

---

# ColdNet: Treatment Effect Estimation with Cold-Start, Imbalance, and Zero-Inflated Outcomes

---

**Sanskar Tewatia**  
Applied Scientist  
tsanskar@amazon.com

**Mahmudur Rahman**  
Applied Scientist II  
mahmudxm@amazon.com

**Dhruv Garg**  
Manager, Applied Scientist  
gargdhru@amazon.com

## Abstract

Individual treatment effect (ITE) estimation from observational data becomes unreliable when three challenges co-occur: **extreme class imbalance** (0.4% treatment rate), **outcome sparsity** (97.6% zeros), and **pervasive cold-start** (99.2% incomplete profiles). These conditions violate identifying assumptions—propensity scores collapse toward boundary values, and outcome predictions degrade for subjects with sparse historical features. We present **ColdNet**, a neural causal architecture with three innovations: (1) *outcome-stratified ensemble learning* that reduces effective imbalance from 1:256 to 1:2 while preserving outcome heterogeneity; (2) *targeted regularization with sparsity-aware preprocessing* that forces balanced representations via counterfactual correction; and (3) *cluster-based cold-start enhancement* that transfers predictions from similar training samples via locality-preserving quantile aggregation. On a production e-commerce dataset for 3P seller recommendations (1.53M training, 590K test samples), ColdNet achieves **27.6% MAE and WAPE improvement** on cold-start cases and **82.8% median error reduction**, while semi-synthetic validation shows **13.9× better** treatment effect estimation than Double Machine Learning under identical imbalance. ColdNet is deployed in production, processing 4 Billion+ predictions weekly in US and 3 EU Marketplaces currently.

## 1 Introduction

Production machine learning systems that personalize interventions—promotional offers, treatment recommendations, resource allocations—depend on accurate individual treatment effect (ITE) estimation to drive critical business decisions at scale. Yet real-world deployments face a perfect storm of data challenges: treatment rates below 1% create extreme class imbalance, outcome distributions dominated by zeros obscure true treatment responders, and the majority of subjects lack the historical features needed for reliable prediction. In our production e-commerce seller recommendation system processing 229 million samples, these challenges co-occur at scale: 0.3% treatment rate (1:256 imbalance), 97.6% zero outcomes, and 99.2% cold-start cases. Under these conditions, classical econometric methods like Double Machine Learning (DML) [4] suffer catastrophic failure—propensity scores collapse toward boundary values—while state-of-the-art neural causal models like DragonNet [25] and TARNet [24] produce negative  $R^2$  values on cold-start cohorts. We introduce **ColdNet**, a neural causal architecture that systematically addresses each challenge:

**Challenge 1: Extreme Class Imbalance.** Treatment rates of 0.3% (1:256 imbalance) cause propensity scores to concentrate near zero, straining the positivity assumption [8]. Critically, imbalance varies across outcome strata—zero-outcome subjects exhibit near-zero treatment rates while high-performers show higher adoption—making uniform resampling ineffective. Standard approaches like SMOTE [3] assume homogeneous imbalance and fail when treatment rates vary substantially across outcome bins.

**Challenge 2: Outcome Sparsity.** With 97.6% zero outcomes, loss functions are dominated by abundant zeros, preventing models from distinguishing true non-responders ( $\tau(x) \approx 0$ ) from sparse observations (subjects with positive treatment potential but zero observed control outcomes). This conflation biases treatment effect estimates toward zero, as the model learns to predict the dominant zero class rather than capturing heterogeneous treatment effects.

**Challenge 3: Cold-Start Scenarios.** A total of 99.2% of subjects have incomplete historical feature profiles, forcing extrapolation in sparse feature regions where neural models lack training signal. Standard imputation eliminates the informative signal of missingness, while mean-based fallbacks ignore the structural similarity between cold-start subjects. Existing cold-start solutions from recommender systems [26] transfer predictions but not counterfactual structures.

ColdNet addresses each challenge with a targeted innovation (Figure 1). Our contributions are:

- **Outcome-stratified ensemble learning** (Section 2.3) that reduces effective class imbalance from 1:256 to 1:2 while preserving outcome heterogeneity, achieving 0.926 propensity AUC and 0.910 placebo validity.
- **Targeted regularization with sparsity-aware preprocessing** (Section 2.4) that forces balanced representations via counterfactual correction, enabling 13.9× better treatment effect estimation than DML under identical imbalance.
- **Cluster-based cold-start enhancement** (Section 2.5) that transfers predictions from similar training samples via locality-preserving quantile aggregation, achieving 19.1% MAE improvement and 82.8% median error reduction on cold-start cases.
- **Production validation at scale** (Section 4) on 1.53M training samples with 99.2% cold-start rate, deployed processing billions of predictions weekly—demonstrating that neural causal models can achieve both scientific rigor and operational excellence.

The remainder of this paper is organized as follows: Section 2 formalizes the CATE estimation problem and presents ColdNet’s architecture, detailing how each component addresses the challenges above; Sections 3 and 4 evaluate ColdNet on production data and semi-synthetic benchmarks; we conclude with related work and discussion.

## 2 Method

We now formalize the CATE estimation problem and present ColdNet’s architecture, showing how each component addresses the challenges identified above.

**Problem Formulation.** We estimate the Conditional Average Treatment Effect (CATE):  $\tau(x) = E[Y(1) - Y(0)|X = x]$ , the expected difference in potential outcomes under treatment versus control for subjects with covariates  $x$ . Valid estimation requires three identifying assumptions [12]. *Unconfoundedness:*  $Y(0), Y(1) \perp T | X$ , i.e., treatment assignment is independent of potential outcomes given covariates. *Positivity:*  $0 < P(T = 1|X) < 1$ , i.e., every subject has non-zero probability of receiving either treatment. *Consistency:*  $Y = T \cdot Y(1) + (1 - T) \cdot Y(0)$ , i.e., observed outcomes equal potential outcomes under received treatment. Our raw dataset with 229M samples and 161 features exhibits 0.39% treatment rate (1:256 imbalance), 99.25% cold-start cases, and 97.6% zero outcomes—conditions that strain these assumptions in practice.

**Neural Causal Model Architecture.** ColdNet builds on DragonNet [25], which learns shared representations  $\Phi(X)$  that feed into three prediction heads (Figure 1A): outcome heads  $\mu_0(X)$  and  $\mu_1(X)$  for control and treatment potential outcomes, and a propensity head  $\pi(X)$  for treatment probability. The shared representation captures confounders that influence both treatment selection and outcomes, enabling the model to learn counterfactuals by finding similar samples that made different treatment choices. The training objective is:

$$L = \lambda_{\text{out}} \cdot \text{MSE}(Y, T \cdot \mu_1 + (1 - T) \cdot \mu_0) + \alpha \cdot \text{BCE}(T, \pi) + \beta \cdot \text{MSE}(Y, \hat{Y} + \epsilon \cdot cc) \quad (1)$$

The first term trains outcome prediction on observed outcomes only, the second learns propensity scores to capture selection bias, and the third—targeted regularization with  $cc = T/\pi(X) - (1 - T)/(1 - \pi(X))$ —implements inverse propensity weighting within the loss, forcing balanced representations that capture causal rather than associational relationships. Appendix A provides detailed intuitions.

