



模式分析与机器智能
工业和信息化部重点实验室
MIT Key Laboratory of
Pattern Analysis & Machine Intelligence



模式识别与神经计算研究组
Pattern Recognition and Neural Computing

Data Valuation Using Reinforcement Learning

Jinsung Yoon*

Department of Electrical and
Computer Engineering, UCLA, CA
jsyoon0823@g.ucla.edu

Sercan Ö. Arık

Google Cloud AI
Sunnyvale, CA
soarik@google.com

Tomas Pfister

Google Cloud AI
Sunnyvale, CA
tpfister@google.com

ICML 2020

□ Background

- REINFORCE algorithm
- Problem setting

□ The method

□ Experiments

- Removing high/low value samples
- Noisy setting
- Transfer setting

□ Conclusions

- The REINFORCE algorithm is a policy gradient method, it **directly update the policy weights**, without learning a value function.

$$E_{\tau \sim p(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

$$\nabla_{\theta} E_{\tau \sim p(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

$$= \int \nabla_{\theta} p(\tau) \left(\sum_{t=1}^T r(s_t, a_t) \right) d\tau$$

$$= \int \underline{p(\tau) \nabla_{\theta} \log p(\tau)} \left(\sum_{t=1}^T r(s_t, a_t) \right) d\tau$$

$$= E_{\tau \sim p(\tau)} \left[\nabla_{\theta} \log p(\tau) \sum_{t=1}^T r(s_t, a_t) \right]$$

$$\nabla_{\theta} \log p(\tau) = \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

- τ is the trajectory
- π is the policy model parameterized by θ

The likelihood ratio trick

- Quantifying the **VALUE** of data is important.
 1. **Incorrect label** (e.g. human labeling errors).
 2. Input comes from a **different distribution** (e.g. different location or time).
 3. Input is noisy or **low quality** (e.g. noisy capturing hardware).
 4. **Usefulness** for target task

- Related works:

