (0.98)	<i>1.30 (0.94)</i>	<b>1.92 (1.47)</b>	0	+	0
20	76.31 (2.10)	0.28 (0.47)	-0.06 (0.48)	<b>0.81 (1.37)</b>	<i>0.72 (1.53)</i>	0	0	0
21	75.45 (3.03)	<b>1.21 (0.93)</b>	<b>1.35 (1.54)</b>	<i>1.88 (1.55)</i>	<b>2.28 (1.50)</b>	0	0	0
22	75.28 (4.13)	<b>1.01 (0.64)</b>	0.49 (0.84)	<b>1.66 (1.08)</b>	<i>1.30 (1.45)</i>	0	+	0
23	69.67 (3.00)	0.12 (0.43)	0.02 (0.72)	<i>1.99 (1.36)</i>	<b>2.07 (1.74)</b>	+	+	0
24	76.24 (2.55)	<b>0.87 (0.57)</b>	<i>1.17 (0.87)</i>	0.77 (0.74)	<b>2.16 (1.49)</b>	+	0	+
25	70.57 (1.93)	-0.04 (0.08)	0.43 (0.76)	<i>0.94 (1.60)</i>	<b>1.21 (1.20)</b>	+	+	0
26	63.59 (2.97)	-0.09 (0.19)	0.07 (0.42)	<i>0.52 (1.15)</i>	<b>0.82 (0.53)</b>	+	+	0
27	72.15 (3.19)	<i>0.39 (0.56)</i>	0.19 (0.49)	-0.06 (0.41)	<b>0.41 (0.49)</b>	0	0	+
28	64.40 (3.29)	0.00 (0.54)	0.29 (0.92)	<i>1.21 (0.97)</i>	<b>1.27 (1.27)</b>	+	+	0
29	57.52 (2.78)	0.09 (0.24)	<b>0.76 (0.96)</b>	<i>2.20 (2.42)</i>	<b>2.50 (2.58)</b>	+	+	0
30	56.73 (2.86)	0.41 (0.72)	<b>0.78 (0.93)</b>	<b>1.98 (1.72)</b>	<i>1.89 (2.30)</i>	0	0	0
31	73.27 (2.96)	<b>1.13 (1.11)</b>	<b>1.14 (0.67)</b>	<i>1.91 (0.97)</i>	<b>2.00 (1.09)</b>	+	+	0
32	<i>52.82 (3.71)</i>	-0.12 (0.53)	-0.12 (0.55)	-0.00 (0.47)	<b>0.64 (0.81)</b>	+	+	+
33	69.13 (2.63)	0.14 (0.34)	0.21 (0.54)	<b>0.52 (0.66)</b>	<i>0.27 (0.98)</i>	0	0	0
34	58.15 (4.74)	-0.07 (0.48)	0.06 (0.44)	<i>0.46 (0.55)</i>	<b>0.90 (1.10)</b>	+	+	0
35	58.61 (3.93)	<i>0.52 (1.35)</i>	0.48 (1.58)	-0.10 (1.16)	<b>1.64 (1.26)</b>	0	+	+
36	58.91 (3.40)	0.35 (0.59)	0.15 (0.79)	<i>0.78 (1.84)</i>	<b>1.73 (1.16)</b>	+	+	0
37	52.60 (4.12)	0.14 (0.42)	0.08 (0.26)	<i>0.27 (0.38)</i>	<b>0.67 (1.01)</b>	0	0	0
38	67.73 (3.67)	<b>2.08 (2.28)</b>	<b>1.64 (0.93)</b>	<b>3.37 (2.68)</b>	<i>3.34 (2.47)</i>	+	+	0
39	76.39 (2.22)	-0.25 (1.44)	0.29 (0.41)	<i>0.67 (0.72)</i>	<b>1.19 (0.92)</b>	0	+	+
40	70.58 (3.12)	-0.04 (0.28)	0.09 (0.38)	<i>1.06 (0.64)</i>	<b>1.95 (0.80)</b>	+	+	+

**Table 5.** Ranking accuracies of different predictor combination algorithms on the first 40 attributes of *AWA2* dataset. We repeated experiments 10 times with different training, validation, and test set splits. For baseline  $\mathbf{f}^0$  (second column), Kendall’s Tau correlations $\times 100$  (standard deviations in parentheses) are presented. For the remaining algorithms (third to sixth columns), the accuracy offsets from  $\mathbf{f}^0$  are presented. The best and second best results are highlighted with **bold** and *italic* fonts, respectively. The results of statistical significance test based on a t-test with  $\alpha=0.95$  are highlighted in **green** (significantly positive) and **orange** (significantly negative). The last three columns show the results of statistical significance test of our algorithm with *GL*, *OPC*, and *MTL*, respectively (+/-: significantly positive/negative).

