

## 1. Introduction

**Motivation:** Probabilistic matrix factorizations are a powerful tool for modelling matrix  $\mathbf{X}$ .

- ▶ They are **robust to overfitting**.
- ▶ They can account for different **data types** (continuous, ordinal, count, etc...).
- ▶ Fast approximate inference is easily implemented using **variational Bayes**.
- ▶ They scale with the number of entries observed in  $\mathbf{X}$ , which is usually low, and not with the size of  $\mathbf{X}$  which can be very large.

**Problem:** Many real-world **binary** matrices are fully observed. Probabilistic approaches are **infeasible** in this case because they are based on **batch variational algorithms** that require processing all the entries in  $\mathbf{X}$  before producing a single parameter update.

**Solution:** A novel stochastic algorithm for variational inference on big binary matrices:

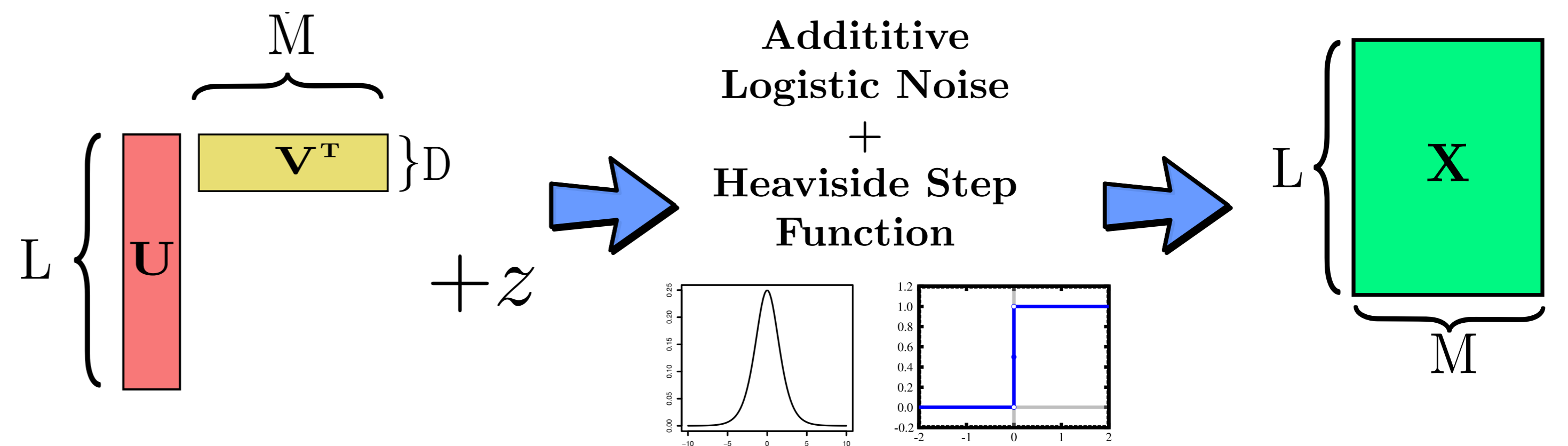
- ▶ We apply the **SVI** method of Hoffman et al., 2013 to matrix factorization models.
- ▶ We subsample **matrix entries** instead of individual data instances.
- ▶ We use **non-uniform** data subsampling strategies which lead to improved predictions.
- ▶ We use minibatches to speed up convergence and **adjust the minibatch size on-line**.

