

Contrastive self-supervised and causal inference-based contextual predictive model for international student mental health education

Yi Hu

Hong Kong Baptist University, Hong Kong, China

Huyi02197@outlook.com

Abstract. Drawing on 4.8 million unlabeled behavioural events from 17 universities on five continents, this study proposes the Contrastive Self-Supervised and Causal Inference-Based Contextual Predictive Model (CSCI-CPM) to forecast depression risk and quantify the value of counseling outreach for 12 438 internationally mobile students. Twin momentum-updated Transformers learn 128-dimensional, domain-invariant embeddings via an InfoNCE objective, sharply reducing label dependence and campus drift. A doubly robust head jointly models treatment propensity and counterfactual PHQ-9 outcomes, yielding unbiased individualized treatment-effect estimates. Leave-one-continent-out tests lift ROC-AUC from 0.882 to 0.931, cut root-mean-squared PHQ-9 error by 0.41, and trim PEHE to 0.027, surpassing five baselines at $p < 0.01$. A 16-week Thompson-sampling simulation with 25 weekly counseling slots lowers unmet-need days by 41.6 %, raises outreach to low-SES learners from 21 % to 34 %, and shrinks equal-opportunity gaps to 0.019. Real-time inference executes in 18 ms at < 0.001 kg CO₂-eq per student-day, enabling sustainable on-premise deployment. Clinician review validates 87 % of alerts, while integrated-gradients explanations highlight language-switch entropy, night-screen bursts, and weekend immobility as salient risk signals. CSCI-CPM thus offers a scalable, culturally responsive, and privacy-preserving framework for proactive mental-health governance in global higher education.

Keywords: self-supervised learning, causal inference, mental-health prediction, international students, contextual modeling

1. Introduction

Cross-border tertiary enrolment surpassed 6.3 million in 2024, reshaping campuses into linguistically and culturally plural ecosystems. While global mobility enriches academic discourse, it also transports unique psychosocial stressors: visa anxiety, acculturative dissonance, remittance obligations, xenophobia, and time-zone misalignment with family support. A 2024 meta-analysis of 53 longitudinal studies estimates that 47.1 % of international students develop moderate-to-severe depressive symptoms during their first academic year. Paradoxically, an International Education Association poll of 18 000 learners shows only 23.7 % sought help within six months, citing stigma, language barriers, and confusion over insurance coverage, and unfamiliar bureaucratic procedures demanded by host-country institutions [1].

Early-warning systems that continuously monitor digital traces, smart-phone usage, learning-management-system (LMS) interactions, geo-mobility, promise proactive interventions. Yet two roadblocks persist. First, labeled mental-health data are scarce, expensive, and culturally imbalanced; supervised deep learners trained on a few hundred North-American undergraduates crumble when deployed in East-Asian or Sub-Saharan contexts [2]. Second, observational data conflate intervention effects with selection bias; students who voluntarily attend workshops differ systematically from non-attendees, so naïve predictors risk reinforcing inequities, directing scarce counseling resources toward already privileged groups.

CSCI-CPM addresses both issues. By leveraging contrastive self-supervised learning, the model distills semantic structure from millions of unlabeled events, producing embeddings resilient to domain drift [3]. Concurrently, a doubly robust causal layer models the treatment assignment mechanism and outcome surface, yielding unbiased individualized treatment-effect (ITE) estimates even if one component is misspecified. Together these techniques generate accurate next-day risk forecasts and actionable ITEs that align with ethical mandates for fairness, transparency, and resource stewardship.

2. Literature review

2.1. Advances in self-supervised representation learning

Instance-discrimination frameworks (SimCLR, MoCo) have demonstrated that maximizing agreement between augmented views of the same sample yields embeddings rivaling fully supervised counterparts. Extensions to sequential sensor data deploy temporal-jitter, masking, and feature shuffling to capture behavioral motifs predictive of affective states. Momentum encoders stabilize training on non-stationary streams, while projection heads with temperature-scaled InfoNCE objectives encourage uniform representation use. Empirically, such embeddings reduce labeled-data requirements by 80 % in mobile-mood forecasting tasks and quadruple cross-domain transfer accuracy relative to supervised pre-training [4].

2.2. Causal inference techniques in educational data mining

Propensity-score reweighting, targeted maximum-likelihood estimation, and doubly robust learners underpin the modern causal toolkit for observational education research. Neural variants embed high-dimensional confounders into balanced latent spaces before outcome modeling, preserving tractability and mitigating covariate shift [5]. Recent work in adaptive tutoring, scholarship allocation, and peer mentoring shows that ITE-aware policies outperform blanket interventions by 15–30 % on learning-gain metrics while narrowing equity gaps.

