A Hybrid Bandit Model with Visual Priors for Creative Ranking in Display Advertising

Shiyao Wang Alibaba Group Beijing, China shiyao.wsy@alibaba-inc.com Qi Liu*
University of Science and Technology
of China
Hefei, China
qiliu67@mail.ustc.edu.cn

Tiezheng Ge Alibaba Group Beijing, China tiezheng.gtz@alibaba-inc.com

Defu Lian
University of Science and Technology
of China
Hefei, China
liandefu@ustc.edu.cn

ABSTRACT

Creative plays a great important role in e-commerce for exhibiting products. Sellers usually create multiple creatives for comprehensive demonstrations, thus it is crucial to display the most appealing design to maximize the Click-Through Rate (CTR). For this purpose, modern recommender systems dynamically rank creatives when a product is proposed for a user. However, this task suffers more cold-start problem than conventional products recommendation since the user-click data is more scarce and creatives potentially change more frequently. In this paper, we propose a hybrid bandit model with visual priors which first makes predictions with a visual evaluation, and then naturally evolves to focus on the specialities through the hybrid bandit model. Our contributions are three-fold: 1) We present a visual-aware ranking model (called VAM) that incorporates a list-wise ranking loss for ordering the creatives according to the visual appearance. 2) Regarding visual evaluation as a prior, the hybrid bandit model (called HBM) is proposed to evolve consistently to make better posteriori estimations by taking more observations into consideration for online scenarios. 3) A first large-scale creative dataset, CreativeRanking¹, is constructed, which contains over 1.7M creatives of 500k products as well as their real impression and click data. Extensive experiments have also been conducted on both our dataset and public Mushroom dataset, demonstrating the effectiveness of the proposed method.

KEYWORDS

Hybrid Bandit Model, Visual Priors, Creative Ranking

ACM Reference Format:

Shiyao Wang, Qi Liu*, Tiezheng Ge, Defu Lian, and Zhiqiang Zhang. 2021. A Hybrid Bandit Model with Visual Priors for Creative Ranking in Display Advertising. In *Proceedings of the Web Conference 2021 (WWW '21), April*

WWW '21, April 19–23, 2021, Ljubljana, Slovenia © 2021 Copyright held by the owner/author(s). ACM ISBN 978-x-xxxx-xxxx-x/YY/MM. https://doi.org/xxx Zhiqiang Zhang Alibaba Group Beijing, China zhang.zhiqiang@alibaba-inc.com



Figure 1: Some examples of ad creatives. Each row presents creatives that display the product in multiple ways. The corresponding CTRs at the bottom row indicate the large CTR gap among creatives.

19–23, 2021, Ljubljana, Slovenia. ACM, New York, NY, USA, 11 pages. https://doi.org/xxx

1 INTRODUCTION

Online display advertising is a rapid growing business and has become an important source of revenue for Internet service providers. The advertisements are delivered to customers through various online channels, e.g. e-commerce platform. Image ads are the most widely used format since they are more compact, intuitive and comprehensible [8]. In Figure 1, each row composes several ad images that describe the same product for comprehensive demonstrations. These images are called creatives. Although the creatives represent the same product, they may have largely different CTRs due to the visual appearance. Thus it is crucial to display the most appealing design to attract the potentially interested customers and maximize the Click-Through Rate(CTR).

In order to explore the most appealing creative, all of the candidates should be displayed to customers. Meanwhile, to ensure the overall performance of advertising, we prefer to display the creative

^{*}This work was done when the author Qi Liu was at Alibaba Group for intern. Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

that has the highest predicted CTR so far. This procedure can be modeled as a typical multi-armed bandit problem (MAB). It not only focuses on maximizing cumulative rewards (clicks) but also balance the exploration-exploitation(E&E) trade-off within a limited exploration resource so that CTR are considered. Epsilon-greedy [12], Thompson sampling [25] and Upper Confidence Bounds (UCB) approaches [2] are widely used strategies to deal with the bandit problem. However, creatives potentially change more frequently than products, and most of them cannot have sufficient impression opportunities to get reliable CTRs throughout their lifetime. So the conventional bandit models may suffer from coldstart problem in the initial random exploration period, hurting the online performance extremely. One potential solution to this problem is incorporating visual prior knowledge to facilitate a better exploration. [3, 8, 9, 21] consider the visual features extracted by deep convolutional networks and make deterministic selections for product recommendation. These deep models are in a heavy computation and cannot be flexibly updated online. Besides, the deterministic and greedy strategy may result in suboptimal solution due to the lack of exploration. Consequently, how to combine both the expressive visual representations and flexible bandit model remains a challenging problem.

In this paper, we propose an elegant method which incorporates visual prior knowledge into bandit model for facilitating a better exploration. It is based on a framework called NeuralLinear [24]. They consider approximate bayesian neural networks in a Thompson sampling framework to utilize both the learning ability of neural networks and the posterior sampling method. By adopting this general framework, we first present a novel convolutional network with a list-wise ranking loss function to select the most attractive creative. The ranking loss concentrates on capturing the visual patterns related to attractiveness, and the learned representations are treated as contextual information for the bandit model. Second, in terms of the bandit model, we make two major improvements: 1) Instead of randomly setting a prior hyperparameter to candidate arms, we use the weights of neural network to initialize the bandit parameters that further enhance the performance in the cold-starting phase. 2) To fit the industrial-scale data, we extend the linear regression model of NeuralLinear to a hybrid model which adopts two individual parameters, i.e. product-wise ones and creative-specific ones. The two components are adaptively combined during the exploring period. Last but not the least, because the creative ranking is a novel problem, it lacks real-world data for further study and comparison. To this end, we contribute a large-scale creative dataset¹ from Alibaba display advertising platform that comprises more than 500k products and 1.7M ad creatives.

In summary, the contributions of this paper include:

- We present a visual-aware ranking model (called VAM) that is capable of evaluating new creatives according to the visual appearance
- Regarding the learned visual predictions as a prior, the improved hybrid bandit model (called HBM) is proposed to make better posteriori estimations by taking more observations into consideration.

- We construct a novel large-scale creative dataset named *CreativeRanking*¹. Extensive experiments have been conducted on both our dataset and public Mushroom dataset, demonstrating the effectiveness of the proposed method.

2 PRELIMINARIES AND RELATED WORK

2.1 Preliminaries

Problem Statement The problem statement is as follows. Given a product, the goal is to determine which creative is the most attractive one and should be displayed. Meanwhile, we need to estimate the uncertainty of the predictions so as to maximize cumulative reward in a long run.

