

Exact Incremental and Decremental Learning for LS-SVM

Wei-Han Lee¹, Bong Jun Ko¹, Shiqiang Wang¹, Changchang Liu¹, Kin K. Leung²

¹ IBM T. J. Watson Research Center, Yorktown Heights, NY, USA

² Imperial College London, UK

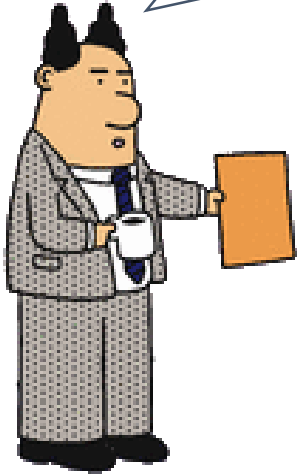
A day in the office of a machine learning (ML) engineer



The ML Engineer

A day in the office of a machine learning (ML) engineer

Here's the
dataset you've
been waiting for.



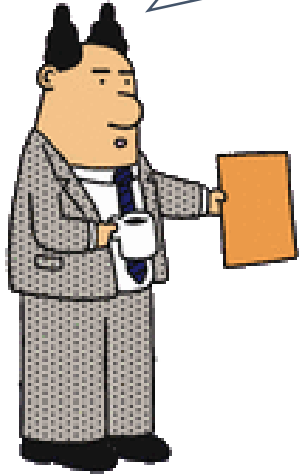
The Manager



The ML Engineer

A day in the office of a machine learning (ML) engineer

Here's the dataset you've been waiting for.



The Manager

Great! Now I can build a model.



The ML Engineer

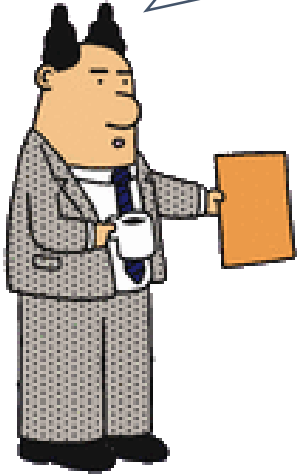


Processing



A day in the office of a machine learning (ML) engineer

Here's the dataset you've been waiting for.



The Manager

Done!



The ML Engineer



The next day

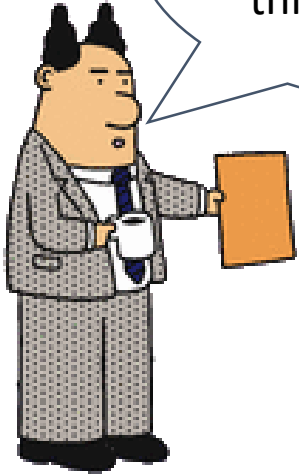


The ML Engineer



The next day

Here's another dataset we just acquired. We need this reflected in the model too.



The Manager

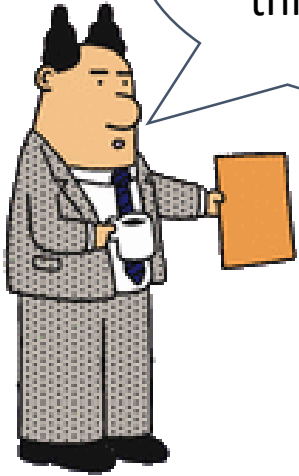


The ML Engineer



The next day

Here's another dataset we just acquired. We need this reflected in the model too.



The Manager

Ok, no problem. I'll just need to add it to the database and run the training job again, so the model doesn't forget the old dataset.



The ML Engineer

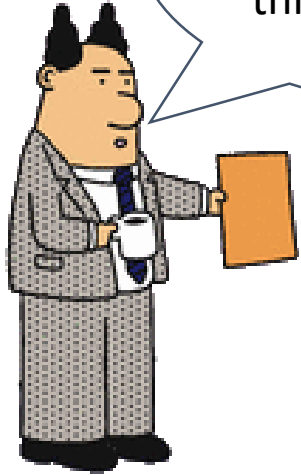


Processing



The next day

Here's another dataset we just acquired. We need this reflected in the model too.



