

Predictive Modeling of Customer Lifetime Value in E-Commerce Using Deep Learning and Causal Inference

Author: Olatunji Olusola Ogundipe, **Affiliation:** Kanpee, **Email:** olatunji.ogundipe@kanpee.com

Abstract

Customer Lifetime Value (CLV) is a central metric in e-commerce for acquisition budgeting, personalization, retention, and strategic planning. Classical statistical and probabilistic models (e.g., Pareto/NBD, BG/NBD) provide principled baselines but struggle with high-dimensional covariates, nonstationary behavior, and counterfactual questions required for causal decisioning. Deep learning has recently delivered substantial gains in CLV point and distributional prediction at industrial scale by integrating representation learning, sequence models, and distributional heads; however, purely predictive models risk conflating correlation with causation when used to inform interventions (promotion allocation, pricing, retention offers). In this article we propose a rigorous, production-ready framework that fuses state-of-the-art deep learning architectures for CLV (sequential encoders, attention/Transformers, mixture/distributional output layers) with modern causal inference techniques (potential outcomes, double/debiased machine learning, causal forests, representation learning for counterfactuals) to deliver accurate, robust, and actionable CLV estimates for e-commerce. We provide formal problem statements, architecture blueprints, objective functions, validation protocols (temporal cross-validation, backtest, uplift evaluation), and risk/ethical governance guidance. We demonstrate how probabilistic deep CLV models (e.g., zero-inflated / mixture output, heteroskedastic heads) can be combined with causal estimators (double ML, causal forests, learned balanced representations) to produce both predictive scores and valid estimates of *causal* effects of marketing actions on CLV enabling prescriptive decisioning with sound uncertainty quantification. We ground our design choices in the literature and present a reproducible experimental

protocol and evaluation suite for industry benchmarking.

Keywords: Customer Lifetime Value (CLV), deep learning, causal inference, double machine learning, causal forests, sequence models, survival analysis, distributional prediction, e-commerce

1. Introduction

Estimating the value a customer will generate over their lifetime CLV is one of the most consequential predictive tasks in e-commerce. Accurate CLV models inform acquisition spend, personalized promotions, segmentation, churn mitigation, inventory and supply planning, and long-term strategy (e.g., which cohorts to prioritize). Classical CLV methods (behavioral probabilistic models and rule-based approaches) provide interpretability and tractable uncertainty quantification for simple transactional data, but modern e-commerce generates high-dimensional, multimodal data (clickstream, product views, text, returns, session behavior, ads exposure) and acts on customers (offers, price changes) which both necessitate flexible predictive models and demand a causal perspective for decisioning. [SSRN+1](#)

This paper develops a unified framework that (i) leverages deep learning to model complex, nonstationary customer behavior for CLV prediction, (ii) incorporates causal inference methods to identify the *effect* of interventions (e.g., coupons, retargeting) on CLV, and (iii) provides practical guidance for deployment and governance in e-commerce settings. We emphasize reproducible experimental design, robust evaluation (including uplift and offline policy evaluation), and auditability all essential for production systems that will be used to execute financial decisions.

In Section 2 we survey relevant literature spanning classical CLV, deep learning methods, and causal inference for treatment effect estimation. Section 3 formalizes the CLV estimation and decision problem. Section 4 describes modeling architectures and loss functions for predictive and distributional CLV. Section 5 integrates causal estimators (double/debiased machine learning, causal forests, representation learning for counterfactuals) into the modeling pipeline. Section 6 details evaluation methodology and robustness checks. Section 7 discusses practical productionization, interpretability, and governance. Section 8 presents limitations and research directions; Section 9 concludes.

2. Related Work

2.1 Probabilistic and statistical CLV models

Early probabilistic approaches model repeat purchase behavior and dropout explicitly (Pareto/NBD; BG/NBD) and estimate expected future transactions from recency/frequency/monetary (RFM) statistics; these methods remain important baselines because of interpretability and closed-form properties in many settings. Fader, Hardie, and Lee (2005) present the BG/NBD model as an accessible alternative to the Pareto/NBD framework with robust empirical performance in many retail contexts.