In the online advertising system, when an ad is shown to a user by displaying a candidate creative, this scenario is counted as an impression. Suppose there are N products, denoted as $\{I^1, I^2, \cdots, I^n, \cdots, I^N\}$, and each product I^n composes a group of creatives, indicated as $\{C_1^n, C_2^n, \cdots, C_m^n, \cdots, C_M^n\}$. For product I^n , the objective is to find the creative that subjects to:

$$C^{n} = \underset{c \in \{C_{1}^{n}, C_{2}^{n}, \dots, C_{M}^{n}\}}{\operatorname{arg\,max}} CTR(c) \tag{1}$$

where $CTR(\cdot)$ denotes the CTR for a given creative. An empirical way to produce CTR is accumulating the current clicks and impressions, and produce the click ratio as:

$$\widehat{CTR}(C_m^n) = \frac{click(C_m^n)}{impression(C_m^n)}$$
 (2)

where $click(\cdot)$ and $impression(\cdot)$ indicate the click and impression number of the creative C_m^n . But it may suffer from insufficient impressions, especially for the cold-start creatives. Another way is to learn a prediction function $\mathcal{N}(\cdot)$ from all the historical data by considering the contextual information (i.e. the image content) as:

$$\widehat{CTR}(C_m^n) = \mathcal{N}(C_m^n) \tag{3}$$

where $\mathcal{N}(\cdot)$ takes the image content of creative C_m^n as input, and learns from the historical data. The collected sequential data can be represented as

$$\mathcal{D} = \{ (C_1, y_1), \cdots, (C_t, y_t), \cdots, (C_{|\mathcal{D}|}, y_{|\mathcal{D}|}) \}$$
 (4)

where $y_t \in \{0,1\}$ is the label denotes whether a click is received. We take both the statistical data and content information into consideration. Subsection 2.2 reviews some product recommendation methods that take visual content as auxiliary information, and subsection 2.3 introduces typical bandit models to estimate uncertainty. Both of above methods will be our strong baselines.

2.2 Visual-aware Recommendation Methods

CTR prediction of image ads is a core task of online display advertising systems. Due to the recent advances in computer vision, visual features are employed to further enhance the recommendation models [3, 6, 8, 9, 14, 20, 21, 31, 33]. [3, 9] quantitatively study the relationship between handcrafted visual features and creative online performance. Different from fixed handcrafted features, [6, 14, 33] apply "off-the-shelf" visual features extracted by deep convolutional neural network[29]. [8, 21, 31] extend these methods by training the CNNs in an end-to-end manner. [20] integrate the category information on top of the CNN embedding to help

 $^{^1{\}rm The~Data}$ and code are publicly available at <code>https://github.com/alimama-creative/A_Hybrid_Bandit_Model_with_Visual_Priors_for_Creative_Ranking.git</code>

(a)	ı		(b) Num products of the Top20 categories	(c) CTR of Poor CTR of Average CTR of Best
Properties	Statistics	80000	#products	Creatives performance Creatives 6% CTRs
Number of Products	500,827	60000	l.	
max candidates(arms)	11	40000	III	3% 2%
min candidates(arms)	3	20000	#categories	1% 0%
mean candidates(arms)	3.4			er e
		4	Age House Man And Age of the Age	Algebra Males Hay Control Hay South Hay Sales Short Hay South Hay

Figure 2: Statistical Analysis of the *Creative Ranking* dataset. (a) summarizes some basic information, while (b) shows the number of products in terms of product categories. (c) conducts CTR analysis by comparing poor and good creatives.

visual modeling. The above works focus on improving the **product ranking** by considering visual information while neglecting the great potential of **creative ranking**. There is a few work so far to address this topic. idealo.de (portal of the German e-commerce market) adopts an aesthetic model[11] to select the most attractive image for each recommended hotel. They believe that photos can be just as important for bookings as reviews. PEAC [34] resembles our method the most and they aim to rank ad creatives based on the visual content. But it is an offline evaluation model that cannot flexibly update the ranking strategy when receiving online observations. Besides, all above methods do *not* model the uncertainty which may lack of exploration ability.

2.3 Multi-armed Bandit Methods

Multi-armed bandits (MAB) problem is a typical sequential decision making process that is also treated as an online decision making problems [32]. A wide range of real world applications can be modeled as MAB problems, such as online recommendation system [16], online advertising [27] and information retrieval [15]. Epsilon-greedy [12], Thompson sampling [25] and UCB [2] are classic context-free algorithms. They use reward/cost from the environment to update their E&E policy without contextual information. It is difficult for model to quickly adjust to new creatives (arms) since the web content undergoes frequent changes. [1, 19, 24] extend these context-free methods by considering side information like user/content representations. They assume that the expected payoff of an arm is linear in its features. The main problem linear algorithms face is their lack of representational power, which they complement with accurate uncertainty estimates. A natural attempt at getting the best of both representation learning ability and accurate uncertainty estimation consists in performing a linear payoffs on top of a neural network. NeuralLinear [24] present a Thompson sampling based framework that simultaneously learn a data representation through neural networks and quantify the uncertainty over Bayesian linear regression. Inspired by this framework, we further improve both the neural network and bandit method that benefit our creative ranking problem.

3 DATASET CONSTRUCTION

In order to promote further research and comparison on creative ranking, we contribute a large-scale creative dataset to the research community. It composes creative images and sequential impression data which can be used for evaluating both visual predictions and E&E strategies. In this section, we first describe how the creatives and online feedbacks are collected in subsection 3.1. Then we provide a statistical analysis of the dataset in subsection 3.2.

3.1 Data Collection

We collect a large and diverse set of creatives from Alibaba display advertising platform during July 1, 2020 to August 1, 2020. The total number of impression is approximately 215 million. There are 500,827 products with 1,707,733 ad creatives. We make this dataset publicly available for further research and comparison. The creative and online feedback collection is subject to the following constraints:

Randomized logging policy. The online system adopts randomized logging policy so that the creatives are randomly drawn to collect an unbiased dataset. Bandit algorithms learn policies through interaction data. Training or evaluation on offline data may suffer from exposure bias called "off-policy evaluation problem" [23]. In [19], they demonstrate that if logging policy chooses each arm uniformly at random, the estimation of bandit algorithms is unbiased. Thus, for each impression of product I^n , the policy will randomly choose a candidate creative, and gather their clicks.

Aligned creative lifetime. Due to the complexity of online environment, the CTRs vary for different time periods, even for the same creative. Creatives will be newly designed or deleted, which will result to inconsistent exposure time (as Figure 3(a)). In order to avoid the noise brought by the different time intervals, we only collect the overlap period among the candidate creatives (see Figure 3(b)). Besides, the overlap should be within 5 to 14 days, which covers the creative lifetime from the cold-starting to relative stable stage. All the filtered creatives are gathered to build the sequential data.

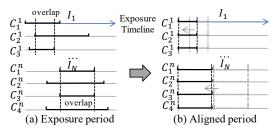


Figure 3: Aligned creative lifetime.