## 2. A Probabilistic Model for Binary Matrices

We use a **logistic likelihood** and a **global bias** parameter.

$$p(\mathbf{X}|\mathbf{U}, \mathbf{V}, z) = \prod_{i=1}^L \prod_{j=1}^M p(x_{i,j}|\mathbf{u}_i, \mathbf{v}_j, z) = \prod_{i=1}^L \prod_{j=1}^M [\sigma(\mathbf{u}_i \mathbf{v}_j^T + z)^{x_{i,j}} \sigma(-\mathbf{u}_i \mathbf{v}_j^T - z)^{1-x_{i,j}}],$$

$$p(\mathbf{U}) = \prod_{i=1}^L \prod_{d=1}^D \mathcal{N}(u_{i,d}|\tilde{u}_{i,d}^0, \tilde{u}_{i,d}^0), \quad p(\mathbf{V}) = \prod_{j=1}^M \prod_{d=1}^D \mathcal{N}(v_{j,d}|\tilde{v}_{j,d}^0, \tilde{v}_{j,d}^0), \quad p(z) = \mathcal{N}(z|\tilde{z}^0, \tilde{z}^0).$$



## 3. Variational Bayes

We approximate the posterior with a tractable  $q(\mathbf{U}, \mathbf{V}, z)$  indexed by variational parameters  $\Phi$ . We optimize  $q$  by maximizing the **Evidence Lower Bound (ELBO)** with respect to  $\Phi$ .

**Jensen's Inequality**

**The ELBO**

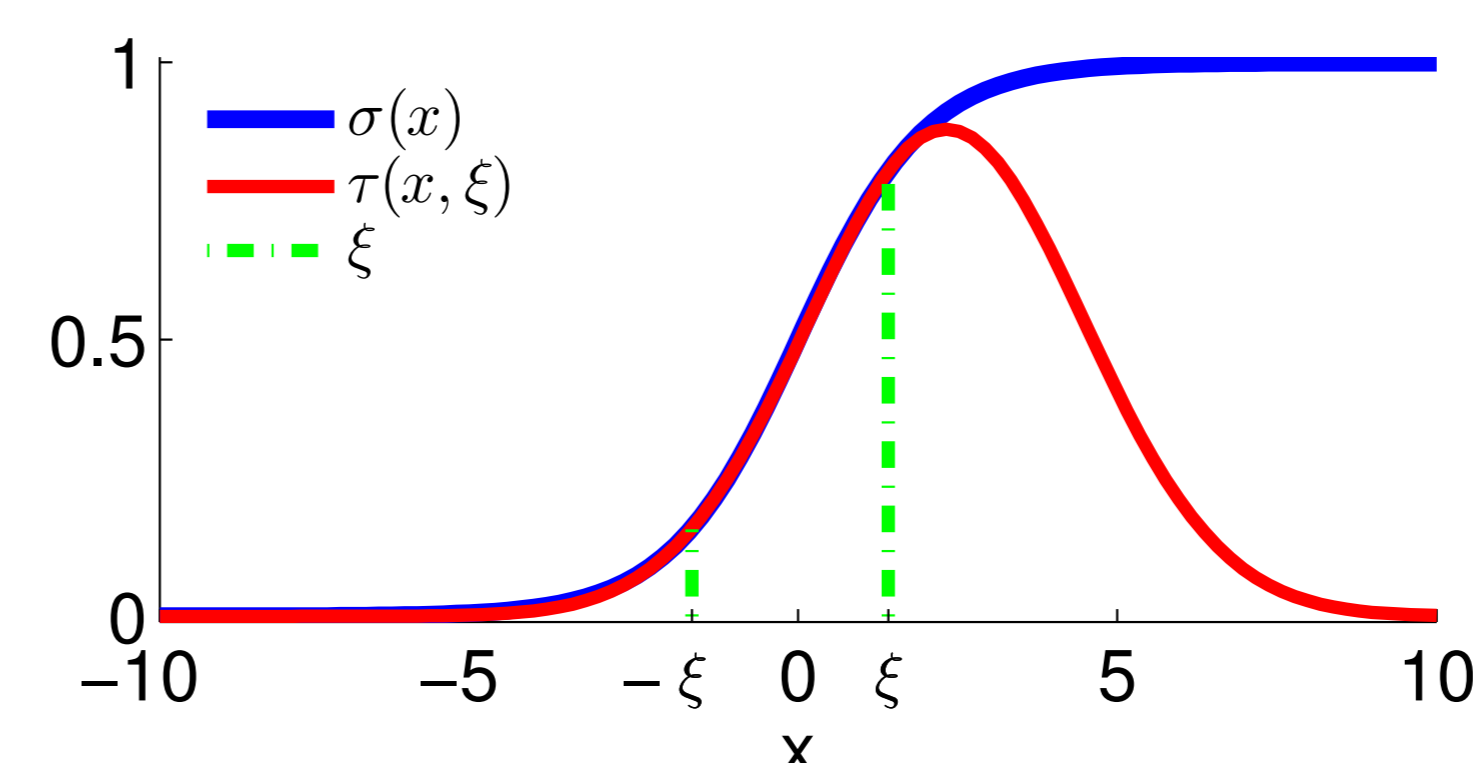
$$\log p(\mathbf{X}) = \log \int p(\mathbf{X}, \mathbf{U}, \mathbf{V}, z) d\mathbf{U} d\mathbf{V} dz \geq \mathbb{E}_{q(\mathbf{U}, \mathbf{V}, z)} \left[ \log \frac{p(\mathbf{X}, \mathbf{U}, \mathbf{V}, z)}{q(\mathbf{U}, \mathbf{V}, z)} \right] \triangleq \mathcal{L}(\Phi).$$