2.2 Machine learning and deep learning for CLV

From the 2010s onward, machine learning (gradient-boosted trees, ensembles) became competitive with parametric CLV models by effectively handling many engineered features. In the late 2010s and 2020s deep learning approaches (sequence encoders, embeddings, mixture output heads) have shown superior performance for complex, high-cardinality e-commerce data; examples include deep probabilistic CLV models that model zero-inflation and heavy tails with mixture losses, and industrial systems that scale to billions of users (e.g., Kuaishou's industrial solution, perCLTV for games). These works demonstrate that neural architectures with careful distributional heads

and domain-specific inductive biases can deliver both accuracy and deployability at scale.

2.3 Causal inference and machine learning

The potential outcomes (Rubin) and structural causal model (Pearl) formalisms provide the foundations for causal claims and counterfactual reasoning; modern machine learning tools have been incorporated into causal estimation via double/debiased machine learning and tree-based heterogeneous treatment effect estimators (e.g., causal forests). Chernozhukov et al. (2017/2018) propose Double/Debiased Machine Learning (DML) to obtain \sqrt{n} -consistent estimates of treatment effects in the presence of high-dimensional nuisance components estimated by ML; Wager & Athey (2018) develop causal forests to estimate heterogeneous treatment effects with valid inference. Representation learning approaches (e.g., Johansson, Shalit & Sontag, 2016) adapt neural nets to produce balanced representations for counterfactual inference in high-dimensional observational data. [arXiv+2arXiv+2](#)

2.4 Hybrid predictive-causal CLV

Recent literature increasingly recognizes that point-predictive models are insufficient when used to support intervention decisions one must separate predictive accuracy from causal identifiability. Hybrid pipelines that combine expressive prediction models with causal estimators (e.g., using powerful feature learners inside double ML or causal forests) produce both accurate CLV forecasts and valid treatment effect estimates, enabling uplift optimization and prescriptive policies. This manuscript synthesizes these threads, focusing on architecture, objectives, and evaluation tailored to e-commerce.

3. Problem Statement and Formal Setup

We consider a merchant with customers indexed by $i=1, \dots, n$. For each customer we observe a history up to time t_0 :

- Transactional sequence $X_i = \{(t_{i,k}, p_{i,k}, q_{i,k}, c_{i,k})\}_{k=1}^{K_i}$

$p_{i,k}, q_{i,k}, c_{i,k}\}_{k=1}^{K_i}$ where t is timestamp, p product, q quantity, c price/revenue.

- Session/clickstream sequences S_i (page views, dwell, categories).
- Static covariates Z_i (demographics, acquisition channel).
- Past marketing exposures $A_i(0:t_0)$ (treatment history: coupons, ads, emails).
- Observed outcomes: cumulative realized revenue up to t_0 , and any labeled churn indicators.

Objective (Predictive CLV): Estimate $\widehat{CLV}_i(\tau | H_i, t_0)$, the expected discounted future revenue from t_0 to $t_0 + \tau$ conditional on observed history H_i, t_0 .

Objective (Causal CLV / Uplift): Estimate the causal effect on CLV of a candidate intervention d (e.g., sending coupon d at time t_0):

$$\Delta_i(d) \equiv E[CLV_i(\tau) | do(A_i, t_0 = d), H_i, t_0 - A] - E[CLV_i(\tau) | do(A_i, t_0 = d_0), H_i, t_0 - A]$$

$$\equiv E[CLV_i(\tau) | do(A_i, t_0 = d), H_i, t_0 - A] - E[CLV_i(\tau) | do(A_i, t_0 = d_0), H_i, t_0 - A]$$

where H_{-A} denotes history excluding the action at t_0 , and d_0 is baseline (no action). Valid estimation requires ignorability assumptions (or instruments / experiment design) or robust semiparametric methods (e.g., DML) when ignorability is plausible conditional on observables.

We target two outputs:

1. A probabilistic predictive distribution $\hat{F}_i(y)$ for future revenue (point estimate + calibrated uncertainty).
2. An estimator of individualized treatment effects (ITE) $\widehat{\Delta}_i(d)$ (with confidence intervals) to support decisioning (who to treat, how much to offer).

4. Modeling: Deep Predictive Architectures for CLV

This section details candidate architecture and loss formulations for CLV prediction in e-commerce.