Attr.	Baseline $\mathbf{f}^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>
1	77.86 (3.70)	<i>0.45 (0.92)</i>	0.12 (0.26)	0.15 (0.29)	<b>7.25 (2.73)</b>	+	+	+
2	83.79 (3.18)	0.05 (0.11)	-0.17 (0.41)	<i>0.33 (0.60)</i>	<b>6.22 (2.14)</b>	+	+	+
3	98.55 (0.65)	<i>0.02 (0.11)</i>	<b>0.04 (0.07)</b>	-0.01 (0.06)	0.02 (0.46)	0	0	0
4	88.21 (3.47)	0.03 (0.32)	<i>0.13 (0.29)</i>	<i>0.05 (0.05)</i>	<b>5.22 (2.27)</b>	+	+	+
5	88.53 (1.90)	<i>0.21 (0.49)</i>	0.12 (0.21)	0.04 (0.29)	<b>3.55 (2.17)</b>	+	+	+
6	<i>97.94 (1.07)</i>	-0.07 (0.17)	-0.02 (0.07)	-0.12 (0.29)	<b>0.69 (0.65)</b>	+	+	+
7	<i>99.22 (0.34)</i>	-0.05 (0.11)	-0.02 (0.06)	-0.00 (0.09)	<b>0.24 (0.23)</b>	+	+	+
8	82.30 (1.69)	-0.03 (0.12)	0.12 (0.39)	<i>0.13 (0.21)</i>	<b>4.32 (1.88)</b>	+	+	+
9	79.33 (4.37)	-0.01 (0.35)	<i>0.17 (0.42)</i>	0.04 (0.14)	<b>7.05 (1.65)</b>	+	+	+
10	98.58 (0.85)	<i>0.08 (0.31)</i>	0.03 (0.13)	0.01 (0.03)	<b>0.26 (0.40)</b>	0	0	0
11	97.44 (0.97)	-0.03 (0.23)	0.01 (0.10)	<i>0.06 (0.17)</i>	<b>0.90 (0.35)</b>	+	+	+
12	94.46 (1.91)	0.00 (0.25)	0.04 (0.31)	<i>0.47 (0.52)</i>	<b>1.13 (0.70)</b>	+	+	+
13	93.52 (1.05)	-0.14 (0.24)	<i>0.08 (0.31)</i>	0.05 (0.19)	<b>2.43 (0.68)</b>	+	+	+
14	94.50 (1.64)	0.04 (0.15)	0.33 (0.49)	<i>0.55 (0.43)</i>	<b>2.29 (0.64)</b>	+	+	+
15	95.04 (1.21)	0.00 (0.07)	<i>0.25 (0.40)</i>	0.13 (0.19)	<b>1.70 (0.68)</b>	+	+	+
16	85.91 (2.85)	-0.02 (0.03)	0.05 (0.23)	<i>1.39 (1.03)</i>	<b>5.88 (2.32)</b>	+	+	+
17	87.00 (2.60)	<i>0.48 (0.64)</i>	0.01 (0.36)	<i>1.64 (0.76)</i>	<b>4.33 (1.84)</b>	+	+	+
18	99.25 (0.55)	0.01 (0.03)	<i>0.08 (0.18)</i>	0.02 (0.10)	<b>0.18 (0.26)</b>	0	0	0
19	<b>99.75 (0.22)</b>	<i>-0.00 (0.01)</i>	-0.02 (0.04)	-0.02 (0.05)	-0.05 (0.31)	0	0	0
20	97.88 (0.99)	0.03 (0.15)	0.07 (0.26)	<i>0.27 (0.40)</i>	<b>0.79 (0.59)</b>	+	+	+
21	92.36 (2.33)	0.08 (0.13)	<i>0.30 (0.62)</i>	0.13 (0.37)	<b>2.44 (1.35)</b>	+	+	+
22	96.81 (1.03)	0.01 (0.06)	<i>0.20 (0.35)</i>	<i>0.18 (0.18)</i>	<b>1.40 (0.77)</b>	+	+	+
23	91.60 (2.67)	0.11 (0.36)	<i>0.15 (0.23)</i>	0.07 (0.14)	<b>3.57 (0.88)</b>	+	+	+
24	95.46 (1.41)	0.02 (0.30)	<i>0.06 (0.23)</i>	0.01 (0.14)	<b>1.43 (1.24)</b>	+	+	+
25	89.94 (2.52)	-0.02 (0.08)	0.07 (0.43)	<i>0.14 (0.45)</i>	<b>3.55 (1.68)</b>	+	+	+