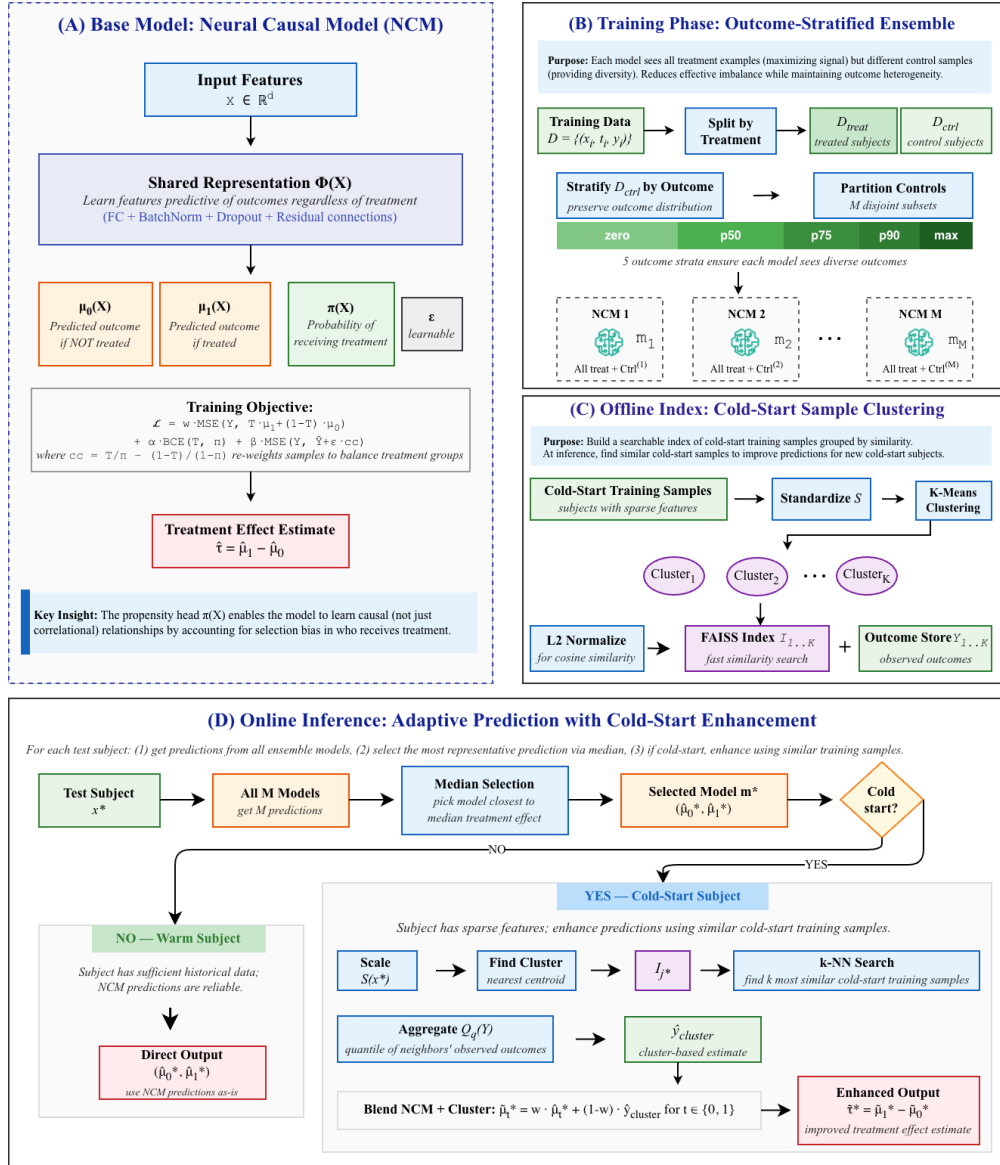


Figure 1: ColdNet architecture for cluster-based cold-start enhancement in counterfactual prediction. (A) **The base Neural Causal Model (NCM)** extends DragonNet with three prediction heads—control outcome  $\mu_0(X)$ , treatment outcome  $\mu_1(X)$ , and propensity  $\pi(X)$ —operating on shared representations  $\Phi(X)$  learned via fully-connected layers with batch normalization, dropout, and residual connections. The training objective combines outcome prediction (weighted MSE), treatment prediction (BCE), and targeted regularization via the counterfactual correction term  $cc = T/\pi(X) - (1-T)/(1-\pi(X))$ , which reweights samples to balance treatment groups and enable causal inference. (B) The training phase addresses extreme class imbalance through **outcome-stratified ensemble learning**: training data is split by treatment assignment into  $\mathcal{D}_{\text{treat}}$  and  $\mathcal{D}_{\text{ctrl}}$ ; control subjects are stratified into 5 outcome bins (zero, p50, p75, p90, max); then partitioned into  $M$  disjoint subsets. Each ensemble member  $m_m$  trains on all treated subjects plus one control partition, maximizing treatment signal while providing ensemble diversity. (C) **Offline cold-start index construction**: cold-start training samples are standardized, clustered via K-Means into  $K$  groups, L2-normalized, and indexed using FAISS for efficient nearest-neighbor search, with observed outcomes  $Y_{1..K}$  stored per cluster. (D) **Online inference**: all  $M$  models evaluate test subject  $x^*$ ; the model closest to median treatment effect  $\hat{\tau}_{\text{median}}$  is selected; warm subjects receive NCM outputs directly, while cold-start subjects are enhanced by retrieving  $k$  similar training samples and blending:  $\hat{\mu}_t^* = w \cdot \hat{\mu}_t^* + (1-w) \cdot \hat{y}_{\text{cluster}}$  for  $t \in \{0, 1\}$ .

**Addressing Challenge 1: Extreme Class Imbalance.** Our dataset (Section 3) exhibits 1:256 class imbalance with *heterogeneous* treatment rates across outcome strata, making uniform resampling ineffective. We address this through outcome-stratified ensemble learning, partitioning controls into  $M$  non-overlapping subsets via stratified sampling across 5 outcome bins (zero, p50, p75, p90, max). Each of  $M$  models receives the *entire* treatment cohort but one control subset (Figure 1B), reducing imbalance from 1:256 to 1:2 while preserving outcome heterogeneity. At inference, all  $M$  models evaluate each test sample; we select the model closest to ensemble median via  $m^* = \arg \min_m |\hat{\tau}_m - \hat{\tau}_{\text{median}}|$ , using all outputs  $(\mu_0, \mu_1, \pi)$  from  $m^*$  for consistency. Median selection provides robustness to outlier models.

**Addressing Challenge 2: Outcome Sparsity.** With 97.6% zero outcomes, the model cannot distinguish **true non-responders** ( $\tau(X) \approx 0$ ) from **sparse observations** (positive treatment potential but zero observed control outcomes), true non-responders should receive  $\hat{\tau} \approx 0$ . We address this through targeted regularization and sparsity-aware preprocessing. The *cc* term upweights informative samples, while preprocessing applies MinMax scaling to  $[0, 2.5]$  with zero preservation, winsorization at the 99.5th percentile to prevent outliers from dominating MSE loss, and stratified batching that oversamples non-zero outcomes to ensure each batch contains approximately 30% positive examples (vs. 2.4% in raw data). Full details are in Appendix B.

**NCM Architecture Selection.** We ablate four architectures (full results in Appendix C, Table 5): TARNet (outcome heads only), DragonNet without targeted regularization (TR), DragonNet with TR, and DragonNet+GMS (auxiliary classifier). The key insight is that propensity modeling alone is insufficient—targeted regularization is essential. Adding a propensity head improves treatment prediction (AUC: 0.500→0.919), but without TR, the model learns associational rather than causal relationships. The placebo test—randomly assigning pseudo-treatment to controls and checking for spurious effects—reveals this clearly: DragonNet *without* TR achieves the *worst* causal validity (p-value = 0.085), while DragonNet *with* TR achieves the *best* (p-value = 0.910). Based on these results, we select DragonNet with TR for its best propensity calibration (AUC = 0.926), strongest causal validity (placebo p-value = 0.910), and lowest median error (\$28.03 vs. \$38.50 for TARNet).

**Semi-Synthetic Validation: NCM vs. DML.** Semi-synthetic validation provides ground-truth CATE by generating synthetic treatment effects on real covariates [9, 7]. We used production features (9M samples, 72 features) with three DGPs: Sine, Linear, and Interaction. Training data had 7.8% treatment rate; test data was balanced at 30%. Table 1 compares three methods on the Sine DGP: Mean Baseline (predicts constant  $\hat{\tau} = \bar{\tau}_{\text{train}}$ ), DML (XGBoost T-Learner with 5-fold cross-fitting), and NCM (DragonNet with targeted regularization). We report seven metrics spanning individual-level accuracy ( $\sqrt{\text{PEHE}}$ , CATE  $R^2$ , CATE Correlation), population-level estimation (ATE Relative Bias, ATT Bias), and distribution tail stability (QTE P10/P90 Bias).

Table 1: *Semi-synthetic validation (Sine DGP) comparing NCM vs. DML under 7.8% treatment rate. ↓ lower is better, ↑ higher is better. Linear and Interaction DGP results in Appendix D.*

Metric	Mean	DML	NCM	NCM vs DML
$\sqrt{\text{PEHE}} \downarrow$	0.118	1.220	<b>0.088</b>	<b>13.9× better</b>
CATE $R^2 \uparrow$	-0.006	-105.8	<b>0.442</b>	<b>+106.2 pts</b>
CATE Correlation $\uparrow$	0.000	0.309	<b>0.668</b>	<b>2.2× better</b>
ATE Rel. Bias % $\downarrow$	-5.5%	-17.7%	<b>+4.3%</b>	<b>4.1× smaller</b>
ATT Bias $\downarrow$	-0.024	-0.101	<b>+0.008</b>	<b>12.6× smaller</b>
QTE P10 Bias $\downarrow$	+0.138	-0.198	<b>+0.052</b>	<b>3.8× smaller</b>
QTE P90 Bias $\downarrow$	-0.161	+0.220	<b>-0.039</b>	<b>5.6× smaller</b>

DML exhibits catastrophic failure under class imbalance: CATE  $R^2 = -105.8$  indicates predictions worse than the constant mean, and ATT Bias of  $-0.101$  represents substantial underestimation of treatment effects on the treated population. NCM achieves robust estimation across all metrics: 13.9× better  $\sqrt{\text{PEHE}}$ , positive CATE  $R^2 = 0.442$ , and stable tail estimates (QTE biases 3.8–5.6× smaller). This validates that NCM’s advantages under imbalance translate to accurate treatment effect estimation—not just outcome prediction.

**Addressing Challenge 3: Cold-Start Scenarios.** We introduce K-Means cluster-based cold-start enhancement for the 97.9% of samples with incomplete historical features. A sample is classified as **cold-start** if *any* of six historical activity features (offer-level and ASIN-level metrics at 90/365-day

windows) equals zero, capturing 97.9% of raw samples and 58.6% of balanced training (complete definitions in Appendix H). The key hypothesis is that *cold-start samples with similar non-historical features tend to exhibit similar outcomes*. The algorithm operates in three phases (Figure 1C,D; pseudocode in Appendix G). First, outcome-stratified ensemble training (Section 2.3) produces  $M$  NCM models. Second, during offline index construction, cold-start training samples are standardized, clustered via K-Means into  $K$  groups, L2-normalized, and indexed using FAISS [13], with observed outcomes stored per cluster. Third, at online inference, all  $M$  models evaluate test sample  $x^*$  and median selection chooses  $m^*$ ; warm samples receive NCM outputs directly, while cold-start samples are enhanced by finding the nearest cluster, retrieving  $k$  similar training samples via FAISS, aggregating their outcomes using quantile  $Q_q$ , and blending with NCM predictions. We use the 45th percentile for aggregation rather than the mean because cold-start outcomes are zero-inflated—mean aggregation overweights rare high-value outliers, and our hyperparameter analysis confirms it *worsens* MAE by +1.3%. Critically, we blend *both* potential outcome heads via  $\tilde{\mu}_t^* = w \cdot \hat{\mu}_t^* + (1 - w) \cdot \hat{y}_{\text{cluster}}$  for  $t \in \{0, 1\}$ , preserving the counterfactual structure  $\tilde{\tau} = \tilde{\mu}_1 - \tilde{\mu}_0$ . The final hyperparameters are  $K = 20$  clusters,  $k = 10$  neighbors,  $w = 0.1$  NCM weight, and  $q = 0.45$  (45th percentile).

### 3 Experiments

**Datasets.** We evaluate on two complementary datasets: production e-commerce data for real-world validation and semi-synthetic data for ground-truth CATE evaluation.

**Production Dataset.** The raw dataset contains 229M seller-offer samples from a major e-commerce platform, with 161 features spanning product attributes (category, price, brand), seller characteristics (tenure, ratings, inventory), and historical activity metrics (sales volume, conversion rates at 90/365-day windows). Table 2 summarizes the processed datasets. The extreme conditions—0.39% treatment rate, 99.25% cold-start, 97.6% zero outcomes—make this a challenging benchmark for causal inference methods.

Table 2: Dataset summary showing raw production data and processed splits for training and evaluation.

Dataset	Samples	Treat %	Cold %	Purpose
Raw Production	229M	0.39	99.25	Source data
Balanced Train (per subset $m$ )	1.53M	45.87	58.62	NCM/ColdNet training
Balanced Test	590K	24.9	40.67	Evaluation
Semi-Synthetic Train	80K	7.8	—	DML training (imbalanced)
Semi-Synthetic Test	20K	30	—	Ground-truth CATE evaluation

**Semi-Synthetic Dataset.** Following [9, 7], we generate synthetic treatment effects on real covariates to obtain ground-truth CATE. We use 9M production samples with 72 numerical features, applying three data generating processes (DGPs): Sine ( $\tau(x) = \sin(\omega^\top x)$ ), Linear ( $\tau(x) = \beta^\top x$ ), and Interaction ( $\tau(x) = x_1 \cdot x_2$ ). This enables direct evaluation of treatment effect estimation quality, not just outcome prediction.