$$q(\mathbf{U}, \mathbf{V}, z) = \left[ \prod_{i=1}^L \prod_{d=1}^D \mathcal{N}(u_{i,d}|\tilde{u}_{i,d}, \tilde{u}_{i,d}) \right] \left[ \prod_{j=1}^M \prod_{d=1}^D \mathcal{N}(v_{j,d}|\tilde{v}_{j,d}, \tilde{v}_{j,d}) \right] \mathcal{N}(z|\tilde{z}, \tilde{z}).$$

$$\Phi = \{ \{ \tilde{u}_{i,d}, \tilde{u}_{i,d} \}_{i=1}^L, \{ \tilde{v}_{j,d}, \tilde{v}_{j,d} \}_{j=1}^M, \tilde{z}, \tilde{z} \}.$$

## 4. Local Variational Approximation

We lower bound each logistic function in the ELBO with a Gaussian:  $\sigma(x) \geq \tau(x, \xi)$ .



(Jaakkola & Jordan, 1997)

Figure: blue = logistic function, red = lower bound, tight at  $x = \pm \xi$ .

We add an **extra variational parameter**  $\xi_{i,j}$  for each matrix entry:  $\Xi = \{ \{ \xi_{i,j} \}_{i=1}^L \}_{j=1}^M$ .

The model is now **conjugate** with Gaussian complete conditionals.

## 5. Stochastic Inference

We use stochastic gradient descent to optimize  $\mathcal{L}(\Phi) \triangleq \arg \max_{\Xi} \mathcal{L}(\Phi, \Xi)$ .

- 1 - Sample a matrix entry  $x_{i,j}$  with probability  $p(i, j)$ .
- 2 - Compute a noisy estimate of  $\mathcal{L}(\Phi)$  which includes only a few of the terms in  $\mathcal{L}(\Phi)$ :

$$\mathcal{L}_{\text{noisy}}(\Phi) = c_{i,j} \underbrace{f(x_{i,j}, \xi_{i,j}, \Phi_{i,j})}_{\text{likelihood}} + \sum_{d=1}^D \underbrace{g(\tilde{u}_{i,d}, \tilde{u}_{i,d})}_{\text{prior on } u_{i,d}} + \sum_{d=1}^D \underbrace{g(\tilde{v}_{j,d}, \tilde{v}_{j,d})}_{\text{prior on } v_{j,d}} + \underbrace{g(\tilde{z}, \tilde{z})}_{\text{prior on } z}.$$

- 3 - Optimize  $\xi_{i,j}$  and choose the values of the scaling constant  $c_{i,j}$ .
- 4 - Update  $\Phi_{i,j} = \{ \{ \tilde{u}_{i,d}, \tilde{u}_{i,d}, \tilde{v}_{j,d}, \tilde{v}_{j,d} \}_{d=1}^D, \{ \tilde{z}, \tilde{z} \} \}$  by taking a small step in the direction of the gradient of  $\mathcal{L}_{\text{noisy}}$ .