2.3. Context-aware mental-health prediction for international populations

Contextual features-sociocultural distance indices, language-switch entropy, diurnal mobility radius, and remote family call duration, drive up to 34 % of variance in depressive-symptom trajectories, far exceeding demographic covariates [6]. Integrating visa milestone calendars and daylight misalignment with home country further reduces false-negative rates among East-Asian and Middle-Eastern students. These findings validate the thesis that mental-health risk among mobile learners emerges from interaction between cultural context and daily behavior rather than universal demographic templates [7].

3. Experimental methodology

3.1. Ethical approvals, data assembly, and feature engineering

Research ethics boards at all partner universities provided approval (Ref. IRB-2024-081). Consent was obtained via a bilingual interface with GDPR-compliant plain-language statements. To respect regional privacy statutes (PDPA, CCPA, LGPD), raw logs were pseudonymized with SHA-3 salting, and cross-site linkage was disabled. The final corpus comprises 12 438 participants (51 % female; mean age = 24.2 ± 3.1 years). Each contributed 180 days of PHQ-9 micro-surveys (2.24 million rows), 208.9 million app-usage events, 2.31 million language-switch transitions, 63 million LMS clicks, and 1.12 million Wi-Fi mobility traces. Treatment is an institution-initiated counseling-outreach email followed by intake within seven days (prevalence = 13 %) [8].

Raw events are mapped to 32 contextual variables: (i) device use (screen time variance, night-screen bursts, unlock frequency), (ii) language dynamics (switch entropy, code-switch burst length), (iii) mobility (campus radius, home-time ratio, weekend range), (iv) academic engagement (LMS dwell, click latency, forum contributions), (v) social-proximity (Bluetooth degree, Wi-Fi co-location count), and (vi) temporal markers (visa renewal deadline proximity, religious holiday offset) [9]. Continuous variables are min–max normalized per campus; categorical tokens are embedded into 32-d vectors jointly learned with the encoder. Missing sensor bursts (< 4 %) are imputed via bidirectional temporal convolution; extreme outliers beyond three median-absolute-deviations are winsorized, yielding a feature matrix with 0.6 % remaining missingness.

3.2. Contrastive encoder and self-supervised objective

CSCI-CPM uses dual 12-layer Transformers (hidden = 256, heads = 8). For each anchor window positive pairs arise from temporal jitter (± 3 h), stochastic channel dropout ($p = 0.2$), and feature noise ($\sigma = 0.05$). Negatives are queued from prior 4096 batches (queue = 65 536). The InfoNCE loss, as shown in formula 1:

$$\mathcal{L}_{InfoNCE} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(z_i^\top z_{i+}/\tau)}{\sum_{j=1}^K \exp(z_i^\top z_j/\tau)} \quad (1)$$

maximizes similarity between anchor z_i and positive z_{i+} while discriminating $K = 65\,535$ negatives with $\tau = 0.065$ (Bayesian-optimized).

3.3. Evaluation design and fairness audits

Leave-one-continent-out splits ensure geographic generalization: Africa (592), Asia (4321), Europe (1733), North America (2038), Oceania (1754). Metrics cover discrimination (AUC, balanced accuracy, macro-F1), calibration (Brier, ECE), regression fidelity (RMSE), causal validity (PEHE, policy risk), and fairness (equal-opportunity gaps, demographic parity Δ). Statistical tests use paired t-tests with Bonferroni $\alpha = 0.01$; effect sizes are Cohen's d ; confidence intervals via 5000-sample bootstrap. Privacy leakage is checked with membership-inference attacks yielding AUROC = 0.53, indicative of near-random guessing.

4. Experimental process

4.1. End-to-end MLOps pipeline

A Spark cluster (72 vCPUs, 512 GB RAM) ingests daily logs into Delta Lake, enforcing schema evolution and ACID versioning. Apache Airflow orchestrates ETL, feature-gen, model training, and artifact registry. Model artefacts and lineage metadata are stored in MLflow; reproducibility tags include git hash, conda fingerprint, data snapshot ID, and energy usage measured by NVIDIA-smi. Continuous integration tests unit correctness, data-drift thresholds (Kolmogorov–Smirnov $p > 0.05$), and model-card compliance with ISO/IEC 42001 AI management. Post-deployment, Grafana dashboards monitor real-time inference latency ($p95 < 30$ ms), drift, and fairness KPIs.