Train/Validation/Test split. We randomly split the 500,827 products into 300,242 training, 100,240 validation and 100,345 test samples, with 1,026,378/340,449/340,906 creatives respectively. We

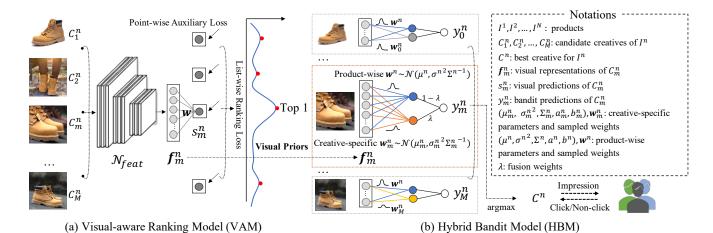


Figure 4: (Better viewed in color) The overall framework of the proposed Hybrid Bandit Model with Visual Priors. It receives several candidate creatives (shown in one column on the left) and try to find the most attractive one through both Visual-aware Ranking Model (VAM) and Hybrid Bandit Model (HBM). (a) VAM is to develop a CNN model that can evaluate creatives base on their visual content. (b) According to the visual priors, HBM aims to estimate the posterior and correct the ranking strategy.

treat each product as a sample, and aim to select the best creative among candidates. The proposed VAM is learned from the training set, while the bandit model HBM is deployed on the validation/test data. This setting is used to prove the effectiveness of visual predictions on the unseen products/creatives, and whether the policy can make a better posterior estimations by using online observations.

3.2 Statistical Analysis

The proposed dataset is collected from ad interaction logs across 32 days. Figure 2(a) gives a summary of our CreativeRanking dataset. It consists of 500,827 products, covering 124 categories. The min and max candidate creatives for a product is 3 and 11, while average number is 3.4. In fact, the number of candidates in the real-world scenarios far exceeds 3.4, but the offline dataset is constrained by conditions introduced by subsection 3.1. Figure 2(b) shows the number of products for top 20 categories, namely Women's tops, Women's bottoms, Men's, Women's shoes, Furniture, and so on. In Figure 2, we make further analysis about creatives for these categories. Suppose we know the CTR for each creative, we select the poorest and best creatives for each product, and accumulate their overall performance, which is visualized as grey and (grep+blue+orange) bins. We find that the CTR of a product can be extremely lifted by selecting a good creative. Specifically, a good creative is capable of lifting CTR by 148% ~ 285% compared to the poorest candidates, while it turns to 41.5% ~ 72.5% compared to averaged performance of all candidates (grep+blue bins).

By proposing this *CreativeRanking* dataset, we would like to draw more attention to this topic which benefits both the research community and website's user experience.

4 METHOD

4.1 Overview

We briefly overview the entire pipeline. Main notations used in this paper are summarized in the right panel of Figure 4. First, as shown in Figure 4(a), feature extraction network N_{feat} will simultaneously receive the creatives of the n-th product as input, and produce the d-dimensional intermediate features $\{f_1^n, f_2^n, \cdots, f_m^n, \cdots, f_M^n\}$. Then, a fully connected layer are employed to calculate the scores for them, indicated as $\{s_1^n, s_2^n, \cdots, s_m^n, \cdots, s_M^n\}$.

Second, the *list-wise ranking loss and auxiliary regression loss* are introduced to guide the learning procedure. Such a multi-objective optimization helps the model not only focus on creative ranking, but also take into account the numerical range of CTR that is benefit for the following bandit model. In addition, due to the data noise that is a common problem in a real-world application, we provide several practical solutions to mitigate casual and malicious noise. Details are described in Subsection 4.2.

After the above steps, the model can evaluate the creative quality directly by its visual content, even a newly uploaded one without any history information. Then we propose a hybrid bandit model that incorporates learned f_m^n as contextual information, and update the policy by interacting with online observations. As in Figure 4(b), the hybrid model combines both product-wise and creative-specific predictions which is more flexible for complex industrial data. The elaborated formulations are in Subsection 4.3.

4.2 VAM: Visual-aware Ranking Model

Given a product I^n , we use feature extraction network \mathcal{N}_{feat} to extract high-level visual representations of creatives. And a linear layer is adopted to produce the attractiveness scores for m-th creative of n-th product:

$$f_m^n = \mathcal{N}_{feat}(C_m^n) \tag{5}$$

$$s_m^n = f_m^{nT} \mathbf{w} \tag{6}$$

where \boldsymbol{w} are learnable parameters of the linear layer.

List-wise Ranking Loss. To learn the relative order of creatives, we need to map a list of predicted scores and ground-truth CTRs to a permutation probability distribution, respectively, and then

take a metric between these distributions as a loss function. The mapping strategy and evaluation metric should guarantee that the candidates with higher scores would be ranked higher. [5] proposed permutation probability and top k probability definitions. Inspired by this work, we simplify the probability of a creative being ranked on the top 1 position as

$$p_m^n = \frac{exp(s_m^n)}{\sum_{i=1}^M exp(s_i^n)} \tag{7}$$

where $exp(\cdot)$ is an exponential function. The exponential function based top-1 probability is both scale invariant and translation invariant. And the corresponding labels are

$$y_{rank}(C_m^n) = \frac{exp(CTR(C_m^n), T)}{\sum_{i=1}^{M} exp(CTR(C_i^n), T)}$$
(8)

where $exp(\cdot,T)$ is exponential function with temperature T. Since the $CTR(C_m^n)$ is about a few percent, we use T to adjust the scale of the value so that make the probability of top1 sample close to 1. With Cross Entropy as metric, the loss for product I^n becomes

$$\mathcal{L}_{rank}^{n} = -\sum_{m} y_{rank}(C_{m}^{n}) log(p_{m}^{n})$$
 (9)

Through such objective function, the model focuses on comparing the creatives within the same product. We concentrate on the top-1 probability since it is consistent with real scenarios which will display only one creative for each impression. Besides, the end-to-end training manner greatly utilizes the learning ability of deep CNNs and boosts the visual prior knowledge extraction.

Point-wise auxiliary regression Loss. In addition to the listwise ranking loss, we expect that the point-wise regression enforce the model to produce more accurate predictions. Actually, the ranking loss function only requires the order of outputs, leaving the numerical scale of the outputs unconstrained. Since the learned representations will be adopted as prior knowledge for the bandit model in Subsection 4.3, making the outputs close to the real CTRs will significantly stabilize the bandit learning procedure. Thus we add the point-wise regression as a regularizer. The formulation is

$$\mathcal{L}_{reg}^{n} = \sum_{m} ||CTR(C_{m}^{n}) - s_{m}^{n}||_{2}$$
 (10)

where $||\cdot||$ denotes L_2 norm. Finally, we add up both the ranking loss and the auxiliary loss to form the final loss:

$$\mathcal{L}^n = \mathcal{L}_{rank}^n + \gamma \mathcal{L}_{req}^n \tag{11}$$

where γ is 0.5 in our experiments.

Noise Mitigation. In both list-wise ranking and point-wise regression in Equation 8 and 10, $CTR(C_m^n)$ can be estimated by Equation 2. But in real-world dataset, some creatives have not sufficient impression opportunities, and the estimation may suffer from huge variance. For example, a creative only get one impression, and a click is accidentally recorded from this impression, the $C\hat{T}R$ will be set to 1, which is inevitably unreliable. To mitigate the problem, we provide two practical solutions, namely label smoothing and weighted sampling.