1. Leave-One-Out

- Computational **expensive**
- Fail when there are two equal data

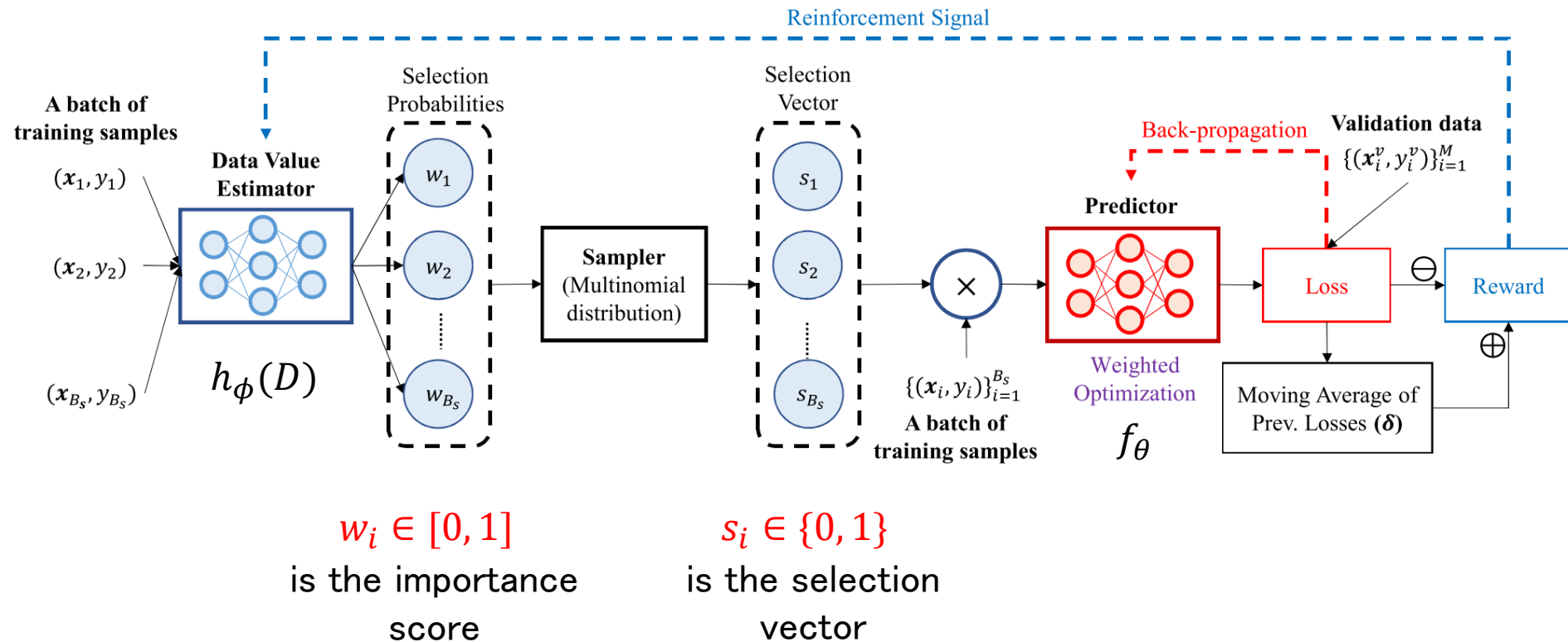
2. Data Shapley

- Computational **expensive**
- Less effective with Monte Carlo approximation

3. Meta-Learning Based Methods

- **Model dependent**
 - **Decoupled from predictor model training**
1. Meta-Weight-Net
 2. Learning to Reweight
 3. MentorNet

- Overall framework



$$s \text{ is sampled from } \pi_\phi(D, s) = \prod_{i=1}^N h_\phi(x_i, y_i)^{s_i} \cdot (1 - h_\phi(x_i, y_i))^{1-s_i}$$

- Overall framework

Algorithm 1 Pseudo-code of DVRL training

- 1: **Inputs:** Learning rates $\alpha, \beta > 0$, mini-batch size $B_p, B_s > 0$, inner iteration count $N_I > 0$, moving average window $T > 0$, training dataset \mathcal{D} , validation dataset $\mathcal{D}^v = \{(\mathbf{x}_k^v, y_k^v)\}_{k=1}^L$
- 2: **Initialize** parameters θ, ϕ , moving average $\delta = 0$
- 3: **while** until convergence **do**
- 4: Sample a mini-batch from the entire training dataset: $\mathcal{D}_B = (\mathbf{x}_j, y_j)_{j=1}^{B_s} \sim \mathcal{D}$
- 5: **for** $j = 1, \dots, B_s$ **do**
- 6: Calculate selection probabilities: $w_j = h_\phi(\mathbf{x}_j, y_j)$ h_ϕ predict scores
- 7: Sample selection vector: $s_j \sim \text{Ber}(w_j)$
- 8: **for** $t = 1, \dots, N_I$ **do**
- 9: Sample a mini-batch $(\tilde{\mathbf{x}}_m, \tilde{y}_m, \tilde{s}_m)_{m=1}^{B_p} \sim (\mathbf{x}_j, y_j, s_j)_{j=1}^{B_s}$
- 10: Update the predictor model network parameters θ

$$\theta \leftarrow \theta - \alpha \frac{1}{B_p} \sum_{m=1}^{B_p} \tilde{s}_m \cdot \nabla_{\theta} \mathcal{L}_f(f_{\theta}(\tilde{\mathbf{x}}_m), \tilde{y}_m)$$

Update prediction model with sampled data

- 11: Update the data value estimator model network parameters ϕ

$$\phi \leftarrow \phi - \beta \left[\frac{1}{L} \sum_{k=1}^L [\mathcal{L}_h(f_{\theta}(\mathbf{x}_k^v), y_k^v)] - \delta \right] \nabla_{\phi} \log \pi_{\phi}(\mathcal{D}_B, (s_1, \dots, s_{B_s}))$$

Update policy model h_{ϕ} with validation set

- 12: Update the moving average baseline (δ): $\delta \leftarrow \frac{T-1}{T} \delta + \frac{1}{LT} \sum_{k=1}^L [\mathcal{L}_h(f_{\theta}(\mathbf{x}_k^v), y_k^v)]$ validation set
-

- Target of the policy model

$$\begin{aligned} \min_{h_\phi} \quad & \mathbb{E}_{(\mathbf{x}^v, y^v) \sim P^t} [\mathcal{L}_h(f_\theta(\mathbf{x}^v), y^v)] \\ \text{s.t.} \quad & f_\theta = \arg \min_{\hat{f} \in \mathcal{F}} \mathbb{E}_{(\mathbf{x}, y) \sim P} [h_\phi(\mathbf{x}, y) \cdot \mathcal{L}_f(\hat{f}(\mathbf{x}), y)] \end{aligned}$$

