

442 **A General Form of Selection Tensor Decomposition**

443 In this section, we further extend the selection tensor decomposition in Section 3.3.1 from a special
 444 case, where $t = 2$, to a more general case, where $2 \leq t \leq n$. The interaction selection tensor \mathbf{A}_v^t for
 445 t -th order features is also semi-positive and symmetric. By extending the Takagi factorization [2], we
 446 have:

$$\mathbf{A}_v^t = \Sigma \times_1 \mathcal{U} \times_2 \mathcal{U} \times_3 \cdots \times_t \mathcal{U}, \quad (13)$$

447 where Σ is a $\underbrace{d' \times \cdots \times d'}_{t \text{ times}}$ diagonal tensor, \times_t denotes the t -mode matrix multiplication [4],

448 $\mathcal{U} \in \mathbb{R}^{m \times d'}$ and $d' < m$. Similar to the special case where $t = 2$, we adopt multi-mode ten-
 449 sor factorization [4] to replace \mathcal{U} as an output of a neural network, denoted as:

$$\mathcal{U} \approx f_{\hat{\theta}}(\hat{\mathbf{E}}), \quad (14)$$

450 where $f_{\hat{\theta}} : \mathbb{R}^{m \times \hat{d}} \rightarrow \mathbb{R}^{m \times d'}$, $\hat{d} \ll d'$ is a neural network with parameter $\hat{\theta}$ and $\hat{\mathbf{E}} \in \mathbb{R}^{m \times \hat{d}}$ is
 451 an additional embedding table for generating feature interaction selection tensor. The element of
 452 architecture metric $\mathbf{A}_{v(k_{i_1}, \dots, k_{i_t})}^t$ can be calculated given the following equation:

$$\mathbf{A}_{v(k_{i_1}, \dots, k_{i_t})}^t = \Sigma \times_1 f_{\hat{\theta}}(\hat{\mathbf{E}}_{k_{i_1}, :}) \times_2 \cdots \times_t f_{\hat{\theta}}^T(\hat{\mathbf{E}}_{k_{i_t}, :}). \quad (15)$$

453 The original value-grained selection tensor \mathbf{A}_v^t consists of $\mathcal{O}(m^t)$ elements. The trainable elements
 454 is reduced to $\mathcal{O}(md')$ after the Takagi factorization [2] and to $\mathcal{O}(\hat{d}(m + d'))$ after the multi-mode
 455 tensor factorization [4].

456 **B Experiment Setup**

457 **B.1 Dataset and Preprocessing**

458 We conduct our experiments on two public real-world benchmark datasets. The statistics of all
 459 datasets are given in Table 3. We describe all these datasets and the pre-processing steps below.

Table 3: Dataset Statistics

Dataset	#samples	#field	#value	pos ratio
Criteo	4.6×10^7	39	6.8×10^6	0.2562
Avazu	4.0×10^7	24	4.4×10^6	0.1698
KDD12	1.5×10^8	11	6.0×10^6	0.0445

Note: *#samples* refers to the total samples in the dataset, *#field* refers to the number of feature fields for original features, *#value* refers to the number of feature values for original features, *pos ratio* refers to the positive ratio.

460 **Criteo** dataset consists of ad click data over a week. It consists of 26 categorical feature fields and
 461 13 numerical feature fields. Following the best practice [26], we discretize each numeric value x
 462 to $\lfloor \log^2(x) \rfloor$, if $x > 2$; $x = 1$ otherwise. We replace infrequent categorical features with a default
 463 "OOV" (i.e. out-of-vocabulary) token, with *min_count*=2.

464 **Avazu** dataset contains 10 days of click logs. It has 24 fields with categorical features. Following the
 465 best practice [26], we remove the *instance_id* field and transform the *timestamp* field into three new
 466 fields: *hour*, *weekday* and *is_weekend*. We replace infrequent categorical features with the "OOV"
 467 token, with *min_count*=2.

468 **KDD12** dataset contains training instances derived from search session logs. It has 11 categorical
 469 fields, and the click field is the number of times the user clicks the ad. We replace infrequent features
 470 with an "OOV" token, with *min_count*=10.

471 **B.2 Parameter Setup**

472 To ensure the reproducibility of experimental results, here we further introduce the implementation
 473 setting in details. We implement our methods using PyTorch. We adopt the Adam optimizer with a

474 mini-batch size of 4096. We set the embedding sizes to 16 in all the models. We set the predictor
 475 as an MLP model with [1024, 512, 256] for all methods. All the hyper-parameters are tuned on
 476 the validation set with a learning rate from [1e-3, 3e-4, 1e-4, 3e-5, 1e-5] and weight decay from
 477 [1e-4, 3e-5, 1e-5, 3e-6, 1e-6]. We also tune the learning ratio for the feature interaction selection
 478 parameters from [1e-4, 3e-5, 1e-5, 3e-6, 1e-6] and while weight decay from [1e-4, 3e-5, 1e-5, 3e-6,
 479 1e-6, 0]. The initialization parameters for the retraining stage is selected from the best-performed
 480 model parameters and randomly initialized ones.