Label smoothing is an empirical Bayes method that is utilized to smoothen the CTR estimation [30]. Suppose the clicks are from a binomial distribution and the CTR follows a prior distribution as

$$\begin{aligned} & clicks(C_m^n) \sim Binomial(Impression(C_m^n), CTR(C_m^n)) \\ & & CTR(C_m^n) \sim Beta(\alpha, \beta) \end{aligned} \tag{12}$$

where $Beta(\alpha,\beta)$ can be regarded as the prior distribution of CTRs. After observing more clicks, the conjugacy between Binomial and Beta allows us to obtain the posterior distribution and the smoothed $C\hat{T}R$ as

$$\hat{CTR}(C_m^n) = \frac{click(C_m^n) + \alpha}{impression(C_m^n) + \alpha + \beta}$$
 (13)

where α and β can be yielded by using maximum likelihood estimate through all the historical data[30]. Compared to the original way, the smoothed \hat{CTR} has smaller variance and benefits the training.

Weighted sampling is a sampling strategy for training process. Instead of treating each product equally, we pay more attention to the products whose impressions are adequate and the CTRs are more reliable. The sampling weights can be produced by

$$p^n = g(impression(I^n)) (14)$$

where $g(\cdot)$ is set to the logarithm of the impressions and p^n denotes the sampling weight of product I^n .

All above modules are integrated in a unified framework and the visual-aware ranking model focuses on learning the general visual patterns about display performance. And then the informative representations are applied as prior knowledge for the bandit algorithm.

4.3 HBM: Hybrid Bandit Model

In this section, we provide an elegant and efficient strategy that tackles the E&E dilemma by utilizing the visual priors and updating the posterior via the hybrid bandit model. Based on NeuralLinear framework [24], we build a Bayesian linear regression on the extracted visual representation. Assume the online feedback data is generated as follows:

$$\mathbf{y} = \mathbf{f}^T \tilde{\mathbf{w}} + \epsilon \tag{15}$$

where \boldsymbol{y} represent clicked/non-clicked data and f is the extracted visual representations by VAM. Different from the deterministic weights \boldsymbol{w} in Equation 6, we need to learn a weight distribution $\tilde{\boldsymbol{w}}$ with the uncertainty that benefits the E&E decision making. ϵ are independent and identically normally distributed random variables:

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$
 (16)

According to Bayes theorem, if the prior distribution of \tilde{w} and σ^2 is conjugate to the data's likelihood function, the posterior probability distributions can be derived analytically. And then Thompson Sampling, as known as Posterior Sampling, is able to tackles the E&E dilemma by maintaining the posterior over models and selecting creatives in proportion to the probability that they are optimal. We model the prior joint distribution of \tilde{w} and σ^2 as:

$$\pi(\tilde{\mathbf{w}}, \sigma^2) = \pi(\tilde{\mathbf{w}}|\sigma^2)\pi(\sigma^2),$$

$$\sigma^2 \sim IG(a, b) \text{ and } \tilde{\mathbf{w}}|\sigma^2 \sim \mathcal{N}(\mu, \sigma^2\Sigma^{-1})$$
(17)

where the $IG(\cdot)$ is an Inverse Gamma whose prior hyperparameters are set to $a_0 = b_0 = \eta > 1$ and $\mathcal{N}(\cdot)$ is a Gaussian distribution with the initial parameters $\Sigma_0 = \lambda Id$. Note that μ_0 is initialized as the learned weights \boldsymbol{w} of VAM in Equation 6. It can provide a better prior hyperparameters that further enhance the performance in

the cold-starting phase. We call it VAM-Warmup and the results is shown in Figure 5(b).

Because we have chosen a conjugate prior, the posterior at time t can be derived as

$$\Sigma(t) = \mathbf{f}^T \mathbf{f} + \Sigma_0$$

$$\mu(t) = \Sigma(t)^{-1} (\Sigma_0 \mu_0 + \mathbf{f}^T \mathbf{y})$$

$$a(t) = a_0 + t/2$$

$$b(t) = b_0 + \frac{1}{2} (\mathbf{y}^T \mathbf{y} + \mu_0^T \Sigma_0 \mu_0 - \mu(t)^T \Sigma(t) \mu(t))$$
(18)

where $f \in \mathbb{R}^{t \times d}$ is a matrix that contain the content features for previous impressions and $\mathbf{y} \in \mathbb{R}^{t \times 1}$ is the feedback rewards. After updating the above parameters at t-th impression, we obtain the weight distributions with uncertainty estimation. We draw the weights $\mathbf{w}(t)$ from the learned distribution $\mathcal{N}(\mu(t), \sigma(t)^2 \Sigma(t)^{-1})$ and select the best creative for product I^n as

$$C^{n} = \underset{c \in \{C_{1}^{n}, C_{2}^{n}, \dots, C_{M}^{n}\}}{\arg \max} \left(\mathcal{N}_{feat}(c) \right)^{T} w(t) \tag{19}$$

The above model makes the weight distributions shared by all the products. This simple linear assumption works well for small datasets, but becomes inferior when dealing with industrial data. For example, bright and vivid colors will be more attractive for women's top while concise colors are more proper for 3C digital accessories. In addition to this product-wise characteristic, a creative may contain a unique designed attribute that is not expressed by the shared weights. Hence, it is helpful to have weights that have both shared and non-shared components.

We extend the Equation 15 to the following hybrid model by combining product-wise and creative-specific linear terms. For creative C_m^n , it can be formulated as

$$y_m^n = f_m^{nT} \mathbf{w}^n + f_m^{nT} \mathbf{w}_m^n \tag{20}$$

where \mathbf{w}^n and \mathbf{w}^n_m are product-wise and creative-specific parameters, and they are disjoinly optimized by Equation 18. Furthermore, we propose an fusion strategy to adaptively combine these two terms instead of the simple addition

$$y_m^n = (1 - \lambda) f_m^{nT} \mathbf{w}^n + \lambda f_m^{nT} \mathbf{w}_m^n$$
 (21)

where $\lambda=(1+e^{\frac{-impression(I^n)+\theta_2}{\theta_1}})^{-1}$ is a sigmoid function with rescale θ_1 and offset θ_2 . We find that if the impressions are inadequate, the product-wise parameters are learned better because it make use of the knowledge among all candidate creatives. Otherwise, the creative-specific term outperforms the shared one due to the sufficient feedback observations. The above procedure is shown in Algorithm 1. Because our hybrid model updates the parameters of each product independently, we take I^n as example and adopt $(a^n(\cdot),b^n(\cdot),\mu^n(\cdot),\Sigma^n(\cdot))$ and $(a^n_m(\cdot),b^n_m(\cdot),\mu^n_m(\cdot),\Sigma^n_m(\cdot))$ to represent the shared and specific parameters. The distributions describe the uncertainty in weights which is related to impressed number: if there is less data, the model relies more on the visual evaluation results; Otherwise, the likelihood will reduce the priori effect so as to converge to the observation data. In order to fit the complex industrial data, we extend the shared linear model to the hybrid version, which consider both product-level knowledge