- Apply the **REINFORCE** method

$$\begin{aligned} \hat{l}(\phi) &= \mathbb{E}_{(\mathbf{x}^v, y^v) \sim P^t} [\mathbb{E}_{\mathbf{s} \sim \pi_\phi(\mathcal{D}, \cdot)} [\mathcal{L}_h(f_\theta(\mathbf{x}^v), y^v)]] \\ &= \int P^t(\mathbf{x}^v) \left[\sum_{\mathbf{s} \in [0,1]^N} \pi_\phi(\mathcal{D}, \mathbf{s}) \cdot [\mathcal{L}_h(f_\theta(\mathbf{x}^v), y^v)] \right] d\mathbf{x}^v, \end{aligned}$$

we directly compute the gradient $\nabla_\phi \hat{l}(\phi)$ as:

$$\begin{aligned} \nabla_\phi \hat{l}(\phi) &= \int P^t(\mathbf{x}^v) \left[\sum_{\mathbf{s} \in [0,1]^N} \nabla_\phi \pi_\phi(\mathcal{D}, \mathbf{s}) \cdot [\mathcal{L}_h(f_\theta(\mathbf{x}^v), y^v)] \right] d\mathbf{x}^v && \text{The likelihood} \\ &= \int P^t(\mathbf{x}^v) \left[\sum_{\mathbf{s} \in [0,1]^N} \nabla_\phi \log(\pi_\phi(\mathcal{D}, \mathbf{s})) \cdot \pi_\phi(\mathcal{D}, \mathbf{s}) \cdot [\mathcal{L}_h(f_\theta(\mathbf{x}^v), y^v)] \right] d\mathbf{x}^v && \text{ratio trick} \\ &= \mathbb{E}_{(\mathbf{x}^v, y^v) \sim P^t} [\mathbb{E}_{\mathbf{s} \sim \pi_\phi(\mathcal{D}, \cdot)} [\mathcal{L}_h(f_\theta(\mathbf{x}^v), y^v)] \nabla_\phi \log(\pi_\phi(\mathcal{D}, \mathbf{s}))], \end{aligned}$$

we directly compute the gradient $\nabla_{\phi} \hat{l}(\phi)$ as:

$$\begin{aligned}\nabla_{\phi} \hat{l}(\phi) &= \int P^t(\mathbf{x}^v) \left[\sum_{\mathbf{s} \in [0,1]^N} \nabla_{\phi} \pi_{\phi}(\mathcal{D}, \mathbf{s}) \cdot [\mathcal{L}_h(f_{\theta}(\mathbf{x}^v), y^v)] \right] d\mathbf{x}^v \\ &= \int P^t(\mathbf{x}^v) \left[\sum_{\mathbf{s} \in [0,1]^N} \nabla_{\phi} \log(\pi_{\phi}(\mathcal{D}, \mathbf{s})) \cdot \pi_{\phi}(\mathcal{D}, \mathbf{s}) \cdot [\mathcal{L}_h(f_{\theta}(\mathbf{x}^v), y^v)] \right] d\mathbf{x}^v \\ &= \mathbb{E}_{(\mathbf{x}^v, y^v) \sim P^t} \left[\mathbb{E}_{\mathbf{s} \sim \pi_{\phi}(\mathcal{D}, \cdot)} [\mathcal{L}_h(f_{\theta}(\mathbf{x}^v), y^v)] \nabla_{\phi} \log(\pi_{\phi}(\mathcal{D}, \mathbf{s})) \right],\end{aligned}$$

where $\nabla_{\phi} \log(\pi_{\phi}(\mathcal{D}, \mathbf{s}))$ is

$$\begin{aligned}\nabla_{\phi} \log(\pi_{\phi}(\mathcal{D}, \mathbf{s})) &= \nabla_{\phi} \sum_{i=1}^N \log \left[h_{\phi}(\mathbf{x}_i, y_i)^{s_i} \cdot (1 - h_{\phi}(\mathbf{x}_i, y_i))^{1-s_i} \right] \\ &= \sum_{i=1}^N s_i \nabla_{\phi} \log [h_{\phi}(\mathbf{x}_i, y_i)] + (1 - s_i) \nabla_{\phi} \log [(1 - h_{\phi}(\mathbf{x}_i, y_i))].\end{aligned}$$

- **Compared methods**

- Random assigned weights
- Leave-One-Out (LOO)
- Data Shapley Value (ICML' 19)
- Learning to Reweight (ICML' 18)
- MentorNet (ICML' 18)
- Influence Function (ICML' 17)