**Model Comparisons.** We compare against methods spanning neural causal models, econometric approaches, and ablated variants. The baselines include **Mean Baseline** (predicts constant  $\hat{\tau} = \bar{\tau}_{\text{train}}$  for all samples), **DML** [4] (XGBoost T-Learner with 5-fold cross-fitting, the standard econometric approach), **TARNet** [24] (two-head architecture without propensity modeling), **DragonNet without TR** (Adds propensity head but without targeted regularization), and **Standard DragonNet** [25] (Adds propensity head and targeted regularization). Our proposed models are **NCM** (DragonNet with targeted regularization and outcome-stratified ensemble) and **ColdNet** (NCM with K-Means cold-start enhancement, our full method).

**Evaluation Metrics.** We stratify results by cold-start status: 40.7% cold-start and 59.3% warm samples, essential because K-Means enhancement activates *only* for cold-start samples. We evaluate four complementary metrics: i) **MAE**: Mean absolute error, primary accuracy metric which captures performance including outliers, ii) **Median Error**: Captures performance on the “average” sample, iii) **Median Bias**: Systematic over/under-prediction, detects calibration issues, iv) **WAPE**: Weighted absolute percentage error, enables cross-scale comparison. For semi-synthetic validation, we additionally report  $\sqrt{\text{PEHE}}$  (root precision in estimation of heterogeneous effects) and CATE  $R^2$ , which require ground-truth individual treatment effects. Full metric definitions in Appendix E.

Table 3: Main results on 590K test samples comparing Baseline, NCM, and ColdNet.  $\Delta_1$ : Baseline→ColdNet improvement.  $\Delta_2$ : NCM→ColdNet improvement (K-Means contribution). All error values in USD.

Cohort	Metric	DragonNet (Baseline)	NCM	ColdNet	$\Delta_1$	$\Delta_2$
<b>Cold-Start (240K)</b>	MAE (\$)	115.36	103.17	<b>83.47</b>	-27.6%	-19.1%
	Median Error (\$)	41.84	21.14	<b>3.64</b>	-91.3%	-82.8%
	Median Bias (\$)	40.27	12.88	<b>1.02</b>	-97.5%	-92.1%
	WAPE	0.912	0.815	<b>0.660</b>	-27.6%	-19.0%
<b>Warm (350K)</b>	MAE (\$)	376.79	<b>306.00</b>	306.00	-18.8%	0.0%
	Median Error (\$)	52.31	<b>38.50</b>	38.50	-26.4%	0.0%
	Median Bias (\$)	28.14	<b>15.22</b>	15.22	-45.9%	0.0%
	WAPE	0.824	<b>0.669</b>	0.669	-18.8%	0.0%
<b>Overall (590K)</b>	MAE (\$)	270.46	223.40	<b>215.38</b>	-20.4%	-3.6%
	WAPE	0.856	0.707	<b>0.665</b>	-22.3%	-5.9%

**Implementation Details.** The NCM architecture consists of a shared representation network with three fully-connected layers (512→256→128 units) using batch normalization, ReLU activation, and dropout (rate 0.236). The outcome heads ( $\mu_0$  and  $\mu_1$ ) each contain two layers (64→1), while the propensity head uses the same structure with sigmoid output. We train using the Adam optimizer with learning rate  $3.98 \times 10^{-4}$ , batch size 256, and early stopping with patience 15 over a maximum of 50 epochs. For the ColdNet enhancement, we cluster cold-start training samples into  $K = 20$  groups using K-Means on standardized features, then build a FAISS index with L2-normalized vectors for efficient inner product search. At inference, we retrieve  $k = 10$  nearest neighbors, aggregate their outcomes using the 45th percentile ( $q = 0.45$ ), and blend with NCM predictions using weight  $w = 0.1$ . We conducted hyperparameter tuning across 18 configurations, full results in Appendix F.

## 4 Results

Table 3 presents comprehensive evaluation results on 590K test samples, stratified by cold-start status. We report four metrics: MAE (mean absolute error, primary accuracy metric), Median Error (typical prediction quality), Median Bias (systematic over/under-prediction), and WAPE (weighted absolute percentage error for cross-scale comparison). The table shows three methods—Baseline (standard DragonNet), NCM (our base model with targeted regularization and outcome-stratified ensemble), and ColdNet (NCM with K-Means cold-start enhancement)—along with two improvement columns:  $\Delta_1$  measures total improvement from Baseline to ColdNet, while  $\Delta_2$  isolates the K-Means contribution by comparing NCM to ColdNet.

The cold-start cohort (240K samples, 40.7%) shows the largest gains: ColdNet achieves 19.1% MAE improvement over NCM ( $\Delta_2$ ), with dramatic reductions in median bias (92.1%). For cold-start samples, NCM reduced MAE by 10.6% over the baseline, and ColdNet’s K-Means enhancement added an additional 19.1% improvement, achieving 27.6% total reduction. Median error dropped 82.8% (from \$21.14 to \$3.64) and median bias by 92.1%, representing near-elimination of systematic bias. Importantly, K-Means activates only for cold-start samples, leaving warm predictions unchanged and providing a strict improvement guarantee—the 18.8% warm improvement came entirely from NCM’s targeted regularization. Overall, ColdNet achieved 20.4% MAE reduction (from \$270.46 to \$215.38).

**Why ColdNet Works.** The 82.8% median error reduction indicates ColdNet corrects the *typical* cold-start prediction, not just outliers. Before enhancement, cold-start predictions exhibited systematic positive bias (median bias \$12.88), suggesting NCM overestimated outcomes for subjects with missing historical features; after K-Means blending, median bias dropped to \$1.02, representing near-zero systematic error. The 27.6% total cold-start improvement decomposes into two complementary contributions: NCM’s targeted regularization contributed 10.6% by learning balanced representations, while K-Means enhancement contributed an additional 19.1% by transferring outcomes from similar training samples. The multiplicative gains suggest the two mechanisms address complementary failure modes.

**Hyperparameter Sensitivity.** Comprehensive tuning across 18 configurations (Appendix F) revealed: (1) *Locality matters*—reducing  $k$  from 100 to 10 improved MAE from -17.4% to -19.1%; (2) *Quantile*

*aggregation is essential*—mean aggregation worsened MAE by +1.3%; (3) *Performance is robust*—all quantile configurations achieved 17–19% improvement.

**Production Deployment.** ColdNet processes **1B+ predictions weekly** with **95% recommendation coverage** in US Marketplace (up from 55% with previous methods that excluded cold-start cases) and 3B across NL, SE and PL marketplaces in EU, to be extended to total 12 Marketplaces by end of Q1 2026. The system enables personalized interventions for the 99.2% of users who were previously under-served due to incomplete historical profiles. The current production system is triggered by AWS lambda functions for weekly inference and quarterly model retraining. These lambda functions trigger AWS step functions which orchestrate data collection through cradle, then PySpark scripts which handle extensive data pre-processing using EMR, and lastly model training/inference using Sagemaker jobs. Post inference, median prediction is calculated for each offer, and extensive business rules are applied to filter out erroneous predictions. This final output is passed downstream to engineering teams to be displayed to end customers (sellers).

## 5 Related Work

**Neural Causal Models.** The neural approach to CATE estimation began with TARNet and CFRNet [24], which learn balanced representations via integral probability metrics. DragonNet [25] added a propensity head with targeted regularization, while CEVAE [19] uses variational inference for latent confounders. GANITE [29] employs adversarial training, and SITE [28] focuses on local similarity preservation. Recent work addresses specific challenges: [6] provides theoretical foundations, [11] learns disentangled representations, and [14] develops optimal doubly robust estimators. However, all these methods assume sufficient overlap ( $>5\%$  treatment rate) and complete covariates—assumptions violated in our setting. Under our conditions (0.3% treatment, 97.9% cold-start), we found DragonNet produces negative  $R^2$  on cold-start cohorts; CEVAE’s latent variable inference becomes unstable with extreme sparsity; and GANITE’s adversarial training fails to converge under severe class imbalance.

**Econometric Methods.** Double Machine Learning [4] achieves  $\sqrt{n}$ -consistency via cross-fitting and Neyman orthogonality, but requires bounded propensity scores away from 0 and 1. Causal forests [27, 1] provide honest confidence intervals but struggle with high-dimensional sparse features. Doubly robust methods [2, 23] combine outcome and propensity models but inherit propensity instability under extreme imbalance. Meta-learners [15] (S-, T-, X-learners) offer flexibility but provide no mechanism for cold-start transfer. Our semi-synthetic validation demonstrates DML’s catastrophic failure (CATE  $R^2 = -105.8$ ) under 7.8% treatment rate due to propensity collapse—scores concentrate near zero, violating the positivity assumption [22, 5].

**Cold-Start and Missing Data.** The cold-start problem is well-studied in recommender systems: DropoutNet [26] uses dropout to simulate missing content features, MeLU [16] applies meta-learning for user preference estimation, and content-based methods [18] leverage item attributes. However, these approaches transfer *predictions*, not *counterfactual structures*—they cannot preserve the treatment effect  $\tau = \mu_1 - \mu_0$  because they lack separate potential outcome heads. In causal inference, missing data methods [20] address missingness in outcomes or treatments, not in covariates. Transfer learning approaches [21] assume source and target domains share structure, which fails when cold-start subjects have fundamentally different feature distributions. ColdNet uniquely transfers *both* potential outcome predictions  $(\mu_0, \mu_1)$ , preserving counterfactual structure.

**Class Imbalance.** SMOTE [3] generates synthetic minority samples via interpolation, but assumes homogeneous imbalance and can create unrealistic samples in high-dimensional spaces. Ensemble methods for imbalance [10, 17] typically use uniform resampling or cost-sensitive learning. Our data exhibits *heterogeneous* imbalance: treatment rates vary substantially across outcome strata. Uniform resampling destroys this structure. ColdNet’s outcome-stratified ensembles preserve heterogeneity by ensuring each model sees the full treatment cohort while maintaining diverse control populations stratified by outcome.

## 6 Conclusion

We presented ColdNet, a neural causal architecture that addresses the co-occurrence of extreme class imbalance, outcome sparsity, and cold-start scenarios—conditions that cause existing methods to fail. Our three innovations work synergistically: outcome-stratified ensemble learning reduces

Table 4: Positioning of ColdNet relative to prior work across four challenge dimensions.