Algorithm 1: Hybrid Bandit Model

```
Input: T > 0, product I^n, visual representations of
               candidate creatives f_0^n, f_1^n, \dots, f_m^n, \dots, f_M^n
 1 Initialize the a_0, b_0, \mu_0 and \Sigma_0;
 a^{n}(0) \leftarrow a_{0}, b^{n}(0) \leftarrow b_{0}, \mu^{n}(0) \leftarrow \mu_{0}, \Sigma^{n}(0) \leftarrow \Sigma_{0};
 a_m^n(0) \leftarrow a_0, b_m^n(0) \leftarrow b_0, \mu_m^n(0) \leftarrow \mu_0, \Sigma_m^n(0) \leftarrow \Sigma_0;
4 for t = 1, 2, 3, \dots, T do
5 \lambda = (1 + e^{\frac{-impression(I^n) + \theta_2}{\theta_1}})^{-1};
          for m = 1, 2, 3, ..., M do
                Sample \sigma^{n2} from IG(a^n(t-1), b^n(t-1));
                Sample w^n from \mathcal{N}(\mu^n(t-1), \sigma^{n^2}\Sigma^n(t-1)^{-1});
              Sample \sigma_m^n from IG(a_m^n(t-1), b_m^n(t-1));

Sample w_m^n from \mathcal{N}(\mu_m^n(t-1), \sigma_m^n {}^2\Sigma_m^n(t-1)^{-1});

y_m^n = (1-\lambda)f_m^{nT}w^n + \lambda f_m^{nT}w_m^n;
10
11
12
          k = \arg\max{(y_1^n, \dots, y_m^n, \dots, y_M^n)};
13
          Display the creative C_k^n, and get the reward;
14
          Update a^n(t), b^n(t), \mu^n(t), \Sigma^n(t) by the historical data
15
           of product I^n and Equation 18;
          Update a_k^n(t), b_k^n(t), \mu_k^n(t), \Sigma_k^n(t) by the historical data
16
           of creative C_k^n and Equation 18;
          Set the other parameters of time t as the same as
17
           previous time (t-1);
          impression(I^n) \leftarrow impression(I^n) + 1;
18
19 end
```

and creative-specific information, and fused by empirical attention weights.

5 EXPERIMENTS

5.1 Dataset preparation and Evaluation Metrics

Dataset Preparation. The description of CreativeRanking data is presented in Section3.1. The original images and rewards for each creative are provided in the order of displaying. For VAM, we aggregate the number of impressions and clicks to produce \hat{CTR} by Equation 13 on training set, and train the VAM using the loss function in Equation 11. For HBM, we update the policy by providing the visual representations extracted by VAM and the impression data like Equation 4. Note the interaction and policy updating procedure (see Algorithm 1) of HBM is conducted on the test set for simulating the online situations. We record the sequential interactions and rewards to measure the performance (see Algorithm 2 and Equation 21). Validation is used for hyperparameter tuning.

In addition to the CreativeRanking data, we also evaluate the methods on a public dataset, called Mushroom. Since there is no public dataset for creative ranking yet, we test the proposed hybrid bandit model on this standard dataset. The Mushroom Dataset [26] contains 22 attributes for each mushroom, and two categories: poisonous and safe. Eating a safe mushroom will receive reward +5 while eating a poisonous one delivers reward +5 with probability 50% and reward -35 otherwise. Not eating will provide no reward. We follow the protocols in [24], and interact for 50000 rounds.

Evaluation Metrics. For CreativeRanking data, we present two evaluation metrics to measure the performance, named simulated CTR (*sCTR*) and cumulative regret (*Regret*), respectively.

Simulated CTR (sCTR) is a practical metric which is quite close to the online performance. The details are shown in Algorithm 2. It replays the recorded impression data for all products. For each product, the policy will play T^n rounds by receiving the recorded data (C, y), and selects the best creative according to the predicted scores. If the selected one is the same as the C, the impression number, click number and policy itself will be updated (see line 3 to 14 in Algorithm 2).

Take HBM as an example, algorithm 1 shows the online update process. To test the HBM by using offline data, we can change the action "display and update" (line 14 to 18 in Algorithm 1) to the conditioned version in the line 8 to 12 in Algorithm 2.

Cumulative regret (*Regret*) is commonly used for evaluating bandit models. It is defined as

$$Regret = E[r^* - r] \tag{22}$$

where r^* is the cumulative reward of the optimal policy, i.e., the policy that always selects the action with highest expected reward given the context [24]. Specifically, we select the optimal creative for our dataset, and calculate the *Regret* as

$$Regret = \frac{\sum_{n=1}^{N} click(C^{n})}{\sum_{n=1}^{N} impression(C^{n})} - sCTR$$
 (23)

where sCTR should be produced by Algorithm 2 first. And the C^n is selected by calculating $C\hat{T}R$ in Equation 2 on the test set.

For Mushroom, we follow the definition of cumulative regret in [24] to evaluate the models.

5.2 Implementation details

The model was implemented with Pytorch [22]. We adopt deep residual network (ResNet-18)[17] pretrained on ImageNet classification [10] as backbone, and the model is finetuned with Creative

Algorithm 2: Evaluation Metrics - sCTR

```
Input: impression data, policy \pi
   Output: sCTR
 1 impressions \leftarrow 0;
2 clicks \leftarrow 0;
 3 for n = 1, 2, 3, ..., N do
        for t = 1, 2, 3, ..., T^n do
 4
             Get next impression (C, y);
 5
             Get predicated scores (y_1^n, \ldots, y_M^n) by policy \pi;
 6
             k = \arg\max(y_1^n, \dots, y_M^n);
 7
             if C_k^n = C then
                  impressions \leftarrow impressions + 1;
                  clicks \leftarrow clicks + y;
10
                  update policy \pi by data (C, y);
             end
12
        end
13
14 end
15 sCTR = \frac{clicks}{impressions};
16 return sCTR
```

Ranking task. For VAM, we use stochastic gradient descent (SGD) with a mini-batch of 64 per GPU. The learning rate is initially set to 0.01 and then gradually decreased to e-4. The training process lasts 30 epochs on the datasets. For HBM, we extract the feature representations f_m^n from VAM, and update the weights distribution w_m^n and w^n by using bayesian regression.

5.3 Comparison with State-of-the-art Systems

In this subsection, we show the performance of the related methods in Table 1 and Figure 5. The methods are divided into some groups: a uniform strategy, context-free bandit models, linear bandit models, neural bandit models and our proposed methods. Table 1 presents the *Regret* and *sCTR* of all above models on both Mushroom and CreativeRanking datasets, and our methods - (NN/VAM-HBM) exhibits state-of-the-art results compared to the related models. We also conduct further analysis by showing the reward tendency of consecutive 15 days in Figure 5. Daily *sCTR* evaluates the model for each day independently, showing the flexibility of the policy when interacting with the feedback. And cumulative *sCTR* presents the cumulative rewards up to the specific day which is used to measure the overall performance.

Uniform: The baseline strategy that randomly selects an action (eat/not eat for Mushroom and one creative for CreativeRanking). Because this strategy has neither prior knowledge nor abilities of learning from the data, it gets poor performance on the test sets.