26	84.78 (3.62)	0.03 (0.28)	<i>0.61 (1.16)</i>	0.23 (0.34)	<b>3.77 (1.15)</b>	+	+	+
27	90.67 (2.26)	-0.27 (1.30)	<i>0.92 (1.09)</i>	<i>1.09 (1.05)</i>	<b>3.91 (1.80)</b>	+	+	+
28	90.67 (1.67)	0.00 (0.03)	<i>0.05 (0.23)</i>	-0.00 (0.10)	<b>4.51 (1.59)</b>	+	+	+
29	95.73 (1.49)	<b>0.16 (0.30)</b>	-0.04 (0.16)	-0.11 (0.34)	<i>0.15 (0.80)</i>	0	0	0
30	94.97 (1.58)	0.04 (0.11)	<i>0.12 (0.42)</i>	0.06 (0.17)	<b>1.21 (1.10)</b>	+	+	+
31	94.39 (2.23)	<i>0.41 (0.45)</i>	<b>0.40 (0.27)</b>	<i>0.41 (0.48)</i>	<b>2.52 (0.99)</b>	+	+	+
32	96.92 (1.17)	-0.02 (0.13)	<i>0.14 (0.33)</i>	0.04 (0.11)	<b>0.86 (0.86)</b>	+	0	+
33	85.97 (4.63)	0.08 (0.21)	<i>0.16 (0.53)</i>	0.02 (0.20)	<b>6.68 (3.30)</b>	+	+	+
34	98.58 (0.83)	0.00 (0.01)	<i>0.07 (0.14)</i>	-0.07 (0.12)	<b>0.25 (0.60)</b>	0	0	0
35	98.69 (0.47)	0.03 (0.07)	<i>0.30 (0.34)</i>	<i>0.26 (0.28)</i>	<b>0.45 (0.45)</b>	+	0	0
36	<i>97.71 (0.82)</i>	-0.01 (0.14)	-0.02 (0.15)	-0.07 (0.13)	<b>0.44 (0.27)</b>	+	+	+
37	97.19 (1.48)	-0.04 (0.12)	0.22 (0.59)	<i>0.34 (0.42)</i>	<b>1.17 (1.18)</b>	+	+	+
38	95.21 (0.97)	<i>0.05 (0.10)</i>	-0.00 (0.16)	0.04 (0.17)	<b>1.95 (0.75)</b>	+	+	+
39	89.66 (1.83)	-0.02 (0.05)	0.13 (0.75)	<i>1.17 (0.94)</i>	<b>4.01 (1.53)</b>	+	+	+
40	95.33 (1.58)	-0.04 (0.17)	<b>0.20 (0.26)</b>	<i>0.64 (0.45)</i>	<b>2.60 (1.27)</b>	+	+	+

**Table 6.** Ranking accuracies of different predictor combination algorithms on the last 40 attributes of *AWA2* dataset. For each dataset, we repeated experiments 10 times with different training, validation, and test set splits. For baseline  $\mathbf{f}^0$  (second column), Kendall’s Tau correlations $\times 100$  (standard deviations in parentheses) are presented. For the remaining algorithms (third to sixth columns), the accuracy offsets from  $\mathbf{f}^0$  are presented. The best and second best results are highlighted with **bold** and *italic* fonts, respectively. The results of statistical significance test based on a t-test with  $\alpha=0.95$  are highlighted in **green** (significantly positive) and **orange** (significantly negative). The last three columns show the results of statistical significance test of our algorithm with *GL*, *OPC*, and *MTL*, respectively (+/-: significantly positive/negative).