4.1 Input representation and embedding

- **Entity embeddings:** high-cardinality categorical features (product IDs, category, campaign id) encoded as learned embeddings.
- **Temporal encodings:** event timestamps converted to time-since-last event, cyclical time features (day-of-week), and positional encodings for sequence models.
- **Multimodal fusion:** structured features, text (product reviews, chat), and image embeddings combined via late or cross attention.

4.2 Sequential encoders

- **RNN/GRU/LSTM:** Good for per-customer purchase sequences when event order matters and sequences are moderate in length.
- **Transformer encoders:** Scalable to long sequences and support self-attention that captures modality interactions; application to CLV prediction has recently been adopted in industrial recipes. [Keras](#)
- **Hybrid:** short-term LSTM + long-context Transformer to capture multi-scale dependencies.

4.3 Output heads: distributional and mixture models

E-commerce CLV distributions are heavy-tailed and zero-inflated (many users never purchase again within horizon). We strongly recommend distributional output layers rather than point regression:

- **Zero-inflated mixture (ZILN):** Mixture of a point mass at zero and a log-normal (or Gamma) for positive revenue; train by maximum likelihood (negative log-likelihood of mixture). Proven effective in prior work. [arXiv](#)
- **Mixture density networks (MDN):** Model CLV as a mixture of parametric components (e.g., lognormals) with neural mixture weights.
- **Quantile regression heads:** Predict multiple quantiles (e.g., 10th, 50th, 90th) for calibration.
- **Bayesian last-layer / MC Dropout:** Uncertainty quantification via approximate Bayesian inference (Monte Carlo Dropout), useful for risk-sensitive decisioning. [arXiv](#)

4.4 Losses and regularization

- Negative log-likelihood of chosen distributional head (ZILN/MDN).
- Weighted combination with downstream business losses (e.g., expected profit after cost of offers), calibration penalties (CRPS), and fairness/regularization constraints.
- Use label smoothing, adversarial noise, and time-aware dropout to prevent overfitting to historical regimes.

4.5 Two-stage vs. joint multitask modeling

- **Two-stage:** Predict intermediate quantities (churn hazard, purchase frequency, average order value) and combine via probabilistic composition to produce CLV. Classical models often adopt this; it offers interpretability.
- **End-to-end:** Direct neural regression to CLV can exploit end-to-end losses and complex features but requires careful calibration. A hybrid approach multitask heads predicting

churn probability + expected spend can be effective and improves robustness.

4.6 Example architecture (pseudocode)

Input: Customer history H_i

Embeddings: e =
 $\text{EmbeddingLayer}(\text{categorical_features})$

Sequence encoding: z_{seq} =
 $\text{TransformerEncoder}(\text{sequence_features})$

Static encoding: $z_{\text{static}} = \text{MLP}(\text{static_features})$

Fusion: $z = \text{Concatenate}([z_{\text{seq}}, z_{\text{static}}, e])$

Heads:

- ProbAliveHead $\rightarrow p_{\text{alive}}$ (sigmoid)

- SpendHead \rightarrow params of ZILN (p_i , μ , σ)

Loss = $-\log_{\text{likelihood}}(\text{ZILN} \mid \text{ground_truth}) + \lambda_1 * \text{calibration_loss} + \lambda_2 * \text{regularization}$

5. Causal Inference for CLV: Estimands, Identification, and Estimation

Predictive CLV alone often misleads decision-making because it does not answer *what will happen if we do X?* To operationalize CLV for interventions (who to target, personalized offers), we must estimate causal effects.

5.1 Target estimands

- **Average Treatment Effect on CLV (ATE):**
 $\tau = E[\text{CLV}_i(1) - \text{CLV}_i(0)]$
 $\tau = E[\text{CLV}_i(1)] - E[\text{CLV}_i(0)]$
- **Conditional / Heterogeneous Treatment Effect (CATE / ITE):**
 $\tau(x) = E[\text{CLV}_i(1) - \text{CLV}_i(0) \mid X_i = x]$
 $\tau(x) = E[\text{CLV}_i(1) \mid X_i = x] - E[\text{CLV}_i(0) \mid X_i = x]$, where X_i includes high-dimensional covariates.