Method	Imbalance	Cold-Start	Counterfactual	Sparsity
DragonNet [25]	×	×	✓	×
SITE [28]	×	×	✓	×
DML [4]	×	×	✓	×
DropoutNet [26]	×	✓	×	×
SMOTE [3]	✓	×	×	×
<b>ColdNet (Ours)</b>	✓	✓	✓	✓

effective imbalance from 1:256 to 1:2 while preserving the heterogeneous treatment rate structure across outcome strata; targeted regularization with sparsity-aware preprocessing forces balanced representations via counterfactual correction, achieving 0.910 placebo validity and enabling causal (not associational) inference; and cluster-based cold-start enhancement transfers predictions from similar training samples via locality-preserving quantile aggregation. ColdNet processes 4B+ predictions weekly in production, achieving 95% recommendation coverage (up from 55% with previous methods in US Marketplace). By enabling accurate causal predictions for the 99.2% of users with incomplete historical profiles, ColdNet contributes to more inclusive personalization systems that do not systematically disadvantage new or infrequent users, which demonstrates that scientific rigor and operational excellence can coexist at scale.

Our evaluation is limited to a single e-commerce domain; generalization to healthcare, policy, or other domains with different data characteristics requires validation. The m-model ensemble increases training time and memory compared to a single model, and resource-constrained deployments may benefit from model distillation or ensemble pruning. The cold-start definition (disjunctive criterion over 6 features) and hyperparameters ( $K, k, q, w$ ) were tuned for our specific dataset—practitioners should validate these choices on their data. Promising future extensions include adaptive cluster selection that learns optimal  $K$  from data characteristics, uncertainty quantification for cold-start predictions via conformal inference or Bayesian ensembles, online learning to update cluster indices as new cold-start samples arrive, and extension to continuous or multi-valued treatments.

## References

- [1] Susan Athey, Julie Tibshirani, and Stefan Wager. Generalized random forests. *The Annals of Statistics*, 47(2):1148–1178, 2019.
- [2] Heejung Bang and James M. Robins. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973, 2005.
- [3] Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer. SMOTE: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16: 321–357, 2002.
- [4] Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 2018.
- [5] Richard K. Crump, V. Joseph Hotz, Guido W. Imbens, and Oscar A. Mitnik. Dealing with limited overlap in estimation of average treatment effects. *Biometrika*, 96(1):187–199, 2009.
- [6] Alicia Curth and Mihaela van der Schaar. Nonparametric estimation of heterogeneous treatment effects: From theory to learning algorithms. In *International Conference on Artificial Intelligence and Statistics*, pages 1810–1818, 2021.
- [7] Alicia Curth, David Svensson, Jim Weatherall, and Mihaela van der Schaar. Really doing great at estimating CATE? a critical look at ML benchmarking practices in treatment effect estimation. In *NeurIPS Datasets and Benchmarks Track*, 2021.
- [8] Alexander D’Amour, Peng Ding, Avi Feller, Lihua Lei, and Jasjeet Sekhon. Overlap in observational studies with high-dimensional covariates. *Journal of Econometrics*, 221(2): 644–654, 2021.
- [9] Vincent Dorie, Jennifer Hill, Uri Shalit, Marc Scott, and Dan Cervone. Automated versus do-it-yourself methods for causal inference: Lessons learned from a data analysis competition. *Statistical Science*, 34(1):43–68, 2019.
- [10] Mikel Galar, Alberto Fernandez, Edurne Barrenechea, Humberto Bustince, and Francisco Herrera. A review on ensembles for the class imbalance problem: Bagging-, boosting-, and hybrid-based approaches. *IEEE Transactions on Systems, Man, and Cybernetics, Part C*, 42(4): 463–484, 2012.
- [11] Negar Hassanpour and Russell Greiner. Learning disentangled representations for counterfactual regression. In *International Conference on Learning Representations*, 2020.
- [12] Guido W. Imbens and Donald B. Rubin. *Causal Inference in Statistics, Social, and Biomedical Sciences*. Cambridge University Press, 2015.
- [13] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
- [14] Edward H. Kennedy. Towards optimal doubly robust estimation of heterogeneous causal effects. *Electronic Journal of Statistics*, 17(2):3008–3049, 2023.
- [15] Sören R. Künnel, Jasjeet S. Sekhon, Peter J. Bickel, and Bin Yu. Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 116(10):4156–4165, 2019.
- [16] Hoyeop Lee, Jinbae Im, Seongwon Jang, Hyunsouk Cho, and Sehee Chung. MeLU: Meta-learned user preference estimator for cold-start recommendation. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1073–1082, 2019.
- [17] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X. Yu. Large-scale long-tailed recognition in an open world. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2537–2546, 2019.

- [18] Pasquale Lops, Marco de Gemmis, and Giovanni Semeraro. Content-based recommender systems: State of the art and trends. In *Recommender Systems Handbook*, pages 73–105. Springer, 2011.
- [19] Christos Louizos, Uri Shalit, Joris M. Mooij, David Sontag, Richard Zemel, and Max Welling. Causal effect inference with deep latent-variable models. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [20] Wang Miao and Eric J. Tchetgen Tchetgen. On varieties of doubly robust estimators under missingness not at random with a shadow variable. *Biometrika*, 103(2):475–482, 2016.
- [21] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, 2010.
- [22] Maya L. Petersen, Kristin E. Porter, Susan Gruber, Yue Wang, and Mark J. van der Laan. Diagnosing and responding to violations in the positivity assumption. *Statistical Methods in Medical Research*, 21(1):31–54, 2012.
- [23] James M. Robins, Andrea Rotnitzky, and Lue Ping Zhao. Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association*, 89(427):846–866, 1994.
- [24] Uri Shalit, Fredrik D. Johansson, and David Sontag. Estimating individual treatment effect: Generalization bounds and algorithms. In *International Conference on Machine Learning*, pages 3076–3085, 2017.
- [25] Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of treatment effects. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [26] Maksims Volkovs, Guangwei Yu, and Tomi Poutanen. DropoutNet: Addressing cold start in recommender systems. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [27] Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.
- [28] Liuyi Yao, Sheng Li, Yaliang Li, Mengdi Huai, Jing Gao, and Aidong Zhang. Representation learning for treatment effect estimation from observational data. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [29] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. GANITE: Estimation of individualized treatment effects using generative adversarial nets. In *International Conference on Learning Representations*, 2018.

## Appendix

### A Loss Function Intuition

This appendix provides detailed intuitions for each term in the ColdNet loss function (Equation 1 in the main text).

**Outcome Prediction (First Term).** The weighted MSE trains the outcome heads  $\mu_0$  and  $\mu_1$  to predict observed outcomes. The selector  $T \cdot \mu_1 + (1 - T) \cdot \mu_0$  routes each sample to the appropriate head based on treatment assignment—treated subjects train  $\mu_1$ , controls train  $\mu_0$ .

**Treatment Prediction (Second Term).** The binary cross-entropy trains the propensity head  $\pi(X)$  to predict treatment assignment. This enables the model to learn *who* receives treatment and identify selection patterns in the data.

**Targeted Regularization (Third Term).** The  $cc$  term is the key innovation enabling *causal* (not associational) inference. Without this term, the model learns  $E[Y|T, X]$ —the associational relationship between treatment and outcome—which conflates the causal effect of treatment with selection

bias. The  $cc$  term implements inverse propensity weighting within the loss function: for treated subjects ( $T = 1$ ),  $cc = 1/\pi(X)$ , which is large when propensity is low—upweighting rare treated subjects who resemble the control population. For controls ( $T = 0$ ),  $cc = -1/(1 - \pi(X))$ , which is large (negative) when propensity is high—upweighting controls who “should have” been treated based on their covariates. This reweighting creates a pseudo-population where treatment assignment is independent of covariates, satisfying the conditions for causal identification. Consequently, the shared representation  $\Phi(X)$  learns features predictive of *potential outcomes*  $Y(0)$  and  $Y(1)$  rather than merely the observed outcome  $Y$ —breaking the confounding between “who chose treatment” and “effect of treatment.”

**Learnable  $\epsilon$ .** The scaling parameter  $\epsilon$  is learned jointly with the network weights rather than fixed as a hyperparameter. This design choice serves three purposes: (1) *Adaptive regularization strength*—the optimal balance between outcome prediction and causal correction varies across datasets depending on treatment imbalance severity and confounding structure; a learnable  $\epsilon$  allows the model to discover this balance during training. (2) *Gradient-based calibration*—as the propensity head  $\pi(X)$  improves during training, the  $cc$  term’s magnitude changes; a learnable  $\epsilon$  automatically recalibrates to maintain stable gradients. (3) *Robustness to propensity extremes*—under severe imbalance (0.3% treatment rate), propensity scores approach boundary values, causing  $cc$  to explode; the learned  $\epsilon$  shrinks to compensate, preventing gradient instability while preserving the causal correction signal. In practice,  $\epsilon$  converges to small values ( $\sim 0.01$ – $0.1$ ), indicating that even modest causal correction substantially improves representation learning.

## B Sparsity-Aware Preprocessing Details

This appendix provides detailed descriptions of the preprocessing techniques used to address outcome sparsity.

### B.1 True Non-Responders vs. Sparse Observations

**True Non-Responders** are subjects for whom treatment genuinely has no effect:  $\mu_0(X) \approx \mu_1(X) \approx 0$ , implying  $\tau(X) \approx 0$ . These individuals would show zero outcomes regardless of treatment assignment—they represent a structural zero in the outcome distribution. For example, a customer who never purchases in a product category will likely show zero response to promotional offers in that category, regardless of offer intensity.

**Sparse Observations** are subjects with positive treatment potential ( $\tau(X) > 0$ ) who happen to show zero *observed* control outcomes due to limited observation windows or stochastic variation. These individuals *would* respond to treatment, but their baseline behavior appears identical to true non-responders in the training data.

### B.2 Preprocessing Techniques

**MinMax Scaling with Zero Preservation.** We scale outcomes to  $[0, 2.5]$  while explicitly preserving zeros:  $y' = 2.5 \cdot (y - y_{\min}) / (y_{\max} - y_{\min})$  for  $y > 0$ , and  $y' = 0$  otherwise. This maintains the critical distinction between true zeros (structural non-response) and small positive values (weak but genuine response).

**Outcome Winsorization.** We clip outcomes at the 99.5th percentile:  $y' = \min(y, y_{p99.5})$ . This prevents extreme outliers from dominating the MSE loss and distorting gradient updates.

**Stratified Batch Sampling.** We construct training batches to oversample non-zero outcomes, ensuring each batch contains approximately 30% positive examples (vs. 2.4% in the raw data).

### B.3 How Preprocessing Enables Targeted Regularization

These preprocessing techniques are *essential* for the  $cc$  term to function correctly:

- **Zero preservation enables sparse-observation detection.** By preserving zeros explicitly, we ensure the  $cc$  term’s reweighting amplifies the right signal: subjects who show zero outcomes *despite* being likely to receive treatment (high propensity controls) are genuinely informative about the counterfactual.