Context-free Bandit Models: Epsilon-greedy [12], Thompson sampling [25] and Upper Confidence Bounds (UCB) approaches [2] are simple yet effective strategies to deal with the bandit problem. They rely on history impression data (click/non-click) and keep updating their strategies. However, for the cold-start stage, they might randomly choose a creative like "Uniform" strategy (orange lines in Figure 5(c) in the first few days). We find that their curves are gradually rising, but without prior information, the overall performance is inferior to the other models.

Linear Bandit Models: The linear bandit model is an extension to the context-free method by incorporating contextual information. For Mushroom, we adopt the 22 attributes to describe a mushroom, such as shape, color and so on. The *Regret* is reduced when combining the side information. For CreativeRanking, we use color distribution [3] to represent a creative, and update the linear payoff functions. From the results in Table 1, the linear models achieve better results than the context-free methods, but they still face the problem of lacking representational power.

Neural Bandit Models: The neural bandit models add a linear regression on top of the neural network. In Table 1, "NN" denotes fully connected layers that used for extracting mushroom representations. For CreativeRanking, all these neural models use our VAM as feature extractor, and adopt different E&E policies. Figure 5(a) reveals some interesting observations: (1) The orange and blue lines represent the E-greedy and VAM-Greedy, respectively. With the visual priors, VAM-Greedy achieves much better performance at the beginning (about 5% CTR lift), which demonstrates the effectiveness of the visual evaluation. (2) Because VAM-Greedy is a greedy strategy that lack of exploration, it becomes mediocre in the long run. When incorporating E&E model - HBM, our VAM-HBM outperforms the other baselines by a significant margin. Besides,

	Mushroom	Creative	Ranking
Evaluation Metrics	Regret (%)	Regret (%)	sCTR (%)
Uniform	100	100	2.950
Context-free Bandit Mode	ls (Orange lin	es)	
E-Greedy[12]	52.99	87.22	3.166
Thompson Sampling[25]	52.49	87.69	3.158
UCB[2]	52.42	87.04	3.169
Linear Bandit Models (Gre	en lines)		
LinGreedy[24]	14.28	91.72	3.090
LinThompson[24]	2.37	85.68	3.192
LinUCB[19]	10.27	85.50	3.195
Neural Bandit Models (Blu	ie lines)		
NN/VAM-Greedy	6.68	84.11	3.219
NN/VAM-Thompson[24]	2.22	83.02	3.237
NN/VAM-UCB	7.51	83.91	3.222
NN/VAM-Dropout[13]	5.57	84.32	3.215
Our Methods (red lines)			
VAM-Warmup	-	79.70	3.293
NN/VAM-HBM	1.93	78.11	3.320

Table 1: Performance comparison with state-of-the-art systems on both *Mushroom* and *CreativeRanking* test set. *Regret* is Normalized with respect to the performance of Uniform.

we also use Dropout as a Bayesian approximation[13], but it is not able estimate the uncertainty as accurate as the other policies.

Our Methods: We propose VAM-Warmup that initialize the μ_0 in bandit model by learned weights in VAM. By comparing red and blue dashed lines in Figure 5(b), we find the parameters with prior distributions improves 1.7% CTR for overall performance. In addition, we extend the model by adding creative-specific parameters, named VAM-HBM, and it further enhances the model capacity and achieves the state-of-the-art result, especially the impressions for creatives become adequate (see solid red line in Figure 5(b)(c)(d)). For Mushroom dataset, in order to demonstrate the idea, we cluster the data into 2 groups by attribute "bruises", each maintaining the individual parameters. When combining the individual and shared parameters by fusion weights in Equation 21, the model reduces the *Regret* to 1.93. Note that we use the default hyperparameters provided by NeuralLinear without carefully tuning.

5.4 Ablation Study

In this subsection, we conduct an ablation study on CreativeRanking dataset so as to validate the effectiveness of each component in the VAM, including list-wise ranking loss, point-wise auxiliary regression loss and noise mitigation. Besides, we also compare our VAM with "learning-to-rank" visual models (including aesthetic

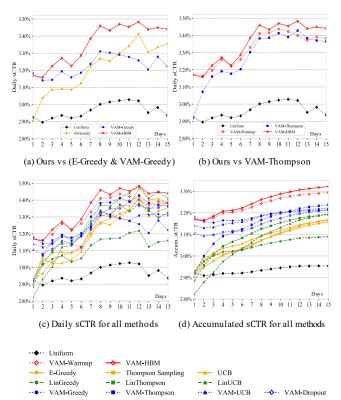


Figure 5: Reward tendency of consecutive 15 days on CreativeRanking.

Methods	Base	(a)	(b)	(c)	(d)
Point-wise Loss?		√		√	√
List-wise Loss?			√	√	√
Noise Mitigation?					√
sCTR(%)	2.950	3.140 ↑ 6.4%	3.167 ↑ 7.4%	3.194 ↑ 8.3%	3.219 ↑ 9.1%

Table 2: Ablation study for each component in the VAM. sCTR are performed on the CreativeRanking test set and $\uparrow sCTR$ lift is calculated by $\frac{(sCTR(*)-sCTR(base))}{sCTR(base)}*100\%$.

models). We show the results in Table 2 and Table 3 to demonstrate the consistent improvements.

Base in Table 2 stands for the baseline result. We adopt "uniform" strategy that randomly choose a creative among the candidates. The baseline is 2.950% for sCTR.

Method (a) and (b): Method (a) and (b) utilize point-wise (Equation 10) and list-wise loss (Equation 9) as the objective function, respectively. Although the model has never seen the products/creatives on the test set before, it has learned general patterns to identify more attractive creatives. Moreover, the ranking loss concentrates on the top-1 probability learning which is more suitable than the point-wise objective for our scenarios. The simple version (b) can improve the *sCTR* by 7.4%.

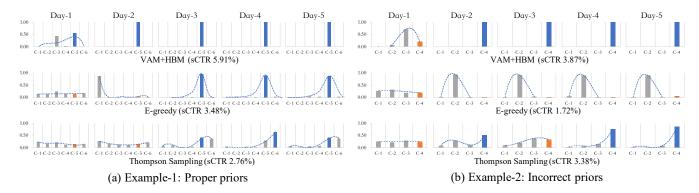


Figure 6: Two typical cases that present the changing of strategies. The horizontal axis shows different creatives while the vertical axis is the probability of being displayed for creatives. "Proper priors" indicates that VAM provides right predictions and "Incorrect priors" otherwise.

Ranking Loss	sCTR (%)
Pairwise Hinge Loss [7]	3.170
Aesthetics Ranking Loss [18]	3.167
Triplet Loss [28]	3.115
Pairwise [34]	3.188
VAM (Ours)	3.219

Table 3: Comparison with other "learn-to-rank" visual models. All above models adopt ResNet-18 as backbone.

Method (c): Method (c) combines the point-wise auxiliary regression loss with the ranking objective. It not only learns the relative order of creative quality, but also the absolute CTRs. We find it is good at fitting the real CTR distributions and achieve the better performance 3.194% (8.3% lift) for *sCTR*.

Method (d): Method (d) contains *label smooth* and *weighted sampler*, both of which are designed for mitigating the label noise. Weighted sampler makes the model pay more attention to the samples whose impression numbers are sufficient while label smooth aims to improve the label reliability. These two practical methods further improve the *sCTR* to 3.216%, lifting 9.1% in total.