5.2 Identification assumptions

- **Ignorability (unconfoundedness):** $\{CLV_i(1), CLV_i(0)\} \perp A_i | X_i$. Plausible when treatment assignment is well measured and includes campaign, timing, exposure, and selection features; otherwise require randomized experiments or instruments.
- **Overlap:** All customers have nonzero probability of receiving each treatment conditional on X_i .

If ignorability is suspect, prioritize randomized experiments or utilize instrumental variables or panel difference-in-differences designs.

5.3 Estimation strategies

5.3.1 Double / Debiased Machine Learning (DML)

DML decomposes estimation into flexible nuisance estimation (propensity, outcome models) using ML and orthogonalization to remove first-order bias, producing \sqrt{n} -consistent estimates under mild conditions. To estimate ATE/CATE for CLV (a continuous, heavy-tailed outcome), one can use DML with robustification (e.g., trimmed propensity scores) to obtain valid inference in high-dimensional settings. [arXiv](#)

Practical recipe:

1. Split data into folds.
2. Fit ML models for propensity $e(X)$ and outcome $m(X, A)$ (here, the deep CLV model can serve as m).
3. Form orthogonal scores and solve for treatment effect; average across folds.

5.3.2 Causal forests and generalized random forests

Causal forests estimate heterogeneous treatment effects nonparametrically with valid confidence intervals and are robust to irrelevant covariates; they are suitable when we need individualized effect

estimates and policy targeting. Use causal forests to complement DML when effect heterogeneity is central. [arXiv](#)

5.3.3 Representation learning for counterfactuals (neural balancing)

If the feature space is high-dimensional (embeddings, sequence encodings), train a neural representation that balances treated and control distributions (minimizes discrepancy) while preserving predictive power for outcomes. Methods like TARNet / Dragonnet / the representation learning framework of Johansson et al. can be adapted for CLV; they have been shown to outperform classical kernels in many observational settings.

5.3.4 Uplift / Two-model and single-model approaches

Traditional uplift uses two separate models (predict outcome under treatment and control), but modern best practice uses doubly robust learners (DML) or targeted learning approaches combining propensity and outcome prediction for improved robustness.

5.4 Combining deep predictive models with causal estimators

A practical pipeline:

1. **Representation stage:** Train a deep encoder $\phi\psi(X)$ that compresses history into a low-dimensional, informative state; optionally use a balancing regularizer (Wasserstein / MMD) to reduce distribution shift across treatments.
2. **Nuisance estimation stage:** Fit propensity $e(\phi\psi(X))$ and outcome $m(\phi\psi(X), A)$ models using flexible learners (XGBoost, deep nets).
3. **Orthogonalization / DML stage:** Compute orthogonal scores and estimate ATE/CATE with sample splitting.

4. **Heterogeneity discovery:** Use causal forest on $\phi\psi(X)\backslash\phi\psi(X)\phi\psi(X)$ or tree-structured learners to discover interpretable heterogeneity segments.
5. **Policy learning:** Solve for a personalized treatment policy (cost-sensitive) using regret-minimization / off-policy evaluation with doubly robust estimators.

This integrated approach returns both (a) distributional CLV forecasts and (b) validated causal uplift estimates with uncertainty for decisioning.

6. Evaluation: Metrics, Validation, and Experimental Protocols

Robust evaluation is essential: temporal leakage, distribution shift, and censored outcomes are pervasive in CLV.

6.1 Predictive metrics

- **Point metrics:** MAE, RMSE (for logged positive revenue), MAPE (with caution for zeros), Gini / area under Lorenz curve for ranking quality.
- **Distributional metrics:** CRPS (continuous ranked probability score), negative log-likelihood, calibration plots (reliability diagrams), sharpness.
- **Business metrics:** expected profit uplift under deployment policy, cost-adjusted CLV, and percent gain over baseline rules.

6.2 Causal metrics and validation

- **Policy evaluation:** Off-policy evaluation using Inverse Probability Weighting (IPW) and Doubly Robust (DR) estimators to estimate deployed policy value.
- **Uplift evaluation:** Qini curves and uplift AUC for ranking uplift capacity.
- **Confidence intervals and coverage:** Evaluate empirical coverage of estimated ITE

intervals via held-out randomized subsets or synthetic experiments.