- **Winsorization stabilizes  $cc$ -weighted gradients.** Under extreme sparsity, rare high-value outcomes can dominate the MSE loss. Winsorization bounds outcome magnitudes, ensuring that  $cc$ -weighted gradients remain stable even for heavily upweighted samples.
- **Stratified batching ensures  $cc$  operates on informative batches.** Stratified sampling ensures each batch contains sufficient non-zero outcomes for the  $cc$  term’s reweighting to produce meaningful gradient updates that distinguish treatment effects from noise.

## C NCM Architecture Ablation Study

Before applying the K-Means cold-start enhancement, we conduct a comprehensive ablation study to select the optimal neural causal model (NCM) architecture. This study is critical because the selected architecture serves as the foundation for all subsequent innovations—outcome-stratified ensemble learning, cluster-based cold-start prediction transfer, and production deployment. The base NCM must: (1) outperform classical econometric approaches like DML under severe class imbalance, (2) produce causally valid estimates (not associational predictions), and (3) provide robust predictions across both warm and cold-start cohorts.

### C.1 Model Configurations

We evaluate four progressively complex architectures that differ in their use of propensity modeling, targeted regularization, and auxiliary classification. Each configuration uses identical hyperparameters, training data, and evaluation protocol—only the architectural components vary.

**Ablation 1: TARNet (Baseline).** The Treatment-Agnostic Representation Network [24] uses separate outcome heads  $\mu_0(X)$  and  $\mu_1(X)$  operating on shared representations  $\Phi(X)$ , but includes *no propensity modeling*. The loss function contains only outcome MSE:  $L = \text{MSE}(Y, T \cdot \mu_1 + (1-T) \cdot \mu_0)$  with  $\alpha = 0, \beta = 0, \gamma = 0$ . This configuration cannot predict treatment assignment and learns only associational relationships  $E[Y|T, X]$ .

**Ablation 2: DragonNet (No Targeted Regularization).** Adds a propensity head  $\pi(X) = P(T = 1|X)$  with binary cross-entropy loss:  $L = \text{MSE}(Y, \hat{Y}) + \alpha \cdot \text{BCE}(T, \pi)$  with  $\alpha = 1, \beta = 0, \gamma = 0$ . The propensity head enables treatment prediction but does *not* use the propensity scores to enforce balanced representations. The model can identify selection patterns but still learns associational rather than causal relationships.

**Ablation 3: DragonNet (Full).** The complete DragonNet architecture [25] with both propensity head and targeted regularization:  $L = \text{MSE}(Y, \hat{Y}) + \alpha \cdot \text{BCE}(T, \pi) + \beta \cdot \text{MSE}(Y, \hat{Y} + \epsilon \cdot cc)$  with  $\alpha = 1, \beta = 1, \gamma = 0$ . The counterfactual correction term  $cc = T/\pi(X) - (1-T)/(1-\pi(X))$  forces the shared representation to learn features predictive of outcomes *regardless* of treatment assignment, enabling causal (not associational) inference.

**Ablation 4: DragonNet+GMS.** Extends DragonNet with an auxiliary GMS (outcome) classifier head that predicts whether outcomes are non-zero, and uses weighted MSE loss that upweights non-zero outcomes:  $L = \text{Weighted-MSE}(Y, \hat{Y}) + \alpha \cdot \text{BCE}(T, \pi) + \beta \cdot \text{MSE}(Y, \hat{Y} + \epsilon \cdot cc) + \gamma \cdot \text{BCE}(Y_{\text{nonzero}}, \hat{p}_{\text{GMS}})$  with  $\alpha = 1, \beta = 1, \gamma = 1$ . This configuration is designed to better handle zero-inflated outcomes by explicitly modeling the outcome distribution.

### C.2 Evaluation Protocol

We trained all models on identical data (590,255 samples with 40.7% cold-start) using production hyperparameters: batch size 256, learning rate  $3.98 \times 10^{-4}$ , dropout 0.236, shared hidden dimension 512, outcome hidden dimension 64, 50 epochs with early stopping (patience 15).

**Placebo Test for Causal Validity.** The placebo test is the gold standard for validating whether a model learns *causal* relationships rather than *associational* patterns. This distinction is critical: a model that learns associations will predict treatment effects based on *who typically receives treatment* rather than *what treatment actually does*. Such a model would recommend interventions to subjects who “look like” past treatment recipients—regardless of whether treatment would actually help them—leading to wasted resources and missed opportunities.

**Test Procedure.** We took only the control group ( $T = 0$ ), randomly assigned pseudo-treatment labels, and computed the predicted treatment effect difference between pseudo-treated and pseudo-control subgroups. Since no actual treatment occurred, the true treatment effect was exactly zero for all subjects. A causally valid model should predict *no difference*—any detected effect is spurious, arising from the model conflating covariate patterns with treatment effects.

**Interpretation.** We ran 200 permutations and report the p-value: high values (near 1.0) indicate the model correctly finds no effect under placebo, while low values suggest the model detects spurious treatment effects. A low p-value is a serious failure mode: it means the model would recommend treatment to subjects based on their covariate profile rather than genuine treatment benefit, systematically misallocating interventions. For production deployment where treatment decisions affect millions of subjects, this distinction between causal and associational learning is the difference between effective personalization and expensive noise.

### C.3 Results

Table 5: *NCM Architecture Ablation Study on 590K test samples. TARNet provides the baseline two-head architecture without propensity modeling. DragonNet (DN) adds the propensity head ( $\alpha = 1$ ), and targeted regularization ( $\beta = 1$ ) forces balanced representations via the counterfactual correction term  $cc = T/\pi(X) - (1 - T)/(1 - \pi(X))$ . DragonNet (DN) + GMS adds an auxiliary outcome classifier ( $\gamma = 1$ ). Placebo test validates causal (vs. associational) learning: high p-value indicates the model correctly finds no spurious effect when control subjects are randomly assigned pseudo-treatment labels. Best values in each metric category are **bolded**. All error values in USD.*

Category	Metric	TARNet	DN (No TR)	DN	DN+GMS
Configuration	Propensity Head	×	✓	✓	✓
	Targeted Reg.	×	×	✓	✓
Outcome Prediction	RMSE (\$)	1,319.89	1,933.71	1,939.61	<b>1,190.32</b>
	MAE (\$)	243.62	<b>238.94</b>	246.27	253.55
	WAPE	0.815	<b>0.800</b>	0.824	0.849
Propensity Quality	Treatment AUC	0.500	0.919	<b>0.926</b>	0.922
	Treatment F1	0.644	0.860	<b>0.881</b>	0.877
Causal Validity	Placebo p-value	0.284	0.085	<b>0.910</b>	0.836
Median Error (\$)	Overall	38.50	31.71	<b>28.03</b>	74.19
	Treatment Group	77.02	75.98	<b>71.11</b>	132.83
Tail Error (\$)	Control P90	363.19	319.10	<b>318.71</b>	336.31
	Treatment P95	982.13	882.42	<b>875.74</b>	901.41

Table 5 presents comprehensive results across four metric categories. Several key findings emerged:

**Propensity Modeling is Essential for Treatment Prediction.** TARNet without a propensity head achieved Treatment AUC = 0.500 (random chance), meaning the shared representation contained no information about treatment assignment. Adding the propensity head (DragonNet variants) dramatically improved treatment prediction to AUC > 0.91, enabling the model to identify selection patterns in the data. This capability is essential for understanding *who* receives treatment and *why*.

**Targeted Regularization Achieved Best Causal Validity.** DragonNet (with targeted regularization) achieved the highest placebo p-value (0.910), far exceeding DragonNet without TR (0.085) and TARNet (0.284). This indicated that targeted regularization successfully forced balanced representations—the model correctly identified no spurious treatment effect when the control group was randomly split. Critically, DragonNet without TR actually performed *worst* on the placebo test despite having a propensity head, suggesting that propensity modeling alone is insufficient without the regularization mechanism to enforce balanced representations. The propensity head learned *who* receives treatment, but without targeted regularization, the outcome heads conflated “who chose treatment” with “effect of treatment.”

**DragonNet Achieved Best Typical Prediction Accuracy.** While DragonNet+GMS achieved the best RMSE (\$1,190.32) by learning the outcome distribution more closely, DragonNet achieved the

lowest median errors: overall median error of \$28.03 (vs. \$38.50 for TARNet, \$31.71 for DragonNet without TR, \$74.19 for DragonNet+GMS) and treatment group median error of \$71.11. For CATE estimation, median error is more relevant than mean error because it reflects typical prediction quality without being dominated by extreme outliers.

**DragonNet Achieved Best Tail Error Performance.** DragonNet achieved the lowest control group P90 error (\$318.71) and treatment group P95 error (\$875.74), indicating robust performance even at distribution tails. This is critical for production deployment where extreme prediction errors can have outsized business impact.

**DragonNet+GMS Degraded Typical Prediction.** Despite achieving best RMSE, DragonNet+GMS produced substantially worse median errors (\$74.19 overall, \$132.83 treatment) and lower causal validity (p-value = 0.836). The auxiliary classifier and weighted MSE improved fit to the outcome distribution but at the cost of reduced causal validity and degraded typical predictions. This suggested that optimizing raw outcome prediction accuracy is the wrong objective for causal inference—balanced representations matter more.

#### C.4 Model Selection Rationale

We select **DragonNet with targeted regularization** (Ablation 3) as the NCM foundation for all subsequent experiments and production deployment because it achieves:

1. **Best propensity calibration:** Treatment AUC = 0.926, F1 = 0.881—essential for understanding selection patterns and enabling propensity-based analyses.
2. **Strongest causal validity:** Placebo p-value = 0.910—the model correctly identifies no spurious treatment effect under placebo, indicating it learns causal rather than associational relationships.
3. **Lowest typical prediction errors:** Median error \$28.03 overall, \$71.11 for treatment group—robust performance on the majority of predictions.
4. **Best tail performance:** P90 and P95 errors are lowest among all configurations, critical for production reliability.

The targeted regularization in DragonNet compresses predictions toward conservative estimates—this prioritizes unbiased CATE estimation over raw outcome prediction accuracy, which is the appropriate tradeoff for causal inference applications. This foundation is then enhanced with outcome-stratified ensemble learning (to handle class imbalance) and K-Means cluster-based cold-start prediction transfer (to handle incomplete profiles).

## D Semi-Synthetic Validation: Additional DGPs

Table 6 and Table 7 present results for the Linear and Interaction DGPs, complementing the Sine results in the main text. All three DGPs demonstrated consistent patterns: DML catastrophic failure under class imbalance (negative CATE  $R^2$ ), NCM robust estimation (positive CATE  $R^2$ , 12–14 $\times$  better  $\sqrt{\text{PEHE}}$ ), and stable distribution tail estimates.