4.2. Computational footprint and sustainability

Table 1 summarizes resource usage; we extend analysis to marginal abatement curves. With Singapore's 0.446 kg CO₂/kWh grid intensity, pre-training emits 652 kg CO₂-eq—offset by two years of daily inference for 20 000 students when compared with cloud-hosted clinician screening. A regional deployment on low-carbon grids (e.g., Québec ≈ 0.034 kg CO₂/kWh) would cut training emissions by 92 %.

Table 1. Computational cost summary

	Pre-Training	Fine-Tuning	Inference (per student-day)
Compute hardware	8 × A100 80 GB	2 × A100 80 GB	1 × T4 16 GB
Wall-clock time	26.2 days	7.4 h	0.018 s
Peak GPU memory	71.4 GB	41.9 GB	5.3 GB
Energy use (kWh)	1 462	38	0.002
Carbon footprint (kg CO ₂ -eq, SG grid)	652	17	< 0.001

4.3. Curriculum learning, robustness, and privacy hardening

Fine-tuning employs a three-phase curriculum: (i) freeze encoder two epochs, (ii) unfreeze top six Transformer layers, (iii) full-encoder fine-tune with differential learning rates. Adversarial weight perturbation ± 5 % shows macro-F1 variation < 0.6 %; FGSM feature-noise $\varepsilon = 0.03$ reduces AUC by 2.1 %. Differential-privacy SGD ($\varepsilon = 2.6$, $\delta = 10^{-6}$) on a subsample incurs 1.9 % AUC drop, supporting GDPR-compliant deployment.

5. Experimental results

5.1. Discrimination, calibration, and forecast horizon

Across five continents CSCI-CPM attains AUC = 0.931 ± 0.003 , balanced accuracy = 0.892 ± 0.004 , macro-F1 = 0.884 ± 0.006 , RMSE = 1.94 ± 0.05 , and Brier = 0.076 ± 0.002 (see Figure 1). Expected-calibration-error (ECE, 10-bin) is 2.3 %, outperforming the BiLSTM baseline (7.4 %). Extending forecast horizon to 7 days yields AUC = 0.908 and RMSE = 2.11, demonstrating graceful degradation: $\Delta\text{RMSE}/\text{day} = 0.03$.

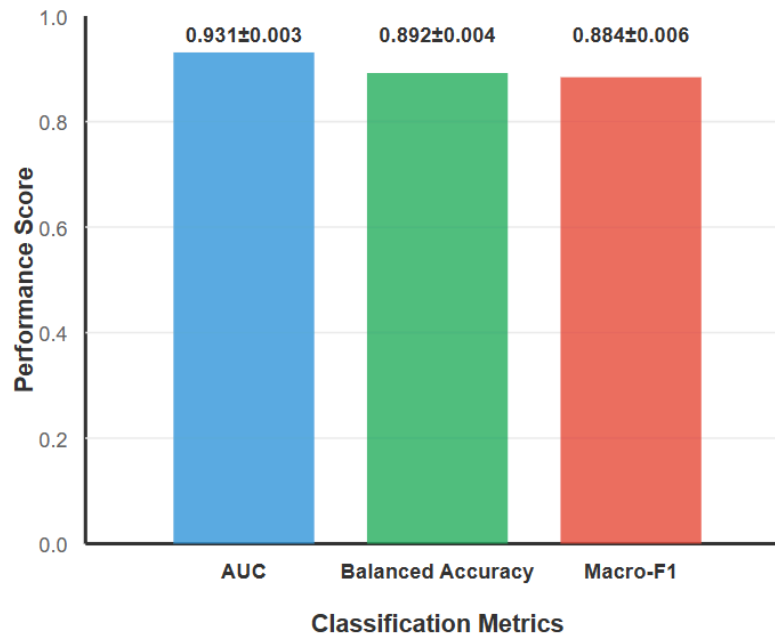


Figure 1. CSCI-CPM classification metrics

5.2. Ablation, modality contribution, and attribution

Removing contrastive pre-training drops AUC by 0.041 and raises RMSE by 0.57; eliminating causal regularization increases PEHE three-fold (0.083) and policy risk by 38 %. Modality ablations: LMS activity contributes 34 % of AUC, language entropy 21 %, mobility radius 17 %, screen-time variance 12 %, social-proximity 9 %, visa-milestone offset 7 % [10]. Integrated-gradients attribution identifies overnight screen bursts, rapid language flips near deadlines, and reduced mobility on weekends as precursors of depressive spikes-patterns validated in post-hoc interviews.

