

Attribute Efficient Linear Regression with Distribution-Dependent Sampling

Doron Kukliansky Ohad Shamir

Weizmann Institute of Science

ICML
July 2015

Medical diagnosis classification task

- Problem: volunteers are willing to undergo only a small number of medical tests

Medical diagnosis classification task

- Problem: volunteers are willing to undergo only a small number of medical tests
- Solution: limit the number of tests each volunteer undergoes, with the cost of using more volunteers

What it is?

- Training set of m examples, each of dimension d
- At training time, can view only $k \ll d$ attributes
- At testing time, the model has access to all d attributes

What it is?

- Training set of m examples, each of dimension d
- At training time, can view only $k \ll d$ attributes
- At testing time, the model has access to all d attributes

What it is not?

- Not missing data
- Not feature selection

Problem Definition

$$m = 4, k = 4$$

$$\begin{pmatrix} X_{11} & X_{12} & X_{13} & X_{14} \\ X_{21} & X_{22} & X_{23} & X_{24} \\ X_{31} & X_{32} & X_{33} & X_{34} \\ X_{41} & X_{42} & X_{43} & X_{44} \\ ? & ? & ? & ? \\ ? & ? & ? & ? \\ ? & ? & ? & ? \\ ? & ? & ? & ? \end{pmatrix}$$

$$m = 8, k = 2$$

$$\begin{pmatrix} X_{11} & ? & X_{13} & ? \\ X_{21} & ? & ? & X_{24} \\ ? & X_{32} & X_{33} & ? \\ X_{41} & ? & ? & X_{44} \\ ? & ? & X_{53} & X_{54} \\ ? & X_{62} & X_{63} & ? \\ X_{71} & ? & X_{73} & ? \\ ? & ? & X_{83} & X_{84} \end{pmatrix}$$

1. Can we learn using the **same** total number of attributes?

Problem Definition

$$m = 4, k = 4$$

$$\begin{pmatrix} X_{11} & X_{12} & X_{13} & X_{14} \\ X_{21} & X_{22} & X_{23} & X_{24} \\ X_{31} & X_{32} & X_{33} & X_{34} \\ X_{41} & X_{42} & X_{43} & X_{44} \\ ? & ? & ? & ? \\ ? & ? & ? & ? \\ ? & ? & ? & ? \\ ? & ? & ? & ? \end{pmatrix}$$

$$m = 6, k = 2$$

$$\begin{pmatrix} X_{11} & ? & X_{13} & ? \\ X_{21} & ? & ? & X_{24} \\ ? & X_{32} & X_{33} & ? \\ X_{41} & ? & ? & X_{44} \\ ? & ? & X_{53} & X_{54} \\ ? & X_{62} & X_{63} & ? \\ ? & ? & ? & ? \\ ? & ? & ? & ? \end{pmatrix}$$

2. Can we learn using the **a smaller** total number of attributes?

We Focus on:

- Stochastic i.i.d. data
- Linear regression: $\hat{y}_t = \langle \mathbf{w}, \mathbf{x}_t \rangle$
 - Under 2-norm constraint: $\|\mathbf{w}\|_2 \leq B$
- Squared loss: $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

We Focus on:

- Stochastic i.i.d. data
- Linear regression: $\hat{y}_t = \langle \mathbf{w}, \mathbf{x}_t \rangle$
 - Under 2-norm constraint: $\|\mathbf{w}\|_2 \leq B$
- Squared loss: $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

Goal:

- Minimize excess risk: $\mathbb{E}[L_{\mathcal{D}}(\mathbf{w}_{\text{output}})] - L_{\mathcal{D}}(\mathbf{w}^*)$

- Goal: loss minimization - optimization problem
- Solution: Stochastic Gradient Descent (SGD)
- Algorithm template [Cesa-Bianchi et al. (2011)]:
 1. Scan the training set
 2. For each example build an unbiased gradient estimator based on sampling k attributes
 3. Feed into SGD

- Squared loss for \mathbf{x}_t :

$$\ell(\mathbf{w}; \mathbf{x}_t, y_t) = \frac{1}{2} (\langle \mathbf{w}, \mathbf{x}_t \rangle - y_t)^2$$

- Gradient for \mathbf{x}_t :

$$\nabla \ell(\mathbf{w}; \mathbf{x}_t, y_t) = (\langle \mathbf{w}, \mathbf{x}_t \rangle - y_t) \cdot \mathbf{x}_t$$

- Squared loss for \mathbf{x}_t :

$$\ell(\mathbf{w}; \mathbf{x}_t, y_t) = \frac{1}{2} (\langle \mathbf{w}, \mathbf{x}_t \rangle - y_t)^2$$

- Gradient for \mathbf{x}_t :

$$\nabla \ell(\mathbf{w}; \mathbf{x}_t, y_t) = (\langle \mathbf{w}, \mathbf{x}_t \rangle - y_t) \cdot \mathbf{x}_t$$

- Estimation:

- Sample $k - 1$ attributes with probabilities q_i and set

$$\tilde{\mathbf{x}}_t = \frac{1}{k-1} \sum_{r=1}^{k-1} \frac{1}{q_{i_t,r}} \mathbf{x}_t [i_t, r] \mathbf{e}_{i_t,r}, \text{ or}$$

$$\tilde{\mathbf{x}}_t = \frac{1}{k-1} (0, \dots, \frac{1}{q_{i_t,1}} \mathbf{x}_t [i_t, 1], 0, \dots, \frac{1}{q_{i_t,k-1}} \mathbf{x}_t [i_t, k-1], 0, \dots, 0)$$

- Sample 1 attribute with probability $p_{j_t} = w_{t,j_t}^2 / \|\mathbf{w}_t\|_2^2$ and set

$$\tilde{\phi}_t = \frac{w_{t,j_t}}{p_{j_t}} \mathbf{x}_t [j_t] - y_t$$

- Set $\tilde{\mathbf{g}}_t = \tilde{\phi}_t \cdot \tilde{\mathbf{x}}_t$

- Use uniform sampling: $q_i = \frac{1}{d} \quad \forall i \in [d]$
- Expected risk bound: $O\left(\sqrt{\frac{d}{km}}\right)$
 - Can learn using the same total number of attributes as the full-information scenario

- Use uniform sampling: $q_i = \frac{1}{d} \quad \forall i \in [d]$
- Expected risk bound: $O\left(\sqrt{\frac{d}{km}}\right)$
 - Can learn using the same total number of attributes as the full-information scenario

Can we do better?

- Use uniform sampling: $q_i = \frac{1}{d} \quad \forall i \in [d]$
- Expected risk bound: $O\left(\sqrt{\frac{d}{km}}\right)$
 - Can learn using the same total number of attributes as the full-information scenario

Can we do better?
No.

- Use uniform sampling: $q_i = \frac{1}{d} \quad \forall i \in [d]$
- Expected risk bound: $O\left(\sqrt{\frac{d}{km}}\right)$
 - Can learn using the same total number of attributes as the full-information scenario

Can we do better?
No.

- They also prove a corresponding lower bound

- How should we pick the q_i -s?

- How should we pick the q_i -s?
- Can show that the risk bound depends on

$$\mathbb{E} \left[\|\tilde{\mathbf{g}}_t\|_2^2 \right] \leq 4B^2 \left(\frac{1}{k-1} \sum_{i=1}^d \frac{1}{q_i} \mathbb{E} \left[x_i^2 \right] + 1 \right)$$

- How should we pick the q_i -s?
- Can show that the risk bound depends on

$$\mathbb{E} \left[\|\tilde{\mathbf{g}}_t\|_2^2 \right] \leq 4B^2 \left(\frac{1}{k-1} \sum_{i=1}^d \frac{1}{q_i} \mathbb{E} \left[x_i^2 \right] + 1 \right)$$

- Try to minimize as a function of the q_i -s

- How should we pick the q_i -s?
- Can show that the risk bound depends on

$$\mathbb{E} \left[\|\tilde{\mathbf{g}}_t\|_2^2 \right] \leq 4B^2 \left(\frac{1}{k-1} \sum_{i=1}^d \frac{1}{q_i} \mathbb{E} [x_i^2] + 1 \right)$$

- Try to minimize as a function of the q_i -s
- Solution:

$$q_i = \frac{\sqrt{\mathbb{E} [x_i^2]}}{\sum_{j=1}^d \sqrt{\mathbb{E} [x_j^2]}}.$$

- The optimal sampling probabilities can be calculated if second moments of attributes are known
- This is the DDAERR (Distribution-Dependent Attribute Efficient Ridge Regression) algorithm
 - Risk bound:

$$O\left(\sqrt{\frac{\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} + k}{km}}\right)$$

- By definition $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} = \left(\sum_{i=1}^d \sqrt{\mathbb{E}[\mathbf{x}_i^2]}\right)^2$
- Assuming $\|\mathbf{x}\|_2 \leq 1$:
 - $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} \leq d$
 - $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} = d$ only when all second moments are equal
 - On real-world data sets, typically $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} < d$
- The example in the lower bound satisfies $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} = d$

- By definition $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} = \left(\sum_{i=1}^d \sqrt{\mathbb{E}[\mathbf{x}_i^2]}\right)^2$
- Assuming $\|\mathbf{x}\|_2 \leq 1$:
 - $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} \leq d$
 - $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} = d$ only when all second moments are equal
 - On real-world data sets, typically $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} < d$
- The example in the lower bound satisfies $\|\mathbb{E}[\mathbf{x}^2]\|_{\frac{1}{2}} = d$

Conclusion

The DDAERR algorithm uses a **smaller** total amount of attributes compared to the full-information algorithms on real-world data sets.

- What if we don't know second moments?