481 B.3 Hardware Platform

482 All experiments are conducted on a Linux server with one Nvidia-Tesla V100-PCIe-32GB GPU,
 483 128GB main memory and 8 Intel(R) Xeon(R) Gold 6140 CPU cores.

484 C Ablation Study

485 C.1 Feature Interaction Operation

486 In this section, we conduct an ablation study on the feature interaction operation, comparing the
 487 performance of the default setting, which uses the *inner product*, with the *outer product* operation.
 488 We evaluate these operations on *OptFeature* and its two variants: *OptFeature-f* and *OptFeature-v*.
 489 The results are summarized in Table 4.

Table 4: Performance Comparison over Different Feature Interaction Operation.

Dataset		Criteo		Avazu		KDD12	
Category	Model	AUC	Logloss	AUC	Logloss	AUC	Logloss
inner product	OptFeature-f	0.8115	0.4404	0.7920	0.3744	0.7978	0.1530
	OptFeature-v	0.8116	0.4403	0.7920	0.3742	0.7981	0.1529
	OptFeature	0.8116	0.4402	0.7925	0.3741	0.7982	0.1529
outer product	OptFeature-f	0.8114	0.4404	0.7896	0.3760	0.7957	0.1535
	OptFeature-v	0.8113	0.4405	0.7902	0.3752	0.7961	0.1533
	OptFeature	0.8115	0.4403	0.7899	0.3753	0.7961	0.1533

490 From the table, we observe that the *inner product* operation outperforms the *outer product* operation.
 491 This performance gap is particularly significant on the Avazu and KDD12 datasets, while it is
 492 relatively insignificant on the Criteo dataset. The drop in performance with the *outer product*
 493 operation is likely due to the introduction of a significantly larger number of inputs into the final
 494 predictor. This makes it more challenging for the predictor to effectively balance the information
 495 from raw inputs and feature interactions.

496 C.2 Dimension Selection

497 In this section, we perform an ablation study on the feature interaction
 498 selection dimension \hat{d} . We compare the AUC performance with
 499 the corresponding dimension \hat{d} and present the results in Figure 5.
 500 From the figure, we can observe that as the dimension \hat{d} increases,
 501 the AUC performance remains relatively consistent over the Criteo
 502 dataset. This suggests that it is relatively easy to distinguish value-
 503 level selection on the Criteo dataset. However, on the Avazu and
 504 KDD12 datasets, the AUC performance improves as the selection
 505 dimension \hat{d} increases. This indicates that distinguishing informative
 506 values is comparatively more challenging on these two datasets.

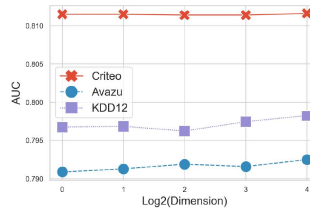


Figure 5: Ablation over feature interaction selection dimension on OptFeature.

507 C.3 Higher-Order Feature Interactions

508 In this section, we investigate the influence of higher-order feature
 509 interactions over the final results on the KDD12 dataset. We compare
 510 the default setting where only considering second-order interactions with two other settings: (i) only

511 third-order interactions and (ii) both second and third-order interactions. We visualize the result in
 512 Figure 6.

513 From the figure, we can draw the following ob-
 514 servations. First, only considering third-order
 515 interactions leads to the worst performance.
 516 This aligns with the common understanding that
 517 second-order interactions are typically consid-
 518 ered the most informative in deep sparse predic-
 519 tion [13]. Second, for field-level selection, the
 520 performance improves when both second and
 521 third-order interactions are incorporated into the
 522 model. This finding is consistent with previ-
 523 ous studies [10, 6], as the inclusion of additional
 524 interactions introduces extra information that en-
 525 hances the performance. In contrast, for value-
 526 level selection, the performance tends to decrease
 527 when both second and third-order interactions are
 528 included. This could be attributed to the fact that
 529 value-level selection operates at a finer-grained
 530 level and might be more challenging to optimize
 directly.

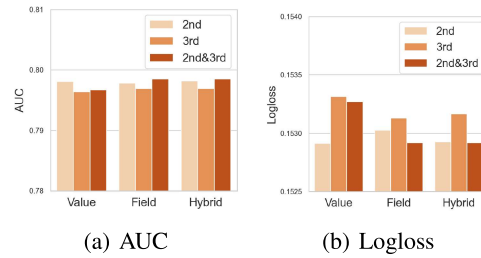


Figure 6: Performance Comparison over Different Feature Interaction Orders.