- **Datasets**

- Blog
- Adult
- Rossmann
- HAM 10000
- MNIST
- USPS
- Flower
- Fashion-MNIST
- CIFAR-10
- CIFAR-100
- Email Spam
- SMS Spam

- **Settings**

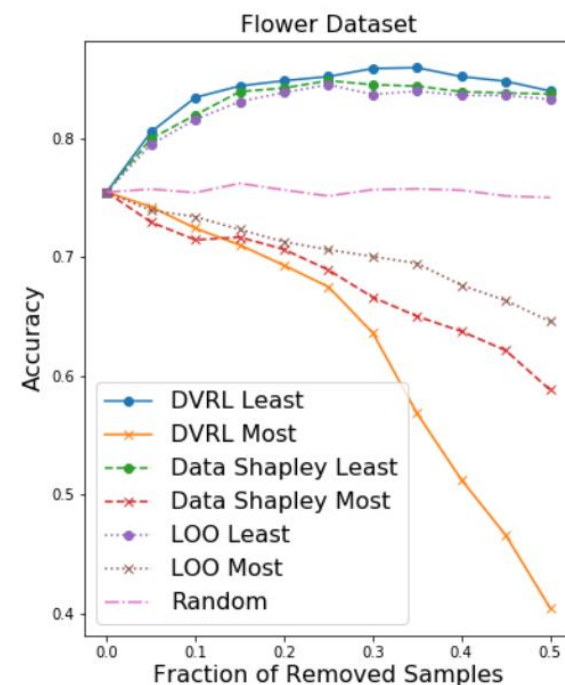
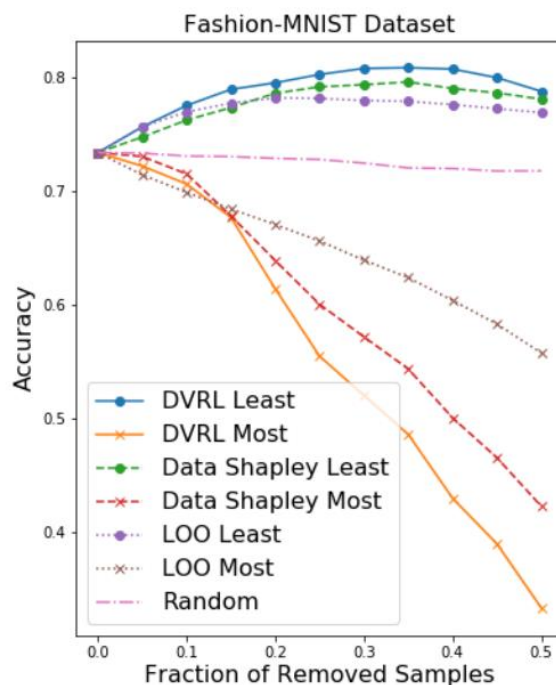
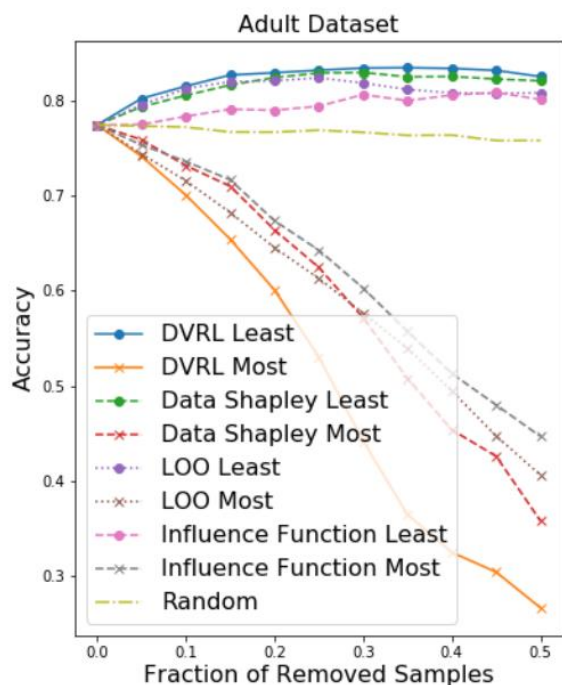
- Removing high/low value samples
- Noisy setting
- Transfer setting

- Clean data



Figure 2: Performance after removing the most (marked with \bigcirc) and least (marked with \times) important samples according to the estimated data values in a conventional supervised learning setting.