Attr.	Baseline $\mathbf{f}^0$	<i>GL</i>	<i>OPC</i>	<i>MTL</i>	<i>NPC (ours)</i>	vs. <i>GL</i>	vs. <i>OPC</i>	vs. <i>MTL</i>
41	95.21 (1.26)	-0.06 (0.16)	0.01 (0.13)	<i>0.07 (0.19)</i>	<b>1.03 (0.60)</b>	+	+	+
42	84.32 (3.88)	<i>0.21 (0.49)</i>	-0.06 (0.17)	-0.07 (0.27)	<b>2.86 (1.72)</b>	+	+	+
43	95.78 (1.49)	0.06 (0.12)	-0.10 (0.19)	<i>0.13 (0.20)</i>	<b>2.57 (1.12)</b>	+	+	+
44	<i>98.22 (0.78)</i>	-0.05 (0.25)	-0.00 (0.04)	-0.00 (0.17)	<b>0.62 (0.50)</b>	+	+	+
45	88.71 (2.33)	0.07 (0.20)	0.42 (0.65)	<i>0.65 (0.83)</i>	<b>2.92 (2.09)</b>	+	+	+
46	83.63 (2.01)	0.01 (0.04)	0.43 (0.74)	<i>0.77 (0.41)</i>	<b>8.15 (1.08)</b>	+	+	+
47	89.57 (2.34)	<i>0.39 (0.59)</i>	0.16 (0.33)	0.13 (0.34)	<b>3.66 (2.15)</b>	+	+	+
48	91.92 (2.16)	<b>0.18 (0.24)</b>	0.51 (0.83)	<i>0.51 (0.43)</i>	<b>2.70 (1.81)</b>	+	+	+
49	89.56 (3.03)	0.05 (0.27)	-0.03 (0.25)	<i>0.34 (0.40)</i>	<b>3.65 (1.79)</b>	+	+	+
50	95.20 (1.24)	-0.01 (0.05)	-0.06 (0.25)	<i>0.11 (0.21)</i>	<b>2.23 (0.82)</b>	+	+	+
51	90.78 (2.72)	-0.34 (1.72)	<i>1.01 (1.61)</i>	0.91 (0.81)	<b>3.16 (1.46)</b>	+	+	+
52	94.00 (1.48)	-0.14 (0.19)	0.08 (0.70)	<i>0.28 (0.29)</i>	<b>3.11 (0.97)</b>	+	+	+
53	94.51 (1.71)	<i>0.10 (0.17)</i>	-0.00 (0.17)	-0.02 (0.18)	<b>1.25 (0.78)</b>	+	+	+
54	91.17 (1.73)	0.15 (0.39)	<i>0.22 (0.51)</i>	0.16 (0.47)	<b>4.19 (1.20)</b>	+	+	+
55	93.07 (1.87)	0.07 (0.29)	-0.00 (0.38)	<i>0.30 (0.33)</i>	<b>2.95 (0.86)</b>	+	+	+
56	88.30 (3.12)	0.29 (0.51)	0.57 (1.21)	<i>1.26 (1.36)</i>	<b>3.32 (2.47)</b>	+	+	+
57	92.34 (2.15)	<i>0.34 (0.61)</i>	0.03 (0.30)	0.05 (0.18)	<b>2.59 (0.92)</b>	+	+	+
58	96.10 (1.78)	<b>0.16 (0.18)</b>	0.11 (0.17)	<i>0.43 (0.75)</i>	<b>1.13 (1.02)</b>	+	+	+
59	92.12 (3.24)	<i>0.30 (0.42)</i>	0.16 (0.25)	<b>0.23 (0.23)</b>	<b>3.07 (1.47)</b>	+	+	+