6.3 Cross-validation and backtesting

- **Temporal cross-validation (rolling / expanding window):** Ensure train-test splits respect time ordering to emulate deployment.
- **Bootstrapped backtests:** Assess variability of performance under resampling.
- **Censoring and survival:** If customer lifetimes are censored, integrate survival models or inverse probability censoring weights when computing horizon CLV.

6.4 A/B testing and hybrid evaluation

For causal claims, randomized experiments remain the gold standard; use experiments to validate model uplift estimates and calibrate the causal pipeline. When experiments are unavailable, use sensitivity analyses (Rosenbaum bounds) to quantify robustness to unobserved confounding.

7. Implementation, Productionization, and Governance

7.1 System architecture and infrastructure

- **Feature store** for preprocessing and serving consistent features across training and inference.
- **Online / offline model parity:** maintain identical preprocessing in offline and online pipelines.
- **Serving constraints:** real-time scoring vs batched scoring tradeoffs; sequence encoding may be precomputed to reduce latency.
- **Logging & audit trails:** store inputs, model outputs, and decision logs for post-hoc analysis and regulatory requirements.

7.2 Operational considerations

- **Cold start:** combine cohort-level priors (population distribution) and rapid fine-tuning for new customers.
- **Retraining cadence:** schedule periodic retraining and monitoring for distribution drift; consider continual learning with conservative updates.
- **Experimentation:** allocate holdout treatment groups for randomized validations and shadow deployments.

7.3 Interpretability and human oversight

- Use SHAP / integrated gradients for feature attribution on predicted CLV; complement with causal forest segment explanations for uplift.
- Provide decision makers with *both* predictive CLV and causal uplift (with CIs), and simple, actionable decision rules (e.g., treat if $\Delta^i > \delta$ and $\widehat{\Delta}_i > \delta$ and $\widehat{\text{CLV}}_i > \delta$ under budget constraints).

7.4 Data quality, privacy, and security

- Ensure robust fraud detection and data integrity (anomalous transactions, bot activity) as these distort CLV estimates; secure data exchange and governance are critical for cross-system feature sharing. User-provided works highlight the importance of secure, cloud-based analytics and AI for fraud detection and secure pipelines in large systems. (See Fatunmbi; Samuel). (*User-provided references included in final bibliography.*)

8. Example Experimental Protocol (Reproducible)

1. **Datasets:** public or proprietary e-commerce transaction logs with customer IDs, timestamps, products, session features, marketing exposures. Partition data temporally: train on [T0, T1], validate on [T1+1,

T2], test on [T2+1, T3] with multiple rolling folds.

2. **Baselines:** BG/NBD + Gamma-Gamma monetary model (classical), XGBoost two-stage model, deep probabilistic model (ZILN), perCLTV / ODMN (industrial baselines).
3. **Predictive models:** Transformer encoder + ZILN head (primary), LSTM baseline, two-stage architecture with churn + spend heads.
4. **Causal estimators:** DML using deep outcome & propensity, causal forest on learned representation, and representation-balanced Dragonnet variant for counterfactuals.
5. **Metrics:** MAE, CRPS, calibration, uplift Qini, expected policy profit (DR estimator), interval coverage.
6. **Ablations:** (a) distributional head vs MSE, (b) representation balancing vs none, (c) DML vs two-model uplift.
7. **Reporting:** include fairness audits, calibration plots, policy confusion matrix (who is targeted vs who should be).

9. Limitations, Risks, and Ethical Considerations

- **Confounding and selection bias:** Observational causal estimates rely on unconfoundedness; if violated, estimates may mislead. Use experiments or IV/DID designs when possible.
- **Distribution shift:** Behavior and economics change; retrain and conduct drift detection.
- **Privacy:** CLV relies on personal data; comply with regulations (GDPR, CCPA) and minimize data. Use federated analytics where appropriate.
- **Fairness / discrimination:** Personalized offers may inadvertently discriminate; audit policies and include fairness constraints in optimization.

- **Fraud and poisoning:** Fraudulent transactions or adversarial manipulation can bias models; integrate fraud detection and data quality controls. User-provided studies emphasize deploying robust fraud detection and secure cloud analytics to mitigate such risks (Samuel). (See References.)