## 6. Natural Gradients and Minibatches

We work with **natural parameters**:  $\hat{\mathbf{u}}_{i,d} = [\tilde{u}_{i,d}/\tilde{u}_{i,d}, \tilde{u}_{i,d}^{-1}]^T$ . Let  $\hat{\mathbf{u}}_{i,d}^*$  be the maximizer of  $\mathcal{L}_{\text{noisy}}$  with respect to  $\hat{\mathbf{u}}_{i,d}$ . The **natural gradient** with respect to this parameter is

$$\hat{\nabla} \mathcal{L}_{\text{noisy}}(\hat{\mathbf{u}}_{i,d}) = \hat{\mathbf{u}}_{i,d}^* - \hat{\mathbf{u}}_{i,d}.$$

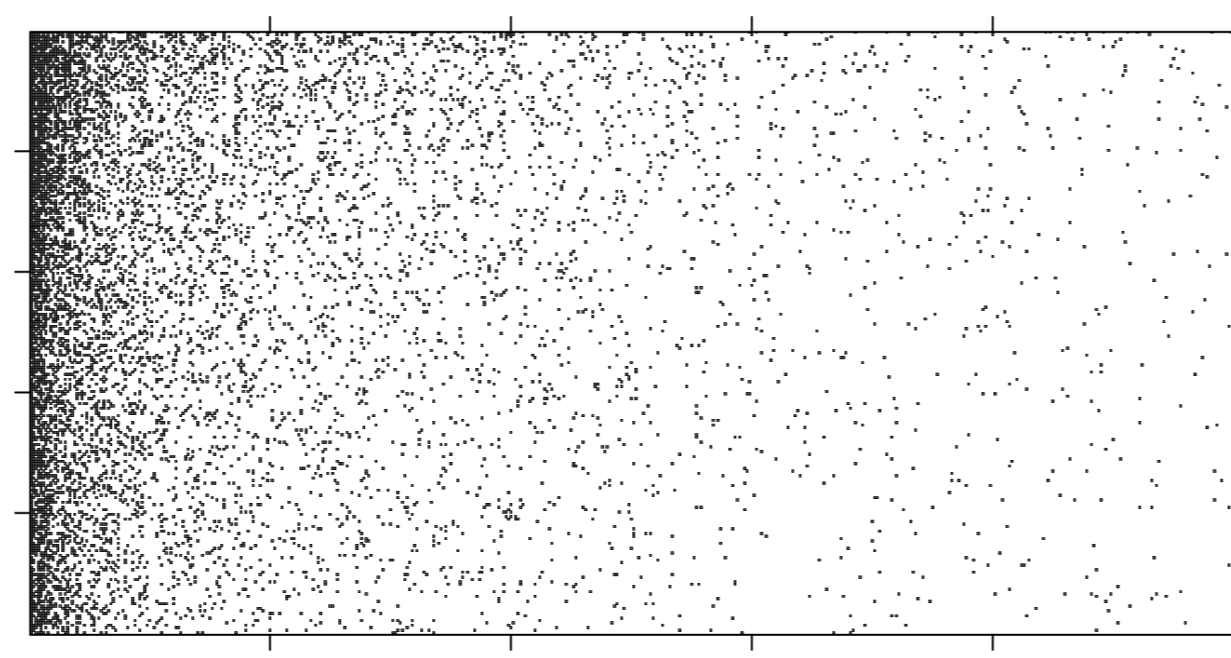
The stochastic update of step size  $\rho$  in the direction of the natural gradient is then

$$\hat{\mathbf{u}}_{i,d}^{\text{new}} = \hat{\mathbf{u}}_{i,d}^{\text{old}} + \rho \hat{\nabla} \mathcal{L}_{\text{noisy}}(\hat{\mathbf{u}}_{i,d}) = (1 - \rho) \hat{\mathbf{u}}_{i,d}^{\text{old}} + \rho \hat{\mathbf{u}}_{i,d}^*.$$