Related Loss functions: Pair-wise and triplet loss are typical loss functions for learning to rank problems. [7, 18, 28] adopt hinge loss that is used for "maximum-margin" classification between the better candidate and the other one. It only requires the better creative to get higher score than the other one by a pre-defined margin, without consideration of the exact difference. Our loss function in Equation 9 and 10 estimate their CTR gaps and produce more accurate differences. [34] employ [4] as their pair-wise framework. Compared to the pair-wise learning, we treat one product as a training sample and use list of creatives as instances. It is more efficient and suitable with real scenarios which will display the best creative for one impression. Thus, our method obtains the leading performance on *sCTR*.

In summary, the proposed list-wise method enables the model focus on learning creative qualities and obtains better generalizability. Incorporating point-wise regression and noise mitigation

γ in Equation 11	0.0	0.1	0.5	1.0	2.0
Validation sCTR(%)	3.15	3.15	3.17	3.16	3.13
Test sCTR(%)	3.17	3.19	3.22	3.18	3.15

Table 4: Val/Test sCTR with different y in Equation 11.

$ heta_1/ heta_2$ in λ	125	150	175
30	3.27%(3.32%)	3.28%(3.33%)	3.28%(3.31%)
50	3.27%(3.31%)	3.28%(3.32%)	3.27%(3.32%)
100	3.27%(3.31%)	3.27%(3.31%)	3.27%(3.31%)

Table 5: x%(x%) denotes val(test) sCTR of different θ_1/θ_2 in λ .

techniques is able to enhance the model capacity of fitting the real-world data.

5.5 Hyperparameter Settings

 γ in Equation 11. We tune hyperparameters in the validation set. γ in Equation 11 is adopted to control the weight of point-wise auxiliary loss. According to the validation results (see Table 4), we take $\gamma = 0.5$. It is consistent with our hypothesis that ranking loss should play a more important role in the creative ranking task. θ_1/θ_2 of λ in Equation 21. θ_1/θ_2 control the rescale and offset of λ in Equation 21. Optimal hyperparameters vary in different real-world platforms(e.g., offset is set to 150, around the mean impression number of each creative). We find the final performance is not sensitive to these hyperparameters (see Table 5). We choose $\theta_1 = 50$ and $\theta_2 = 150$ in our experiments.

5.6 Case Study

Strategy Visualization. We show two typical cases that exhibit the changing of strategies. Figure 6 (a) shows the proper prior of HBM. We believe that the best creative should have the largest displaying probability among candidates. If this expectation is satisfied, a blue bar is shown; otherwise, orange bars are shown. It grants most impression opportunities to creative C-5 from the first day, while the other two methods spend 2 days to find the best creative. For another case that receives incorrect prior in Figure

Attention to models Okiss Okiss Okiss Attention to text Attention to text

Figure 7: Visualization of the learned VAM. The model pays attention to different regions adaptively, including products, models and the text on the creative.

6(b), the HBM adjusts the decision by considering the online feedback. The interactions help to revise the prior knowledge and fit to the real-world feedback. Form this comparison, we find the HBM makes good use of visual priors, and adjusts flexibly according to the feedback signals.

CNN Visualization. Besides ranking performance, we would like to attain further insight into the learned VAM. To this end, we visualize the response of our VAM according to the activations on the high-level feature maps, and the resulting visualization is shown in Figure 7. By learning from the creative ranking, we find that the CNN pays attention to different regions adaptively, including products, models and the text on the creative. As shown in the second row Figure 7, the VAM draw higher attention to the models rather than the products. It may caused by the reason that products endorsed by celebrities are more attractive than simply displaying the products. Besides, some textual information, such as description and discount information, can also attract customers.

6 CONCLUSIONS

In this paper, we propose a hybrid bandit model with visual priors. To the best of our knowledge, this is the first time that formulates the creative ranking as a E&E problem with visual priors. The VAM adopts a list-wise ranking loss function for ordering the creative quality only by their contents. In addition to the ability of visual evaluation, we extend the model to be updated when receiving feedback from online scenarios called HBM. Last but not the least, we construct and release a novel large-scale creative dataset named *CreativeRanking*. We would like to draw more attention to this topic which benefits both the research community and website's user experience. We carried out extensive experiments, including performance comparison, ablation study and case study, demonstrating the solid improvements of the proposed model.