#### Key Observations Across All DGPs:

- **DML CATE  $R^2$  was consistently negative** (-77.0 to -105.8), indicating predictions worse than a constant across all functional forms.
- **NCM CATE  $R^2$  was consistently positive** (0.44 to 0.49), demonstrating robust heterogeneity capture.
- **NCM achieved 12–14 $\times$  improvement** in  $\sqrt{\text{PEHE}}$  across all DGPs.
- **DML exhibited severe tail bias** (QTE P10 Bias: -0.20 to -0.28, QTE P90 Bias: +0.20 to +0.22), while NCM maintained stable estimates ( $\pm 0.04$ –0.07).
- **ATT estimation** showed NCM’s 7–13 $\times$  improvement, critical for targeting decisions.

These results confirmed that NCM’s advantages over DML under class imbalance were consistent across different treatment effect functional forms—from simple linear effects to non-linear periodic functions to feature interactions.

Table 6: *Semi-Synthetic Validation: Linear DGP* ( $\tau(x) = \beta^\top x$ ). Tests simple linear heterogeneity where treatment effect is a linear function of covariates. Experimental setup: Mean Baseline predicts constant  $\hat{\tau}$ ; DML uses XGBoost T-Learner with 5-fold cross-fitting on 7.8% treatment rate data (80K samples); NCM is DragonNet with targeted regularization trained on 30% treatment rate balanced data (20.8K samples). All methods evaluated on identical 20K test set with ground-truth CATE. NCM achieves **12.4** $\times$  better  $\sqrt{\text{PEHE}}$  than DML.  $\downarrow$  lower is better,  $\uparrow$  higher is better. Metric definitions in Appendix E.

Metric	Mean Baseline	DML	NCM	NCM vs DML
$\sqrt{\text{PEHE}} \downarrow$	0.121	1.068	<b>0.086</b>	<b>12.4</b> $\times$ better
CATE $R^2 \uparrow$	0.000	-77.0	<b>0.494</b>	<b>+77.5</b> pts
CATE Correlation $\uparrow$	0.000	0.341	<b>0.708</b>	<b>2.1</b> $\times$ better
ATE Rel. Bias % $\downarrow$	+0.2%	-9.2%	<b>-3.8%</b>	<b>2.4</b> $\times$ smaller
ATT Bias $\downarrow$	+0.013	-0.053	<b>-0.004</b>	<b>13.3</b> $\times$ smaller
QTE P10 Bias $\downarrow$	+0.140	-0.276	<b>+0.036</b>	<b>7.7</b> $\times$ smaller
QTE P90 Bias $\downarrow$	-0.135	+0.195	<b>-0.069</b>	<b>2.8</b> $\times$ smaller

Table 7: *Semi-Synthetic Validation: Interaction DGP* ( $\tau(x) = x_1 \cdot x_2$ ). Tests feature interaction effects where treatment effect depends on the product of two covariates. Experimental setup: Mean Baseline predicts constant  $\hat{\tau}$ ; DML uses XGBoost T-Learner with 5-fold cross-fitting on 7.8% treatment rate data (80K samples); NCM is DragonNet with targeted regularization trained on 30% treatment rate balanced data (20.8K samples). All methods evaluated on identical 20K test set with ground-truth CATE. NCM achieves **12.8** $\times$  better  $\sqrt{\text{PEHE}}$  than DML.  $\downarrow$  lower is better,  $\uparrow$  higher is better. Metric definitions in Appendix E.

Metric	Mean Baseline	DML	NCM	NCM vs DML
$\sqrt{\text{PEHE}} \downarrow$	0.123	1.138	<b>0.089</b>	<b>12.8</b> $\times$ better
CATE $R^2 \uparrow$	0.000	-84.6	<b>0.471</b>	<b>+85.1</b> pts
CATE Correlation $\uparrow$	0.000	0.258	<b>0.694</b>	<b>2.7</b> $\times$ better
ATE Rel. Bias % $\downarrow$	+0.4%	+7.5%	<b>+0.9%</b>	<b>8.3</b> $\times$ smaller
ATT Bias $\downarrow$	+0.023	+0.042	<b>+0.006</b>	<b>7.0</b> $\times$ smaller
QTE P10 Bias $\downarrow$	+0.153	-0.284	<b>+0.073</b>	<b>3.9</b> $\times$ smaller
QTE P90 Bias $\downarrow$	-0.146	+0.203	<b>-0.064</b>	<b>3.2</b> $\times$ smaller

## E Evaluation Metrics for Causal Inference

We evaluate causal models using comprehensive metrics organized into three categories: individual-level treatment effect metrics (requiring ground-truth CATE, available only in semi-synthetic settings), population-level treatment effect metrics, and production outcome prediction metrics.

### E.1 Individual-Level Treatment Effect Metrics

These metrics require ground-truth individual treatment effects  $\tau_i = Y_i(1) - Y_i(0)$ , available only in semi-synthetic validation where both potential outcomes are generated.

- $\sqrt{\text{PEHE}}$  (Root Precision in Estimation of Heterogeneous Effects): The root mean squared error of individual treatment effect predictions:

$$\sqrt{\text{PEHE}} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{\tau}(x_i) - \tau(x_i))^2} \quad (2)$$

where  $\hat{\tau}(x_i) = \hat{\mu}_1(x_i) - \hat{\mu}_0(x_i)$  is the predicted CATE and  $\tau(x_i) = Y_i(1) - Y_i(0)$  is the true CATE. Lower is better. This is the primary metric for individual-level accuracy, measuring how well the model captures heterogeneous treatment effects across the population.

- CATE  $R^2$ : The coefficient of determination for treatment effect predictions, measuring explained variance:

$$R_{\text{CATE}}^2 = 1 - \frac{\sum_{i=1}^n (\hat{\tau}_i - \tau_i)^2}{\sum_{i=1}^n (\tau_i - \bar{\tau})^2} \quad (3)$$

where  $\bar{\tau} = \frac{1}{n} \sum_i \tau_i$  is the mean true treatment effect. Values range from  $-\infty$  to 1, where 1 indicates perfect prediction, 0 indicates performance equivalent to predicting the constant mean, and *negative values indicate predictions worse than the constant mean*—a critical failure mode signaling model collapse.

- **CATE Correlation:** Pearson correlation between predicted and true individual treatment effects:

$$\rho_{\text{CATE}} = \frac{\text{Cov}(\hat{\tau}, \tau)}{\sigma_{\hat{\tau}} \cdot \sigma_{\tau}} \quad (4)$$

Measures ranking accuracy for targeting decisions—important even when magnitude estimates are biased, since targeting often requires only correct ordering of treatment effects.

## E.2 Population-Level Treatment Effect Metrics

These metrics assess aggregate treatment effect estimation quality.

- **ATE** (Average Treatment Effect): The population-level expected treatment effect:

$$\text{ATE} = E[Y(1) - Y(0)] = \frac{1}{n} \sum_{i=1}^n \tau_i \quad (5)$$

- **ATE Relative Bias %:** Percentage deviation of estimated ATE from true ATE:

$$\text{ATE Rel. Bias} = \frac{\widehat{\text{ATE}} - \text{ATE}}{\text{ATE}} \times 100\% \quad (6)$$

where  $\widehat{\text{ATE}} = \frac{1}{n} \sum_i \hat{\tau}_i$ . Measures systematic over- or under-estimation of population-level effects.

- **ATT** (Average Treatment effect on the Treated): The expected treatment effect among those who actually received treatment:

$$\text{ATT} = E[Y(1) - Y(0)|T = 1] = \frac{1}{n_1} \sum_{i:T_i=1} \tau_i \quad (7)$$

where  $n_1 = \sum_i T_i$  is the number of treated subjects.

- **ATT Bias:** Absolute deviation of estimated ATT from true ATT:

$$\text{ATT Bias} = \widehat{\text{ATT}} - \text{ATT} \quad (8)$$

Critical for targeting applications—measures effect estimation accuracy on the population actually targeted for intervention.

- **QTE** (Quantile Treatment Effects): Treatment effects at specific quantiles of the effect distribution:

$$\text{QTE}_q = Q_q(\tau) - Q_q(0) \approx Q_q(\{\tau_i\}_{i=1}^n) \quad (9)$$

where  $Q_q$  denotes the  $q$ -th quantile. We report QTE at P10 (10th percentile) and P90 (90th percentile) to assess estimation quality at distribution tails where extreme responders lie.

- **QTE P10/P90 Bias:** Deviation of estimated quantile treatment effects from true values:

$$\text{QTE}_q \text{ Bias} = \widehat{\text{QTE}}_q - \text{QTE}_q \quad (10)$$

Tests stability at distribution tails—models that perform well on average may exhibit severe instability for extreme responders.

## E.3 Propensity Score Metrics

- **Propensity AUC:** Area Under the ROC Curve for treatment prediction:

$$\text{AUC} = P(\pi(X_i) > \pi(X_j) | T_i = 1, T_j = 0) \quad (11)$$

Measures the propensity head's ability to discriminate between treated and control subjects. Higher AUC indicates better propensity estimation, which is essential for the counterfactual correction term *cc* to function correctly.

## E.4 Production Outcome Prediction Metrics

These metrics evaluate outcome prediction quality on observational data where ground-truth treatment effects are unavailable. Accurate outcome predictions are necessary (though not sufficient) for accurate CATE estimates.

- **MAE** (Mean Absolute Error): Average absolute prediction error:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \tag{12}$$

The primary metric for overall prediction accuracy. MAE is robust to outliers (unlike RMSE) and interpretable in outcome units.

- **Median Absolute Error**: The 50th percentile of absolute errors:

$$\text{Median Error} = \text{median}(\{|y_i - \hat{y}_i|\}_{i=1}^n) \tag{13}$$

Captures typical prediction quality for the “average” sample, unaffected by extreme errors. Under zero-inflated distributions, median error reveals whether the model correctly predicts zeros for the majority of samples.

- **Median Bias**: The 50th percentile of signed errors:

$$\text{Median Bias} = \text{median}(\{y_i - \hat{y}_i\}_{i=1}^n) \tag{14}$$

Detects systematic over- or under-prediction. A bias near zero indicates the model learns the correct outcome distribution. Particularly important for treatment effect estimation, where biased outcome predictions propagate to biased CATE estimates.

- **WAPE** (Weighted Absolute Percentage Error): Error normalized by outcome magnitude:

$$\text{WAPE} = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{\sum_{i=1}^n |y_i|} \tag{15}$$

Enables comparison across different outcome scales.  $\text{WAPE} \approx 1.0$  indicates errors are comparable to outcome magnitudes—acceptable for highly sparse distributions where most outcomes are zero.