To use **minibatches** of size  $S$ , we replace  $\hat{\mathbf{u}}_{i,d}^*$  with  $\hat{\mathbf{u}}_{i,d}^{*,\text{avg}} = \frac{1}{n(i)} \sum_{s=1}^{n(i)} \hat{\mathbf{u}}_{i,d}^{*,s}$ , where  $n(i)$  is the number of entries from the  $i$ -th row found in the last  $S$  subsampled entries and  $\hat{\mathbf{u}}_{i,d}^{*,s}$  is the maximizer of  $\mathcal{L}_{\text{noisy}}$  when the  $s$ -th of those entries in the  $i$ -th row is subsampled.

## 7. Non-uniform Data Subsampling Strategies

Real-world binary matrices are usually very **sparse**, with frequencies for ones and zeros that **change considerably** across rows and across columns.



We use different subsampling strategies:

- **S-Uniform**:  $p(i, j) = 1/(LM)$ .
- **S-Balanced**:  $p(i, j) = 1/(2 \sum_{a=1}^L \sum_{b=1}^M \mathbb{I}[x_{i,j} = x_{a,b}])$ .
- **S-Biased**:  $p(i, j) = r_i^{(1-x_{i,j})} c_j^{(1-x_{i,j})} [2 \sum_{a=1}^L \sum_{b=1}^M \mathbb{I}[x_{i,j} = x_{a,b}] r_a^{(1-x_{a,b})} c_b^{(1-x_{a,b})}]^{-1}$ .

$r_i^{(0)}$  and  $r_i^{(1)}$  are the number of zeros and ones in the  $i$ -th row of  $\mathbf{X}$  and likewise  $c_j^{(0)}$  and  $c_j^{(1)}$  count the number of zeros and ones in the  $j$ -th column.

## 8. Automatically Adjusting the Minibatch Size Online

The minibatch size  $S$  is important. **Trade off**: noise reduction vs. frequency of updates. We bound the relative error of  $\hat{\mathbf{u}}_{i,d}^{*,\text{avg}}$  with respect to its expectation  $\hat{\mathbf{u}}_{i,d}^* = \mathbf{E}[\hat{\mathbf{u}}_{i,d}^{*,\text{avg}}]$ .

**Markov's Inequality**

$$\delta = p \left( \frac{\|\hat{\mathbf{u}}_{i,d}^{*,\text{avg}} - \hat{\mathbf{u}}_{i,d}^*\|_2^2}{\|\hat{\mathbf{u}}_{i,d}^*\|_2^2} \geq \theta \right) \leq \frac{\mathbb{E} \left[ \|\hat{\mathbf{u}}_{i,d}^{*,\text{avg}} - \hat{\mathbf{u}}_{i,d}^*\|_2^2 \right]}{\theta \|\hat{\mathbf{u}}_{i,d}^*\|_2^2}$$