60	89.37 (1.95)	0.12 (0.25)	<i>0.44 (0.75)</i>	0.14 (0.23)	<b>2.37 (1.45)</b>	+	+	+
61	92.43 (1.97)	-0.10 (0.22)	0.04 (0.27)	<i>0.17 (0.28)</i>	<b>2.37 (0.88)</b>	+	+	+
62	98.16 (1.08)	0.08 (0.27)	0.07 (0.24)	<i>0.18 (0.23)</i>	<b>0.39 (0.53)</b>	0	0	0
63	90.95 (3.52)	-0.04 (0.10)	<i>0.33 (0.51)</i>	0.15 (0.49)	<b>3.91 (2.55)</b>	+	+	+
64	92.31 (2.46)	-0.06 (0.13)	<b>0.19 (0.24)</b>	<i>0.39 (0.52)</i>	<b>2.80 (1.80)</b>	+	+	+
65	90.41 (3.04)	0.14 (0.30)	-0.01 (0.29)	<i>0.57 (0.49)</i>	<b>3.53 (1.25)</b>	+	+	+
66	89.81 (2.71)	<i>0.80 (0.89)</i>	0.58 (1.30)	0.00 (0.27)	<b>3.12 (1.38)</b>	+	+	+
67	94.72 (2.26)	-0.02 (0.13)	<i>0.11 (0.77)</i>	0.07 (0.15)	<b>2.38 (0.83)</b>	+	+	+
68	85.60 (2.43)	0.18 (0.83)	<i>0.47 (1.06)</i>	0.11 (0.14)	<b>6.41 (1.36)</b>	+	+	+
69	98.79 (0.41)	0.02 (0.05)	0.14 (0.27)	<i>0.14 (0.22)</i>	<b>0.50 (0.41)</b>	+	+	+
70	97.09 (1.19)	-0.04 (0.12)	<b>0.14 (0.18)</b>	<i>0.27 (0.35)</i>	<b>0.93 (0.83)</b>	+	+	+
71	99.02 (0.47)	0.01 (0.02)	<i>0.26 (0.20)</i>	0.16 (0.11)	<b>0.49 (0.32)</b>	+	+	+
72	96.69 (0.57)	-0.11 (0.29)	<i>0.23 (0.44)</i>	0.18 (0.19)	<b>1.79 (0.72)</b>	+	+	+
73	94.18 (1.34)	<i>0.09 (0.28)</i>	-0.00 (0.06)	0.00 (0.13)	<b>2.48 (0.98)</b>	+	+	+
74	83.62 (3.96)	-0.02 (0.41)	0.16 (0.84)	<i>0.56 (0.62)</i>	<b>5.06 (3.19)</b>	+	+	+
75	82.00 (3.99)	-0.06 (0.11)	<b>0.54 (0.64)</b>	<i>1.43 (1.62)</i>	<b>6.32 (2.62)</b>	+	+	+
76	83.13 (3.21)	0.08 (0.45)	0.23 (0.64)	<i>0.29 (0.65)</i>	<b>6.38 (2.76)</b>	+	+	+
77	85.16 (1.97)	0.09 (0.29)	0.08 (0.40)	<i>0.76 (0.67)</i>	<b>5.58 (2.22)</b>	+	+	+
78	85.83 (2.74)	0.14 (0.48)	-0.11 (0.50)	<i>0.43 (0.81)</i>	<b>4.48 (2.19)</b>	+	+	+
79	87.37 (2.77)	0.08 (0.34)	0.05 (0.29)	<i>0.50 (0.30)</i>	<b>5.25 (1.86)</b>	+	+	+
80	85.30 (2.93)	0.04 (0.21)	0.18 (0.53)	<i>0.44 (0.30)</i>	<b>6.05 (2.39)</b>	+	+	+