REFERENCES

- [1] Shipra Agrawal and Navin Goyal. 2013. Thompson Sampling for Contextual Bandits with Linear Payoffs. In Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013 (JMLR Workshop and Conference Proceedings, Vol. 28). JMLR.org, 127-135. http://proceedings.mlr. press/v28/agrawal13.html
- [2] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. 2002. Finite-time Analysis of the Multiarmed Bandit Problem. Mach. Learn. 47, 2-3 (2002), 235–256. https://doi.org/10.1023/A:1013689704352
- [3] Javad Azimi, Ruofei Zhang, Yang Zhou, Vidhya Navalpakkam, Jianchang Mao, and Xiaoli Z. Fern. 2012. The impact of visual appearance on user response in online display advertising. In Proceedings of the 21st World Wide Web Conference, WWW 2012, Lyon, France, April 16-20, 2012 (Companion Volume), Alain Mille, Fabien L. Gandon, Jacques Misselis, Michael Rabinovich, and Steffen Staab (Eds.). ACM, 457-458. https://doi.org/10.1145/2187980.2188075
- [4] Christopher J. C. Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Gregory N. Hullender. 2005. Learning to rank using gradient descent. In Machine Learning, Proceedings of the Twenty-Second International Conference (ICML 2005), Bonn, Germany, August 7-11, 2005 (ACM International Conference Proceeding Series, Vol. 119), Luc De Raedt and Stefan Wrobel (Eds.). ACM, 89-96.
- [5] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. Learning to rank: from pairwise approach to listwise approach. In Machine Learning, Proceedings of the Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007 (ACM International Conference Proceeding Series, Vol. 227), Zoubin Ghahramani (Ed.). ACM, 129–136. https://doi.org/10.1145/ 1273496.1273513
- [6] Mark Capelo, Karan Aggarwal, and Pranjul Yadav. 2019. Combining Text and Image data for Product Recommendability Modeling. In 2019 IEEE International Conference on Big Data (Big Data), Los Angeles, CA, USA, December 9-12, 2019. IEEE, 5992–5994. https://doi.org/10.1109/BigData47090.2019.9006197
- [7] Parag S. Chandakkar, Vijetha Gattupalli, and Baoxin Li. 2017. A Computational Approach to Relative Aesthetics. CoRR abs/1704.01248 (2017). arXiv:1704.01248 http://arxiv.org/abs/1704.01248
- [8] Junxuan Chen, Baigui Sun, Hao Li, Hongtao Lu, and Xian-Sheng Hua. 2016. Deep CTR Prediction in Display Advertising. In Proceedings of the 2016 ACM Conference on Multimedia Conference, MM 2016, Amsterdam, The Netherlands, October 15-19, 2016, Alan Hanjalic, Cees Snoek, Marcel Worring, Dick C. A. Bulterman, Benoit Huet, Aisling Kelliher, Yiannis Kompatsiaris, and Jin Li (Eds.). ACM, 811-820. https://doi.org/10.1145/2964284.2964325
- [9] Haibin Cheng, Roelof van Zwol, Javad Azimi, Eren Manavoglu, Ruofei Zhang, Yang Zhou, and Vidhya Navalpakkam. 2012. Multimedia features for click prediction of new ads in display advertising. In The 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD, Qiang Yang, Deepak Agarwal, and Jian Pei (Eds.). ACM, 777–785.
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. ImageNet: A large-scale hierarchical image database. In 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009), 20-25 June 2009, Miami, Florida, USA. IEEE Computer Society, 248–255. https://doi.org/10.1109/CVPR.2009.5206848
- [11] Hossein Talebi Esfandarani and Peyman Milanfar. 2018. NIMA: Neural Image Assessment. IEEE Trans. Image Process. 27, 8 (2018), 3998–4011. https://doi.org/ 10.1109/TIP.2018.2831899
- [12] Vincent François-Lavet, Peter Henderson, Riashat Islam, Marc G. Bellemare, and Joelle Pineau. 2018. An Introduction to Deep Reinforcement Learning. Found. Trends Mach. Learn. 11, 3-4 (2018), 219–354. https://doi.org/10.1561/2200000071
- [13] Yarin Gal and Zoubin Ghahramani. 2016. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In Proceedings of the 33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016 (JMLR Workshop and Conference Proceedings, Vol. 48), Maria-Florina Balcan and Kilian Q. Weinberger (Eds.). JMLR.org, 1050–1059. http://proceedings.mlr.press/v48/gal16.html
- [14] Tiezheng Ge, Liqin Zhao, Guorui Zhou, Keyu Chen, Shuying Liu, Huimin Yi, Zelin Hu, Bochao Liu, Peng Sun, Haoyu Liu, et al. 2018. Image matters: Visually modeling user behaviors using advanced model server. In Proceedings of the 27th ACM International Conference on Information and Knowledge Management. 2087–2095.
- [15] Dorota Glowacka. 2017. Bandit Algorithms in Interactive Information Retrieval. In Proceedings of the ACM SIGIR International Conference on Theory of Information Retrieval, ICTIR 2017, Amsterdam, The Netherlands, October 1-4, 2017, Jaap Kamps, Evangelos Kanoulas, Maarten de Rijke, Hui Fang, and Emine Yilmaz (Eds.). ACM, 327–328. https://doi.org/10.1145/3121050.3121108
- [16] Dorota Glowacka. 2019. Bandit algorithms in recommender systems. In Proceedings of the 13th ACM Conference on Recommender Systems. 574–575.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition. 770–778.

- [18] Shu Kong, Xiaohui Shen, Zhe L. Lin, Radomír Mech, and Charless C. Fowlkes. 2016. Photo Aesthetics Ranking Network with Attributes and Content Adaptation. In Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 9905), Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (Eds.). Springer, 662–679. https://doi.org/10.1007/978-3-319-46448-0_40
- [19] Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. 2010. A contextual-bandit approach to personalized news article recommendation. In Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26-30, 2010, Michael Rappa, Paul Jones, Juliana Freire, and Soumen Chakrabarti (Eds.). ACM, 661-670. https://doi.org/10.1145/1772690. 1772758
- [20] Hu Liu, Jing Lu, Hao Yang, Xiwei Zhao, Sulong Xu, Hao Peng, Zehua Zhang, Wenjie Niu, Xiaokun Zhu, Yongjun Bao, et al. 2020. Category-Specific CNN for Visual-aware CTR Prediction at JD. com. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2686–2696.
- [21] Kaixiang Mo, Bo Liu, Lei Xiao, Yong Li, and Jie Jiang. 2015. Image Feature Learning for Cold Start Problem in Display Advertising. In Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015, Qiang Yang and Michael J. Wooldridge (Eds.). AAAI Press, 3728-3734. http://ijcai.org/Abstract/15/524
- [22] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch. (2017).
- [23] Doina Precup. 2000. Eligibility traces for off-policy policy evaluation. Computer Science Department Faculty Publication Series (2000), 80.
- [24] Carlos Riquelme, George Tucker, and Jasper Snoek. 2018. Deep Bayesian Bandits Showdown: An Empirical Comparison of Bayesian Deep Networks for Thompson Sampling. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net. https://openreview.net/forum?id=SyYe6k-CW
- [25] Daniel Russo and Benjamin Van Roy. 2014. Learning to Optimize via Posterior Sampling. Math. Oper. Res. 39, 4 (2014), 1221–1243. https://doi.org/10.1287/moor. 2014.0650

- [26] Jeff Schlimmer. 1981. Mushroom records drawn from the audubon society field guide to north american mushrooms. GH Lincoff (Pres), New York (1981).
- [27] Eric M Schwartz, Eric T Bradlow, and Peter S Fader. 2017. Customer acquisition via display advertising using multi-armed bandit experiments. *Marketing Science* 36, 4 (2017), 500–522.
- [28] Katharina Schwarz, Patrick Wieschollek, and Hendrik PA Lensch. 2018. Will people like your image? learning the aesthetic space. In 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). IEEE, 2048–2057.
- [29] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1409.1556
- [30] Xuerui Wang, Wei Li, Ying Cui, Ruofei Zhang, and Jianchang Mao. 2011. Click-through rate estimation for rare events in online advertising. In Online multimedia advertising: Techniques and technologies. IGI Global, 1–12.
- [31] Yu Wang, Jixing Xu, Aohan Wu, Mantian Li, Yang He, Jinghe Hu, and Weipeng P Yan. 2018. Telepath: Understanding users from a human vision perspective in large-scale recommender systems. In Thirty-Second AAAI Conference on Artificial Intelligence.
- [32] Mengyue Yang, Qingyang Li, Zhiwei Qin, and Jieping Ye. 2020. Hierarchical Adaptive Contextual Bandits for Resource Constraint based Recommendation. In Proceedings of The Web Conference 2020. 292–302.
- [33] Wenhui Yu, Huidi Zhang, Xiangnan He, Xu Chen, Li Xiong, and Zheng Qin. 2018. Aesthetic-based Clothing Recommendation. In Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018, Pierre-Antoine Champin, Fabien L. Gandon, Mounia Lalmas, and Panagiotis G. Ipeirotis (Eds.). ACM, 649-658. https://doi.org/10.1145/3178876.3186146
- [34] Zhichen Zhao, Lei Li, Bowen Zhang, Meng Wang, Yuning Jiang, Li Xu, Fengkun Wang, and Wei-Ying Ma. 2019. What You Look Matters?: Offline Evaluation of Advertising Creatives for Cold-start Problem. In Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019, Wenwu Zhu, Dacheng Tao, Xueqi Cheng, Peng Cui, Elke A. Rundensteiner, David Carmel, Qi He, and Jeffrey Xu Yu (Eds.). ACM, 2605–2613. https://doi.org/10.1145/3357384.3357813