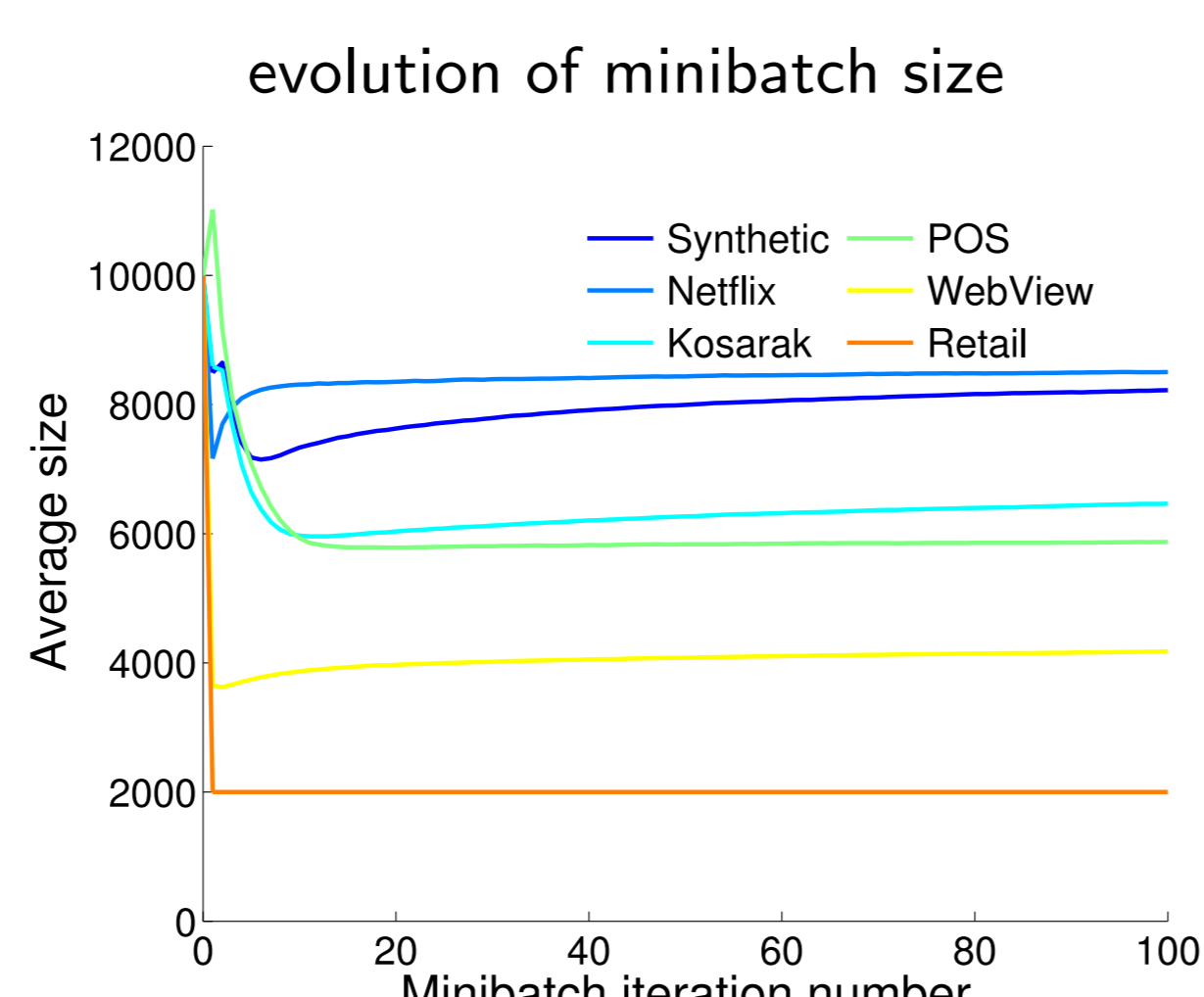
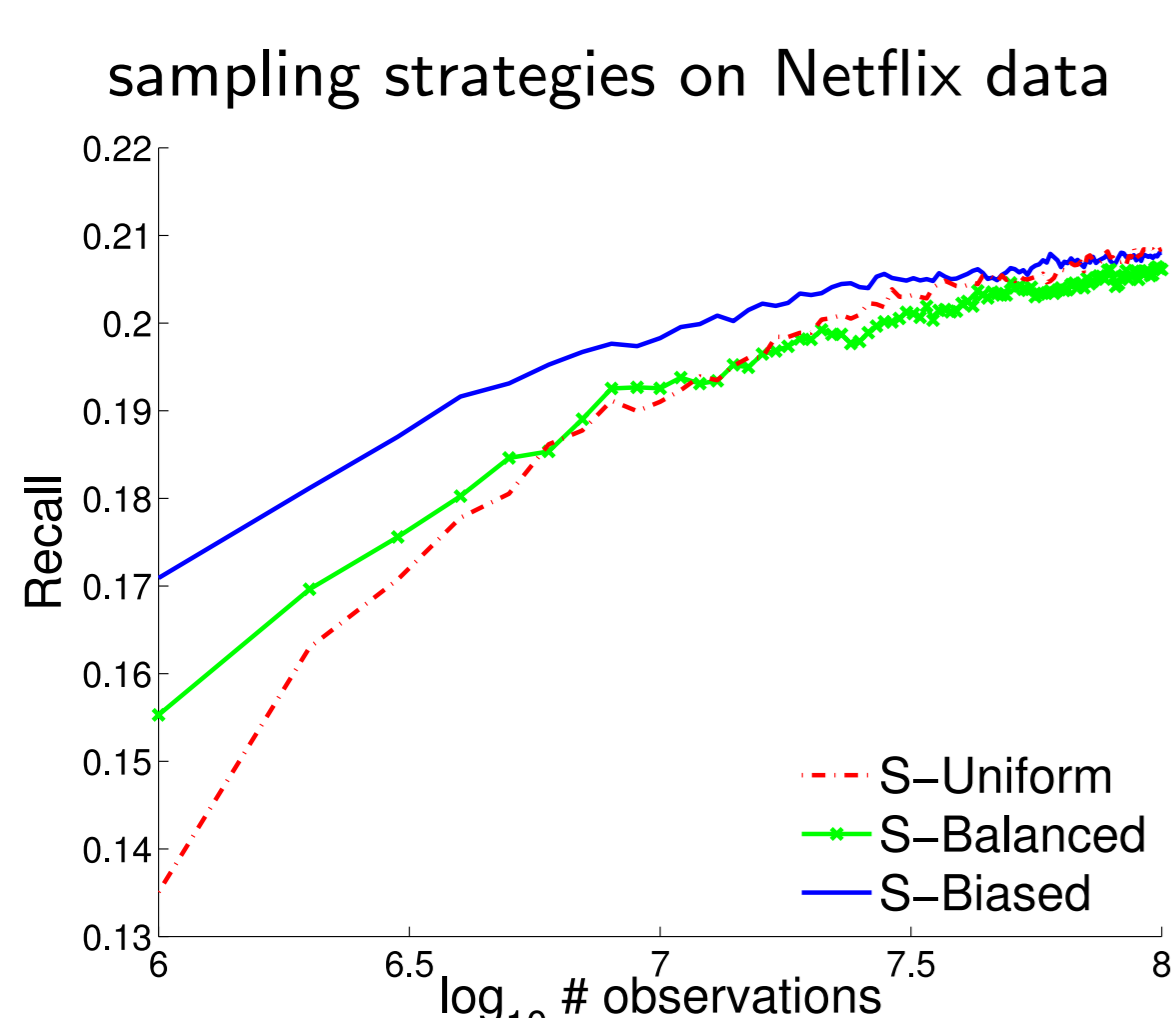
Solving for  $S$ , we obtain that  $S$  should be proportional to the **noise to signal ratio** in  $\hat{\mathbf{u}}_{i,d}^*$ .

$$S = \frac{\|\text{Var}[\hat{\mathbf{u}}_{i,d}^*]\|_1}{\theta \delta p(i) \|\mathbf{E}[\hat{\mathbf{u}}_{i,d}^*]\|_2^2}.$$

- ▶ Only a single effective parameter  $\theta \delta$ .
- ▶ We estimate  $\mathbf{E}[\hat{\mathbf{u}}_{i,d}^*]$  and  $\text{Var}[\hat{\mathbf{u}}_{i,d}^*]$  online.
- ▶ We re-update  $S$  after  $S$  samples have been drawn.

## 9. Sampling Strategies and Evolution of MiniBatchSize

The strategy **S-Biased** performs best. The minibatch size  $S$  converges very quickly.



## 10. Results on Synthetic and Real-world Datasets