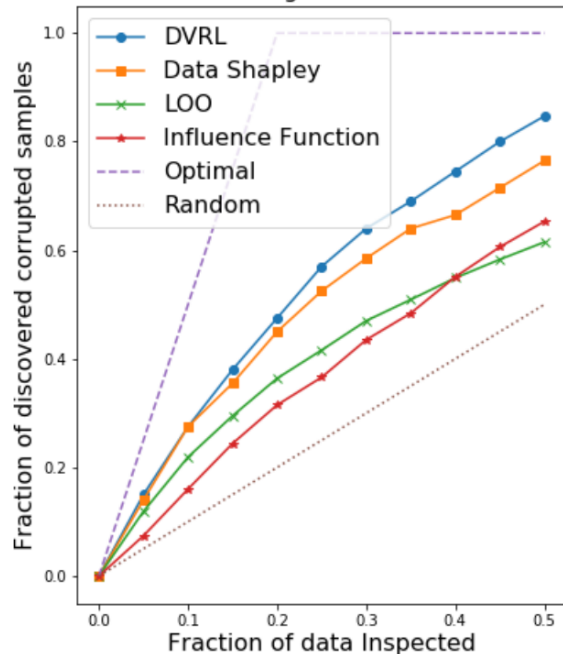
- Noisy data (20%)



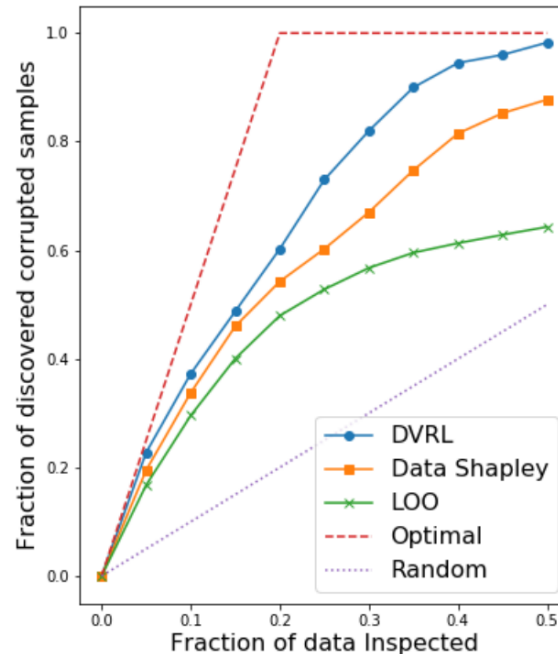
Note that Influence Function is a computationally-efficient approximation of LOO; thus, the performance of Influence Function is similar or slightly worse than LOO in most cases

- Noisy data (20%)

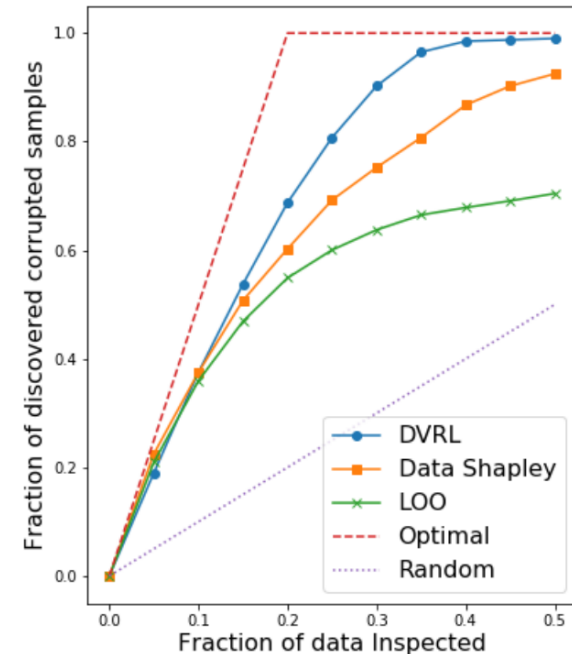
Blog Dataset



HAM-10000 Dataset



CIFAR-10 Dataset



X-axis: The fraction of data with the lowest scores

Y-axis: How many noisy examples are discovered

- Learning from a different domain

Source	Target	Task	Baseline	Data Shapley	DVRL
Google	HAM10000	Skin Lesion Classification	.296	.378	.448
MNIST	USPS	Digit Recognition	.308	.391	.472
Email	SMS	Spam Detection	.684	.864	.903