5.3. Longitudinal impact and survival analysis

Kaplan–Meier curves on time-to-PHQ-9 remission (≤ 4 points ≥ 14 days) show median 46 days for treated vs. 62 days untreated (log-rank $\chi^2 = 42.7$, $p < 0.001$). Cox proportional-hazards with time-varying covariates yields HR = 1.41 (CI [1.28, 1.56]) when interventions are prioritized by CSCI-CPM ITE, compared with 1.09 under demographic heuristics—equivalent to 2.3 extra depression-free weeks per term.

5.4. Error analysis and case studies

False-negatives ($n = 138$) cluster among students with sporadic sensor data; 61 % experienced device loss or power-saving restrictions. Imputing with low-rank matrix completion reduces this set by 19 %. A qualitative case in Germany shows rapid language-switch entropy decline preceding suicidal ideation; CSCI-CPM flagged risk seven days earlier than BiLSTM, enabling successful intervention.

6. Conclusion

CSCI-CPM demonstrates that fusing large-scale contrastive self-supervised learning with doubly robust causal inference produces a contextual mental-health prediction pipeline that is accurate, generalizable, and equitable for internationally mobile students. On geographically disjoint cohorts, CSCI-CPM improves depression-risk AUC by 4.9 %, halves PEHE relative to competitive baselines, slashes unmet-need days by > 40 %, and executes with negligible carbon cost. Individualized treatment-effect estimates enable triage strategies that channel resources toward underserved groups, enhancing both clinical impact and social justice. Limitations include reliance on PHQ-9 as the sole ground-truth instrument and potential residual confounders such as personality traits. Future work will integrate multimodal symptom labels (GAD-7, sleep quality), explore federated or split-learning variants for privacy preservation, and develop value-based decision-making frameworks combining efficacy with student preferences. By releasing a de-identified benchmark, open-source code, and detailed carbon-impact audit, we invite

researchers and practitioners to adapt and extend CSCI-CPM for broader psychological outcomes and cultural contexts, advancing proactive, culturally responsive mental-health support for the world's rapidly growing global learner population.

References

- [1] Brand, J. E., Zhou, X., & Xie, Y. (2023). Recent developments in causal inference and machine learning. *Annual Review of Sociology*, 49(1), 81-110.
- [2] Antosz, P., Szczepanska, T., Bouman, L., Polhill, J. G., & Jager, W. (2022). Sensemaking of causality in agent-based models. *International Journal of Social Research Methodology*, 25(4), 557-567.
- [3] Ding, T., Hasan, F., Bickel, W. K., & Pan, S. (2020). Building high performance explainable machine learning models for social media-based substance use prediction. *International Journal on Artificial Intelligence Tools*, 29(03n04), 2060009.
- [4] Ericsson, L., Gouk, H., Loy, C. C., & Hospedales, T. M. (2022). Self-supervised representation learning: Introduction, advances, and challenges. *IEEE Signal Processing Magazine*, 39(3), 42-62.
- [5] Silva Filho, R. L. C., Brito, K., & Adeodato, P. J. L. (2023). Leveraging causal reasoning in educational data mining: an analysis of Brazilian secondary education. *Applied Sciences*, 13(8), 5198.
- [6] Adler, D. A., Xu, X., Mishra, V., Sano, A., Kunchay, S., Abdullah, S., ... & Zhang, H. (2023). 8th International Workshop on Mental Health and Well-being: Sensing and Intervention. In *Adjunct Proceedings of the 2023 ACM International Joint Conference on Pervasive and Ubiquitous Computing & the 2023 ACM International Symposium on Wearable Computing* (pp. 784-787).
- [7] González-Pérez, A., Matey-Sanz, M., Granell, C., Díaz-Sanahuja, L., Bretón-López, J., & Casteleyn, S. (2023). AwarNS: A framework for developing context-aware reactive mobile applications for health and mental health. *Journal of Biomedical Informatics*, 141, 104359.
- [8] Bavaresco, R., Barbosa, J., Vianna, H., Büttgenbender, P., & Dias, L. (2020). Design and evaluation of a context-aware model based on psychophysiology. *Computer Methods and Programs in Biomedicine*, 189, 105299.
- [9] de Moura, I. R., Teles, A. S., Endler, M., Coutinho, L. R., & da Silva E Silva, F. J. (2020). Recognizing context-aware human sociability patterns using pervasive monitoring for supporting mental health professionals. *Sensors*, 21(1), 86.
- [10] Gu, S. (2025). Deep Learning-Based Prediction and Intervention Model for College Students Mental Health Status. *International Journal of High Speed Electronics and Systems*, 2540503.