# A Contrastive Estimation with the Picard iteration

Letting $\beta = |\mathcal{A}^+| - |\mathcal{A}^-| \geq 0$ and writing $U_A$ as the $M \times |A|$ indicator matrix such that $L_A = U_A^\top L U_A$, we have

$$\phi(L) \propto \underbrace{-\beta \log \det(I + X) + \sum_{A \in \mathcal{A}^+} \log \det(U_A^\top X^{-1} U_A)}_{f \text{ convex}}$$
$$\underbrace{+ \beta \log \det(X) - \sum_{A \in \mathcal{A}^-} \log \det(U_A^\top X^{-1} U_A)}_{g \text{ concave}}$$

where the convexity/concavity results follow immediately from [Mariet and Sra, 2015, Lemma 2.3]. Then, the update rule $\nabla f(L_{k+1}) = -\nabla g(L_k)$ requires

$$\beta L_{k+1} + \sum_{A \in \mathcal{A}^-} L_{k+1} U_A (U_A^\top L_{k+1} U_A)^{-1} U_A^\top L_{k+1}$$
$$\leftarrow \beta (I + L_k^{-1})^{-1} + \sum_{A \in \mathcal{A}^+} L_k U_A (U_A^\top L_k U_A)^{-1} U_A^\top L_k$$

which cannot be evaluated due to the $\sum_{A \in \mathcal{A}^-}$ term.

# B Amazon Baby registries experiments

## B.1 Amazon Baby Registries description

Table 3: Description of the Amazon Baby registries dataset.

| REGISTRY | $M$ | TRAIN SIZE | TEST SIZE |
|---|---|---|---|
| HEALTH | 62 | 5278 | 1320 |
| BATH | 100 | 5510 | 1377 |
| APPAREL | 100 | 6482 | 1620 |
| BEDDING | 100 | 7119 | 1780 |
| DIAPER | 100 | 8403 | 2101 |
| GEAR | 100 | 7089 | 1772 |
| FEEDING | 100 | 10,090 | 2522 |

## B.2 Experimental results

In Tab. 4(a), we compare the performance of the various algorithms with rank $K = 30$. The regularization strength $\alpha$ is set to its optimal value for the LR algorithm, and $|\mathcal{A}^-|/|\mathcal{A}^+| = 1/2$. This allows us to compare the LR algorithm to its "augmented" negative versions without hyperparameter tuning. As PROD performs much worse than LR, it is not included in further experiments.

We evaluate the precision at $k$ as

$$p@k = \frac{1}{|\mathcal{T}|} \sum_{A \in \mathcal{T}} \frac{1}{|A|} \sum_{i \in A} \mathbb{1}[\text{rank}(i \mid A \backslash \{i\}) \leq k].$$

Table 4: MPR, p@$k$, and AUC values for LR, and baseline improvement over LR for other methods. Positive values indicate the algorithm performs better than LR, and bold values indicate improvement over LR that lies outside the standard deviation. Experiments were run 5 times, with $|\mathcal{A}^+|/|\mathcal{A}^-| = \frac{1}{2}$; $\alpha$ is set to its optimal LR value.

| Metric | LR | Improvement over LR | | |
|---|---|---|---|---|
| | | DYN | EXP | NCE |
| MPR | 70.50 | $\mathbf{0.92 \pm 0.56}$ | $\mathbf{0.68 \pm 0.62}$ | $\mathbf{0.86 \pm 0.55}$ |
| p@1 | 9.96 | $0.67 \pm 0.75$ | $0.58 \pm 0.76$ | $0.20 \pm 1.75$ |
| p@5 | 25.36 | $\mathbf{1.04 \pm 0.82}$ | $\mathbf{0.78 \pm 0.67}$ | $0.67 \pm 1.09$ |
| p@10 | 36.50 | $\mathbf{1.39 \pm 0.85}$ | $\mathbf{1.13 \pm 0.79}$ | $0.97 \pm 1.18$ |
| p@20 | 51.22 | $\mathbf{1.38 \pm 0.97}$ | $\mathbf{1.28 \pm 1.11}$ | $\mathbf{1.35 \pm 1.20}$ |
| AUC | 0.630 | $\mathbf{0.027 \pm 0.017}$ | $\mathbf{0.026 \pm 0.016}$ | $0.009 \pm 0.017$ |

Compared to traditional SGA methods, algorithms that use inferred negatives perform (PROD excepted) better across all metrics and datasets. DYN and EXP provide consistent improvements compared to the other methods, whereas NCE shows a higher variance and slightly worse performance. Improvements observed using DYN and EXP are larger than the loss in performance due to going from full-rank to low-rank kernels reported in [Gartrell et al., 2017].

Finally, we also compared all methods when tuning both the regularization $\alpha$ and the negative to positive ratio $|\mathcal{A}^-|/|\mathcal{A}^+|$, but did not see any significant improvements. As this suggests there is no need to do additional hyperparameter tuning when using CE, we fix $\frac{|\mathcal{A}^-|}{|\mathcal{A}^+|} = \frac{1}{2}$ for all experiments.