10. Conclusion and Roadmap

We presented an end-to-end blueprint for combining deep learning and causal inference to deliver both accurate predictive CLV models and valid causal effect estimates for e-commerce decisioning. The central idea is to use expressive deep encoders and distributional heads for prediction, and to place these learned components within rigorous causal estimation frameworks (DML, causal forests, balanced representations) so that treatment policies (e.g., who receives coupons) are guided by estimated uplift with quantifiable uncertainty. This hybrid approach enables prescriptive decisions that are statistically defensible and operationally scalable.

Key research and deployment priorities include: (1) improving representation balancing for disentangling treatment selection from outcome mechanisms; (2) scaling distributional deep models to billion-user regimes (industrial solutions show promising patterns); (3) standardizing uplift policy evaluation frameworks in production; and (4) building robust governance to mitigate fairness, privacy, and fraud risks.

References

1. Fader, P. S., Hardie, B. G. S., & Lee, K. L. (2005). "Counting Your Customers" the Easy Way: An Alternative to the Pareto/NBD Model. *Marketing Science*, 24(2), 275–284.
2. Wang, X., Liu, T., & Miao, J. (2019). A Deep Probabilistic Model for Customer Lifetime Value Prediction. arXiv:1912.07753.
3. Li, K., Shao, G., Yang, N., Fang, X., & Song, Y. (2022). Billion-user Customer Lifetime Value Prediction: An Industrial-scale Solution from Kuaishou. Proceedings of CIKM 2022.
4. Fatunmbi, T. O. (2022). Deep learning, artificial intelligence, and machine learning in healthcare: Applications and future directions. *World Journal of Advanced Research and Reviews*, 15(2), 1–12. <https://doi.org/10.30574/wjarr.2022.15.2.0359>
5. Zhao, S., Wu, R., Tao, J., Qu, M., Zhao, M., Fan, C., & Zhao, H. (2023). perCLTV: A General System for Personalized Customer Lifetime Value Prediction in Online Games. *ACM Transactions on Information Systems*, 41(1).
6. Cao, X., Xu, Y., & Yang, X. (2024). Customer Lifetime Value Prediction with Uncertainty Estimation Using Monte Carlo Dropout. arXiv:2411.15944.
7. Fatunmbi, T. O. (2024). Advanced frameworks for fraud detection leveraging quantum machine learning and data science in fintech ecosystems. *World Journal of Advanced Engineering Technology and Sciences*, 12(01), 495–513. <https://doi.org/10.30574/wjaets.2024.12.1.0057>
8. Wager, S., & Athey, S. (2018). Estimation and Inference of Heterogeneous Treatment Effects using Random Forests. (Causal forests) arXiv / JASA related work.
9. Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., & Robins, J. (2017/2018). Double/Debiased Machine Learning for Treatment and Structural Parameters. arXiv / NBER.
10. Johansson, F. D., Shalit, U., & Sontag, D. (2016). Learning Representations for Counterfactual Inference. Proceedings of ICML 2016.
11. OptDist: Weng, Y., Tang, X., Xu, Z., et al. (2024). OptDist: Learning Optimal Distribution for Customer Lifetime Value Prediction. arXiv:2408.08585.

12. Additional survey and application works on CLV and deep learning (selected): Sun, Y. et al. (2023) review on CLV with ML; perCLTV and ODMN industrial publications; various arXiv and conference papers cited inline.
13. Pearl, J. (2009). *Causality: Models, Reasoning and Inference* (2nd ed.). Cambridge University Press.
14. Rubin, D. B. (1974). *Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies*. *Journal of Educational Psychology*, 66(5), 688–701.
15. Samuel, A. J. (2022). *AI and machine learning for secure data exchange in decentralized energy markets on the cloud*. *World Journal of Advanced Research and Reviews*, 16(2), 1269–1287. <https://doi.org/10.30574/wjarr.2022.16.2.1282>
16. Samuel, A. J. (2023). *Enhancing financial fraud detection with AI and cloud-based big data analytics: Security implications*. *World Journal of Advanced Engineering Technology and Sciences*, 9(02), 417–434. <https://doi.org/10.30574/wjaets.2023.9.2.0208>