- What if we don't know second moments?
- Solution - two phased method:
 1. Estimate second moments
 2. Use UCB-style moment estimates

- What if we don't know second moments?
- Solution - two phased method:
 1. Estimate second moments
 2. Use UCB-style moment estimates
- Estimation phase samples attributes uniformly
 - Run the AERR algorithm \Rightarrow no data is wasted

1-norm Constraint

- Again linear regression and squared loss, but this time with a different norm constraint: $\|\mathbf{w}\|_1 \leq B$

1-norm Constraint

- Again linear regression and squared loss, but this time with a different norm constraint: $\|\mathbf{w}\|_1 \leq B$
- Solution similar to the 2-norm constraint:
 - Based on the Exponentiated Gradient algorithm instead of the Stochastic Gradient Descent algorithm
 - Uses different sampling distributions:

$$q_i = \frac{\mathbb{E} [x_i^2]}{\sum_{j=1}^d \mathbb{E} [x_j^2]}$$

- Risk bound:

$$O\left(\sqrt{\frac{(\|\mathbb{E} [\mathbf{x}^2]\|_1 + k) \log d}{km}}\right)$$

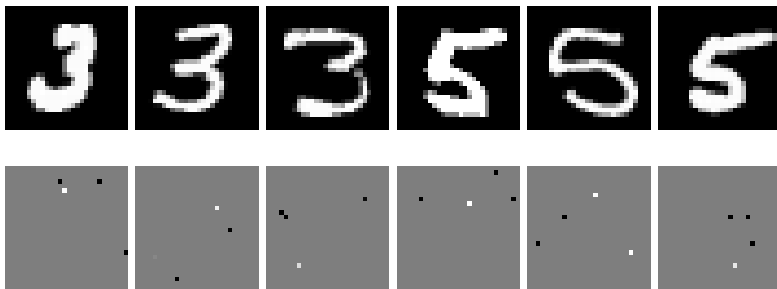
Experiment - MNIST Data Set

- Consist of gray scale hand-written digits
- Standard machine learning data set
- We focused on a classification problem: '3' vs. '5'

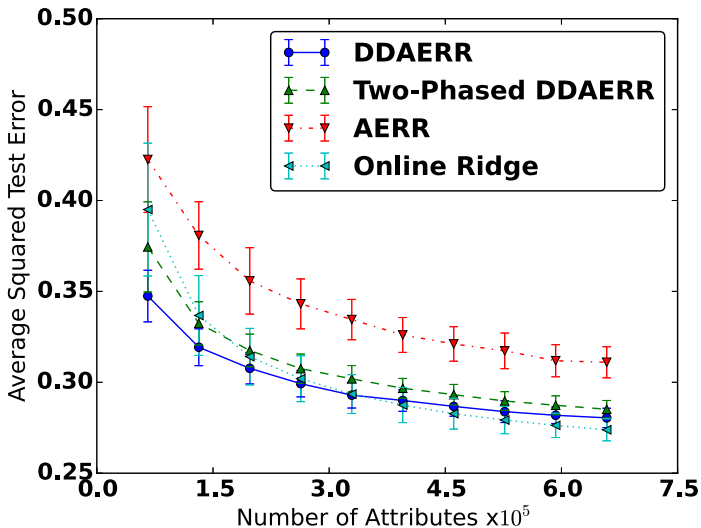


Experiment - MNIST Data Set

- Consist of gray scale hand-written digits
- Standard machine learning data set
- We focused on a classification problem: '3' vs. '5'

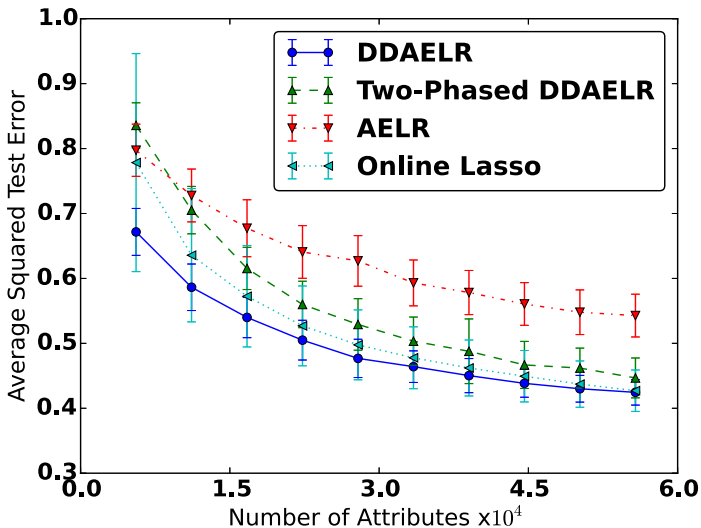


MNIST Data Set - 2-norm Constraint



Test error for the algorithms with $k = 56$ in the 2-norm scenario over the classification task "3" vs. "5" in the MNIST data set.

MNIST Data Set - 1-norm Constraint



Test error for the algorithms with $k = 4$ in the 1-norm scenario over the classification task "3" vs. "5" in the MNIST data set.

Summary

- Distribution-dependent sampling and bound
- Uses **less** attributes than full-information algorithms - useful also in general budgeted learning settings
- "Broken" lower bound for real world data sets
- Method can potentially be effective in other partial-information learning scenarios

Summary

- Distribution-dependent sampling and bound
- Uses **less** attributes than full-information algorithms - useful also in general budgeted learning settings
- "Broken" lower bound for real world data sets
- Method can potentially be effective in other partial-information learning scenarios

Thanks!