531 C.4 Selection Visualization

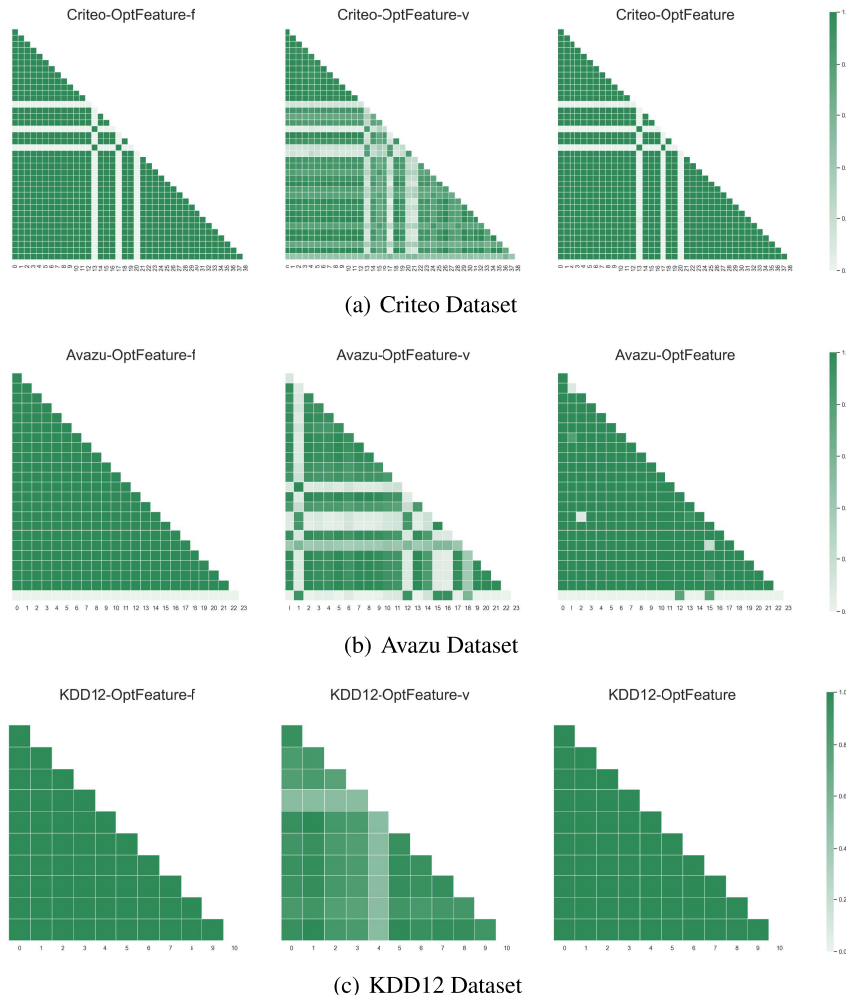


Figure 7: Visualization of the Feature Interaction Selection Results.

532 In this section, we present the visualization of the interaction selection results for OptFeature and its
533 two variants in Figure 7. OptFeature-f performs a binary selection for each interaction field, allowing
534 for easy visualization through a heatmap representation where one indicates keep and zero indicates
535 drop. On the other hand, OptFeature-v and OptFeature involve value-level interaction selection.
536 Hence, we visualize them by setting each element as the percentage of being selected over the training
537 set. The detailed equation for calculating the value for interaction field (i, j) is shown in Equation 16.

$$\mathbf{P}_{(i,j)} = \frac{\# \text{Samples keeping interaction field (i, j)}}{\# \text{Training Samples}} \quad (16)$$

538 From the visualization, we can observe that OptFeature acts as a hybrid approach, exhibiting a
539 combination of both field-level and value-level interactions. Interestingly, we note significant
540 differences between certain interaction fields in the KDD12 and Avazu datasets. OptFeature-f retains
541 all of its interactions, while OptFeature-v only keeps a proportion of the value-level interactions. This
542 observation further emphasizes the importance of exploring interactions at a finer-grained level.

543 **D Broader Impact**

544 Successfully identifying informative feature interactions could be a double-edged sword. One the
545 one hand, by proving that introducing noisy features into model could harm the performance, feature
546 interaction selection could be used as supporting evidences in preventing certain business, such as
547 advertisement recommendation, from over-collecting users' information, thereby protecting user
548 privacy. On the other hand, these tools, if acquired by malicious people, can be used for filtering out
549 potential victims, such as individuals susceptible to email fraud. As researchers, it is crucial for us to
550 remain vigilant and ensure that our work is directed towards noble causes and societal benefits.