## F Hyperparameter Sensitivity Analysis

To understand the contribution of each hyperparameter and validate our design choices, we conduct a systematic sweep across 18 configurations. Table 8 presents comprehensive results where each row varies one hyperparameter while holding others at default values ( $n_{\text{clusters}} = 20$ ,  $k_{\text{neighbors}} = 25$ ,  $w_{\text{NCM}} = 0.1$ ,  $q = 0.45$ ).

Several key findings emerged from this analysis:

**Neighbor Count ( $k_{\text{neighbors}}$ ) was the Most Impactful.** Reducing  $k_{\text{neighbors}}$  from 100 to 10 improved cold-start MAE from -17.4% to -19.1%, demonstrating that more localized predictions outperformed broader averages. This suggested cold-start samples benefited from tightly matched neighbors rather than smoothed estimates from larger neighborhoods. The improvement was monotonic: smaller neighborhoods consistently yielded better accuracy.

**Quantile Aggregation was Essential.** Using mean aggregation instead of quantile aggregation *worsened* cold-start MAE by +1.3% (shown in **red**), making the K-Means approach harmful rather than helpful. This critical failure occurred because cold-start outcome distributions are zero-inflated—mean aggregation overweighted rare high-value neighbors, producing systematically inflated predictions. Quantile aggregation (45th percentile) produced conservative estimates aligned with the majority of zero or low-value outcomes. This finding has important implications: naive implementations using mean aggregation would degrade performance.

**Cluster Count was Insensitive.**  $n_{\text{clusters}} \in \{10, 20, 50, 100\}$  all achieved nearly identical performance (-18.2% to -18.3% cold-start MAE improvement). This robustness suggested the K-Means partitioning primarily served to accelerate neighbor search rather than learning meaningful cluster

Table 8: Cold-Start MAE Improvement (%) Across Hyperparameter Configurations. Evaluated on 240K cold-start test samples. Each row varies one hyperparameter while holding others at default values ( $n_{\text{clusters}} = 20$ ,  $k_{\text{neighbors}} = 25$ ,  $w_{\text{NCM}} = 0.1$ ,  $q = 0.45$ ). Negative values indicate improvement (lower MAE is better). Cold MAE  $\Delta\% = (\text{ColdNet MAE} - \text{NCM MAE}) / \text{NCM MAE} \times 100$ . Best configuration ( $k=10$ ) achieves **-19.1%** cold-start MAE improvement. Mean aggregation (**+1.3%**) worsens performance due to overweighting rare high-value outliers in zero-inflated distributions.

Hyperparameter	Value	Cold MAE $\Delta\%$	Overall MAE $\Delta\%$	Cold Median $\Delta\%$
$n_{\text{clusters}}$	10	-18.3%	-3.4%	-83.4%
	20	-18.2%	-3.4%	-83.4%
	50	-18.3%	-3.4%	-83.2%
	100	-18.2%	-3.4%	-83.3%
$k_{\text{neighbors}}$	<b>10</b>	<b>-19.1%</b>	<b>-3.6%</b>	<b>-82.8%</b>
	25	-18.2%	-3.4%	-83.4%
	50	-17.7%	-3.3%	-83.4%
	100	-17.4%	-3.3%	-83.6%
$w_{\text{NCM}}$	0.0	-17.7%	-3.3%	-100.0%
	0.1	-18.2%	-3.4%	-83.4%
	0.2	-18.0%	-3.4%	-73.4%
	0.3	-17.3%	-3.3%	-64.4%
Aggregation	mean	<b>+1.3%</b>	-0.3%	-45.4%
	$q = 0.35$	-17.9%	-3.4%	-84.1%
	$q = 0.40$	-18.2%	-3.4%	-83.9%
	$q = 0.45$	-18.2%	-3.4%	-83.4%
	$q = 0.50$	-17.9%	-3.4%	-82.6%

structure—the FAISS nearest-neighbor retrieval within clusters did the heavy lifting. This insensitivity simplifies deployment: practitioners need not tune cluster count carefully.

**NCM Weight Had Minimal Impact on MAE.**  $w_{\text{NCM}} \in \{0.0, 0.1, 0.2, 0.3\}$  all achieved 17–18% cold-start MAE improvement. However,  $w_{\text{NCM}} = 0.0$  (pure cluster prediction) eliminated any NCM contribution and may sacrifice prediction diversity. We selected  $w_{\text{NCM}} = 0.1$  to retain minimal NCM signal while primarily relying on cluster-based estimates. Note that  $w_{\text{NCM}} = 0.0$  achieved -100% median error reduction because the NCM’s systematic overprediction was completely removed.

**Robust Performance Across Configurations.** All quantile-based configurations achieved 17–19% cold-start MAE improvement, demonstrating that the K-Means approach was robust to hyperparameter choices within reasonable ranges. This stability is critical for production deployment—the method can be deployed without extensive tuning, reducing operational complexity and making it accessible to practitioners without deep expertise in hyperparameter optimization.

**Final Configuration Selection.** Based on this comprehensive analysis, we selected the configuration ( $n_{\text{clusters}} = 20$ ,  $k_{\text{neighbors}} = 10$ ,  $w_{\text{NCM}} = 0.1$ ,  $q = 0.45$ ) as it achieved the best cold-start MAE improvement (-19.1%) while maintaining interpretable settings and computational efficiency. The 20 clusters provided sufficient granularity without excessive computational overhead, and 10 neighbors ensured localized prediction transfer.

## G ColdNet Algorithm

Algorithm 1 presents the complete ColdNet procedure for cluster-based cold-start enhancement.

## H Dataset Details

This appendix provides comprehensive documentation of all datasets used in our experiments, including data sources, preprocessing steps, and feature definitions.

## H.1 Production E-commerce Dataset

### H.1.1 Raw Production Data

The raw observational dataset comprises seller-offer enrollment decisions and subsequent performance outcomes from a major e-commerce platform.

**Key Insight:** Treated samples are  $29\times$  more likely to have positive outcomes (64.71% vs. 2.20%), demonstrating strong selection bias in observational data—precisely the challenge that causal inference methods must address.

### H.1.2 Balanced Training Set (A)

Used for NCM and ColdNet training. Created via outcome-stratified sampling that enriches for treated samples while maintaining cold-start representation.

### H.1.3 Imbalanced Training Set (B)

Used for Baseline DragonNet to demonstrate that scale alone cannot overcome cold-start challenges.

### H.1.4 Balanced Test Set (C)

All methods are evaluated on this identical test set for fair comparison.

## H.2 Cold-Start Feature Definition

We classify a sample as **cold-start** if ANY of the following six historical activity features equals zero. These features span two granularities (offer-level and ASIN-level) and two time windows (90-day and 365-day):

**Feature Semantics.** The six features capture historical sales activity at two levels of aggregation:

- **Offer-level features** (4 features): Track the specific seller-product combination’s performance. An offer with zero offer-level history has never sold this exact product listing, even if the underlying product (ASIN) has sales from other sellers.
- **ASIN-level features** (2 features): Track the product’s aggregate performance across all sellers. An offer with zero ASIN-level history represents a product with no marketplace sales history, regardless of seller.

**Disjunctive Criterion.** The OR-logic classification ensures comprehensive coverage of incomplete profiles. An offer is flagged as cold-start if it lacks historical activity at *either* the offer level *or* the ASIN level, since either gap indicates insufficient historical signal for reliable neural model predictions. This criterion captures:

- 97.9% of raw samples (before stratified sampling)
- 58.6% of the balanced training set (after outcome-stratified sampling)
- 40.7% of the balanced test set

**Rationale for Feature Selection.** These six features were selected because they directly measure historical sales activity—the primary signal that neural causal models use to predict treatment response. Features like product category, price, and seller characteristics remain available for cold-start samples and enable the cluster-based enhancement (Section 2.5), but the absence of *historical activity* features fundamentally limits the model’s ability to learn from past behavior patterns.

## H.3 Semi-Synthetic Datasets

For NCM vs. DML validation with ground-truth CATE.

### H.3.1 Data Generation

**Source:** 9M production samples with 72 numerical features (after removing IDs, categoricals, and pattern-matched columns).

**Feature Processing:** MinMax scaling to [0, 1] range.

**Treatment Effect Generation:** Three data generating processes (DGPs):

- **Sine:**  $\tau(x) = \sin(\text{weighted features})$  — non-linear, periodic heterogeneity
- **Linear:**  $\tau(x) = \beta \cdot x$  — simple linear heterogeneity
- **Interaction:**  $\tau(x) = x_1 \cdot x_2$  — feature interaction effects

**Treatment Effect Scale:**  $\tau = 2.0$  (scaled by outcome standard deviation).

### H.3.2 Dataset Statistics

### H.3.3 Schema

## H.4 Experimental Design Summary

Table 16 summarizes which datasets are used for training and evaluation of each method.

**Key Design Decisions:**

1. **Stratified Sampling:** Raw data has extreme imbalance (0.39% treatment, 99.25% cold-start). Outcome-based stratified sampling creates balanced training data (45.9% treatment) while preserving cold-start distribution (58.6%).
2. **Same Test Set:** All methods (NCM, ColdNet, Baseline DragonNet) are evaluated on the identical 590K balanced test set for fair comparison.
3. **Baseline Uses Imbalanced:** DragonNet baseline is trained on 9.6M imbalanced data to demonstrate that pure scale ( $6\times$  more data) does not overcome the cold-start problem without proper architecture.
4. **Semi-Synthetic Validation:** Ground-truth CATE enables validation of treatment effect estimation quality, not just prediction accuracy—essential for causal inference.

---

**Algorithm 1:** *ColdNet: Cluster-Based Cold-Start Enhancement for Counterfactual Prediction.* The algorithm operates in three phases: (0) Outcome-stratified ensemble preprocessing addresses extreme class imbalance by training  $M$  models on the complete treatment cohort with stratified control subsets; (1) Offline index construction builds a searchable FAISS index of cold-start training samples clustered via K-Means; (2) Online inference routes warm samples directly to NCM outputs while cold-start samples are enhanced via locality-preserving prediction transfer from similar training samples. Hyperparameters:  $M=10$  ensemble models,  $K=20$  clusters,  $k=10$  neighbors,  $w=0.1$  NCM weight,  $q=0.45$  quantile for aggregation.