Predictor Model (Metric: RMSPE)	Store Type	<i>Train on All</i>		<i>Train on Rest</i>		<i>Train on Specific</i>	
		<i>Baseline</i>	DVRL	<i>Baseline</i>	DVRL	<i>Baseline</i>	DVRL
XGBoost	A	0.1736	0.1594	0.2369	0.2109	0.1454	0.1430
	B	0.1996	0.1422	0.7716	0.3607	0.0880	0.0824
	C	0.1839	0.1502	0.2083	0.1551	0.1186	0.1170
	D	0.1504	0.1441	0.1922	0.1535	0.1349	0.1221
Neural Networks	A	0.1531	0.1428	0.3124	0.2014	0.1181	0.1066
	B	0.1529	0.0979	0.8072	0.5461	0.0683	0.0682
	C	0.1620	0.1437	0.2153	0.1804	0.0682	0.0677
	D	0.1459	0.1295	0.2625	0.1624	0.0759	0.0708

Rossmann Store Sales dataset. Metric is Root Mean Squared Percentage Error (RMSPE, lower the better). We use 79% of the data as training, 1% as validation, and 20% as testing

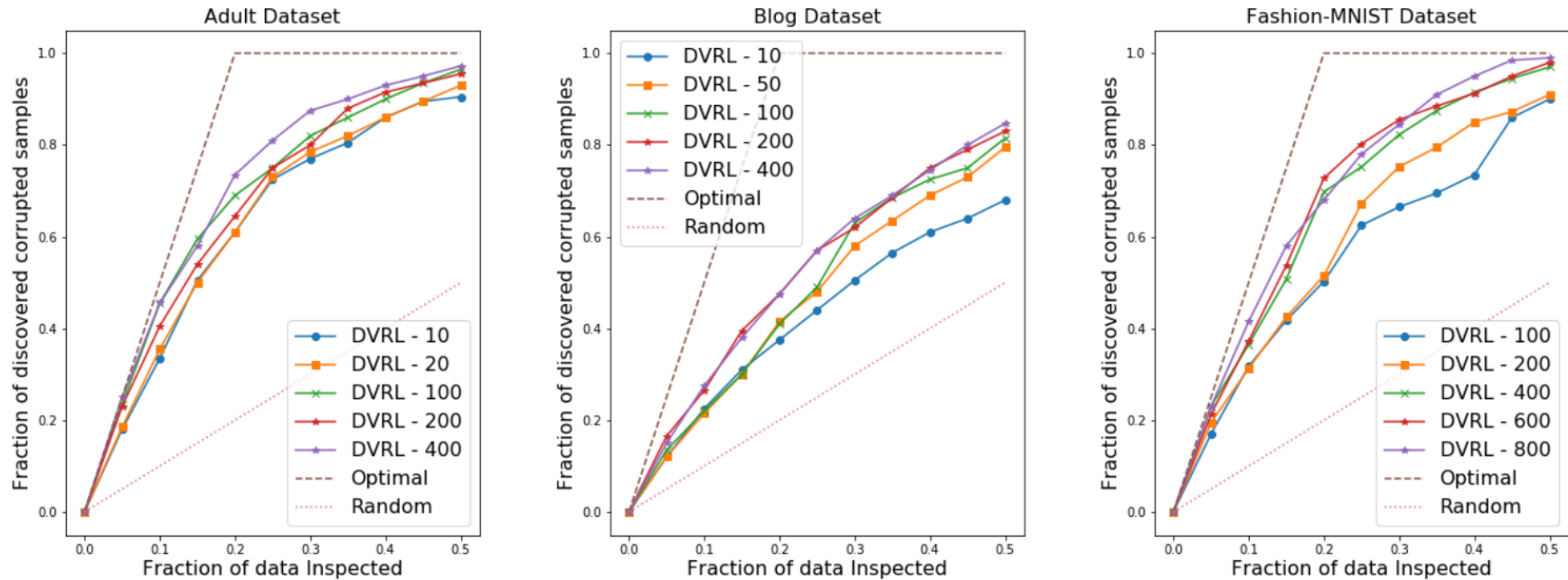


Figure 5: Number of validation samples needed for DVRL. Discovering corrupted samples in three datasets (Adult, Blog and Fashion MNIST) with various number of validation samples.

DVRL achieves reasonable performance with 100 to 400 validation samples. In the Adult dataset, even 10 validation samples are sufficient

- This work propose a meta learning framework, named DVRL, that adaptively **learns data values jointly with a target task predictor model.**
- DVRL uses **REINFORCE** method to update the policy model, and validate the effectiveness.
- DVRL is **model agnostic**, and even applicable to non-differentiable target objectives.