---

**Hyperparameters:**  $M$  (ensemble size),  $K$  (clusters),  $k$  (neighbors),  $w$  (NCM weight),  $q$  (quantile)

**Input** : Raw data  $\mathcal{D} = \{(x_i, t_i, y_i)\}_{i=1}^N$ , cold-start mask  $\mathcal{C}$

// Phase 0: Outcome-Stratified Ensemble Preprocessing

- 1  $\mathcal{D}_{\text{treat}} \leftarrow \{(x_i, y_i) : t_i = 1\}$ ; // Extract treatment cohort (0.3% of data)
- 2  $\mathcal{D}_{\text{ctrl}} \leftarrow \{(x_i, y_i) : t_i = 0\}$ ; // Extract control cohort
- 3 Stratify  $\mathcal{D}_{\text{ctrl}}$  into 5 outcome bins: zero, p50, p75, p90, max
- 4 Partition  $\mathcal{D}_{\text{ctrl}}$  into  $M$  disjoint subsets  $\{\mathcal{D}_{\text{ctrl}}^{(1)}, \dots, \mathcal{D}_{\text{ctrl}}^{(M)}\}$  preserving strata
- 5 **foreach**  $m \in \{1, \dots, M\}$  **do**
- 6  $\mathcal{D}_{\text{train}}^{(m)} \leftarrow \mathcal{D}_{\text{treat}} \cup \mathcal{D}_{\text{ctrl}}^{(m)}$ ; // All treatment + one control subset
- 7 Train NCM model  $\mathcal{M}_m$  on  $\mathcal{D}_{\text{train}}^{(m)}$  with Equation (1)
- 8 **end**

// Phase 1: Offline Cold-Start Index Construction

- 9  $\mathcal{D}_{\text{cold}} \leftarrow \{(x_i, y_i) : i \in \mathcal{C}_{\text{train}}\}$ ; // Extract cold-start cohort
- 10  $\mathcal{S} \leftarrow \text{StandardScaler.fit}(\{x_i : (x_i, y_i) \in \mathcal{D}_{\text{cold}}\})$ ; // Fit feature scaler
- 11  $X_{\text{scaled}} \leftarrow \mathcal{S.transform}(\{x_i\})$ ; // Standardize features
- 12  $\{c_1, \dots, c_K\} \leftarrow \text{KMeans.fit}(X_{\text{scaled}}, K)$ ; // Cluster centroids
- 13 **foreach** cluster  $j \in \{1, \dots, K\}$  **do**
- 14  $\mathcal{D}_j \leftarrow \{(x_i, y_i) \in \mathcal{D}_{\text{cold}} : \text{assign}(x_i) = j\}$ ; // Samples in cluster  $j$
- 15  $X_j \leftarrow \{x_i : (x_i, y_i) \in \mathcal{D}_j\}$ ; // Feature vectors
- 16  $X_j^{\text{norm}} \leftarrow \text{L2Normalize}(X_j)$ ; // Unit norm for cosine similarity
- 17  $\mathcal{I}_j \leftarrow \text{FAISS.IndexFlatIP}(X_j^{\text{norm}})$ ; // Inner product index
- 18  $Y_j \leftarrow \{y_i : (x_i, y_i) \in \mathcal{D}_j\}$ ; // Store outcomes
- 19 **end**

// Phase 2: Online Counterfactual Inference

**Input** : Test sample  $x^*$ , NCM predictions  $(\hat{\mu}_0^*, \hat{\mu}_1^*)$ , cold-start indicator  $c^*$

**Output** : Enhanced predictions  $(\tilde{\mu}_0^*, \tilde{\mu}_1^*)$

- 20 **if**  $c^* = \text{False}$ ; // Warm sample
- 21 **then**
- 22 **return**  $(\hat{\mu}_0^*, \hat{\mu}_1^*)$ ; // NCM predictions unchanged
- 23 **else**
- // Cold-start sample: apply locality-preserving transfer
- 24  $x_{\text{scaled}}^* \leftarrow \mathcal{S.transform}(x^*)$ ; // Apply fitted scaler
- 25  $j^* \leftarrow \arg \min_j \|x_{\text{scaled}}^* - c_j\|_2$ ; // Assign to nearest centroid
- 26  $x_{\text{norm}}^* \leftarrow \text{L2Normalize}(x_{\text{scaled}}^*)$ ; // Normalize query
- 27  $\mathcal{N}_k \leftarrow \mathcal{I}_{j^*}.\text{search}(x_{\text{norm}}^*, k)$ ; // Retrieve  $k$  nearest neighbors
- 28  $Y_{\mathcal{N}} \leftarrow \{Y_{j^*}[i] : i \in \mathcal{N}_k\}$ ; // Neighbor outcomes
- 29  $\hat{y}_{\text{cluster}} \leftarrow Q_q(Y_{\mathcal{N}})$ ; // Quantile-robust aggregation
- // Blend both potential outcome heads
- 30  $\tilde{\mu}_0^* \leftarrow w \cdot \hat{\mu}_0^* + (1 - w) \cdot \hat{y}_{\text{cluster}}$
- 31  $\tilde{\mu}_1^* \leftarrow w \cdot \hat{\mu}_1^* + (1 - w) \cdot \hat{y}_{\text{cluster}}$
- 32 **return**  $(\tilde{\mu}_0^*, \tilde{\mu}_1^*)$
- 33 **end**

---

Table 9: *Raw Production Data Statistics. Original e-commerce dataset before outcome-stratified sampling. Treatment rate of 0.39% creates 1:256 class imbalance. Cold-start rate of 99.25% indicates nearly all samples lack complete historical features. Treated samples show 29× higher positive outcome rate (64.71% vs. 2.20%), demonstrating strong selection bias.*

Attribute	Train	Test
Total Samples	180,111,727	48,996,865
Features	161	161
Treatment Rate	0.39%	0.39%
Cold-Start Rate	99.25%	99.25%
Outcome > Threshold	2.44%	2.44%
Outcome Rate (Treated)	64.71%	–
Outcome Rate (Control)	2.20%	–

Table 10: *Balanced Training Set (A) Statistics. Created via outcome-stratified sampling from raw data. Used for NCM and ColdNet training. Treatment rate increased from 0.39% to 45.87% to enable stable gradient updates. Cold-start representation preserved at 58.62%.*

Attribute	Value
Total Samples	1,530,903
Features	161
Treatment Rate	45.87%
Treated Samples	702,286
Control Samples	828,617
Cold-Start Samples	897,443 (58.62%)
Warm Samples	633,460 (41.38%)

Table 11: *Imbalanced Training Set (B) Statistics. Used for Baseline DragonNet to demonstrate that scale alone (6× more data than balanced set) cannot overcome cold-start challenges without proper architecture. Treatment rate of 8.9% is higher than raw data but still imbalanced.*

Attribute	Value
Total Samples	9,587,881
Features	161
Treatment Rate	8.9%
Cold-Start Samples	5,178,941 (54.0%)
Warm Samples	4,408,940 (46.0%)

Table 12: *Balanced Test Set (C) Statistics. All methods (NCM, ColdNet, Baseline DragonNet) are evaluated on this identical test set for fair comparison. Contains 240K cold-start samples (40.67%) and 350K warm samples (59.33%).*

Attribute	Value
Total Samples	590,255
Features	161
Treatment Rate	24.9%
Treated Samples	147,005
Control Samples	443,250
Cold-Start Samples	240,064 (40.67%)
Warm Samples	350,191 (59.33%)

Table 13: *Cold-Start Definition Features.* We classify a sample as cold-start if ANY of these six historical activity features equals zero (disjunctive criterion). Features span two granularities (offer-level: specific seller-product combination; ASIN-level: product aggregate across all sellers) and two time windows (90-day and 365-day). This criterion captures 97.9% of raw samples and 58.6% of the balanced training set.

Feature Name	Granularity	Window	Description
offer_net_ordered_gms_amt_sum_90d	Offer	90-day	Gross merchandise sales (\$)
offer_net_ordered_gms_amt_sum_365d	Offer	365-day	Gross merchandise sales (\$)
offer_net_shipped_units_sum_90d	Offer	90-day	Number of units shipped
offer_net_shipped_units_sum_365d	Offer	365-day	Number of units shipped
offer_asin_net_ord_gms_sum_90d	ASIN	90-day	ASIN-aggregated GMS (\$)
offer_asin_net_ship_units_sum_90d	ASIN	90-day	ASIN-aggregated units shipped

Table 14: *Semi-Synthetic Dataset Statistics.* Generated from 9M production samples with 72 numerical features. Treatment effects synthesized via three DGPs (Sine, Linear, Interaction) with ground-truth CATE available for direct evaluation. Imbalanced train mimics production conditions (7.8% treatment); balanced train enables stable NCM training (30% treatment).

Dataset	Samples	Treatment %	Purpose
Imbalanced Train	80,000	7.8%	DML/Mean Baseline
Balanced Train	20,839	30%	NCM
Test (All)	20,000	30%	All methods

Table 15: *Semi-Synthetic Data Schema.* Each sample contains observed outcome (factual), ground-truth CATE (for evaluation), true propensity score, both potential outcomes (counterfactual), and 72 MinMax-scaled features. Ground-truth CATE =  $y_1 - y_0$  enables direct evaluation of treatment effect estimation quality.

Column	Description
treatment	Binary treatment indicator (0/1)
y_observed	Observed outcome = $T \times y_1 + (1 - T) \times y_0$
true_tau	<b>Ground truth CATE</b> = $y_1 - y_0$
propensity_score	True $P(T = 1 X)$
y0_potential	Counterfactual outcome under control
y1_potential	Counterfactual outcome under treatment
feature_0...feature_71	72 features (0-1 scaled)

Table 16: *Experimental Design: Training and Evaluation Data by Method.* All production methods are evaluated on the identical Balanced Test (C) for fair comparison. NCM and ColdNet use balanced training data (45.9% treatment) while Baseline DragonNet uses 6x more imbalanced data (8.9% treatment) to demonstrate that scale alone cannot overcome cold-start challenges. Semi-synthetic validation uses ground-truth CATE to evaluate treatment effect estimation quality directly.

Method	Training Data	Samples	Treatment %	Test Data
<i>Production Evaluation (Table 3)</i>				
NCM	Balanced Train (A)	1.53M	45.9%	Balanced Test (C)
ColdNet	Balanced Train (A)	1.53M	45.9%	Balanced Test (C)
Baseline DragonNet	Imbalanced Train (B)	9.59M	8.9%	Balanced Test (C)
<i>Semi-Synthetic Validation (Table 1)</i>				
Mean Baseline	Imbalanced Semi-Synth	80K	7.8%	Semi-Synth Test
DML (XGBoost)	Imbalanced Semi-Synth	80K	7.8%	Semi-Synth Test
NCM (DragonNet)	Balanced Semi-Synth	20.8K	30%	Semi-Synth Test