The Manager

Done! It took a little longer, but not a big deal...

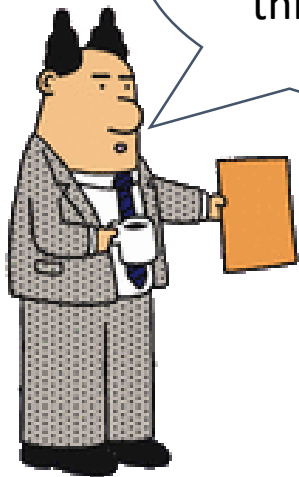


The ML Engineer



And so it goes on...

Here's **yet another** dataset we just acquired. We need this reflected in the model too.



The Manager

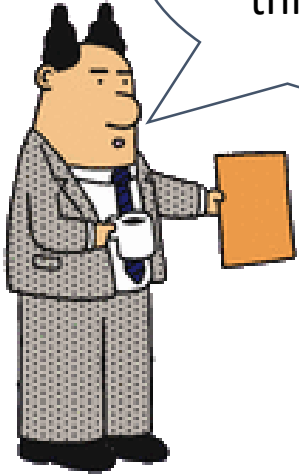


The ML Engineer



And so it goes on...

Here's **yet another** dataset we just acquired. We need this reflected in the model too.



The Manager

Aarggggh! Not again...



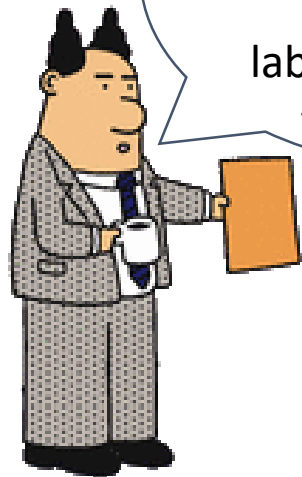
The ML Engineer



Processing



And yet, even worse...



The Manager

Oh, by the way,
remember the data
a few weeks ago?
Turns out their
labels are wrong, so
take them out.

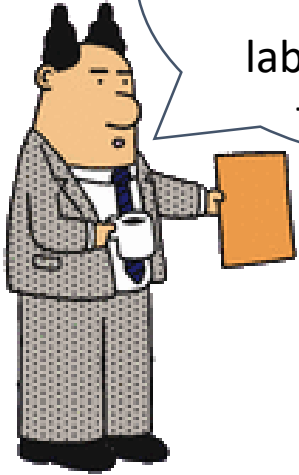


The ML Engineer



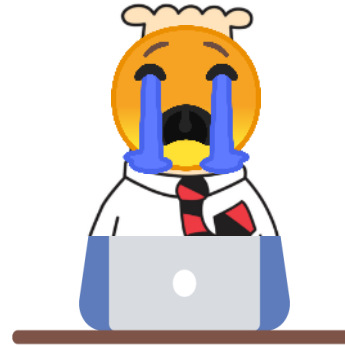
And yet, even worse...

Oh, by the way,
remember the data
a few weeks ago?
Turns out their
labels are wrong, so
take them out.



The Manager

What???
Nooooooooooooo!!



The ML Engineer



Processing



How to update the model after adding, correcting, or removing data, without having to retrain the entire model?

- Add data
 - Incremental learning
- Remove data
 - Decremental learning
 - Why? – e.g., erroneous data, per user request (GDPR regulation)
- Correction to data
 - Both incremental and decremental learning

State of the Art

- Incremental learning has been considered in various contexts
 - Heuristics for deep neural networks to avoid catastrophic forgetting
 - Exact update methods for shallow models such as support vector machine (SVM)
- A few existing approaches on decremental learning
 - Not for neural networks
 - For SVM but requires old data
- We propose an exact (provably optimal) method for joint incremental and decremental learning for least squares SVM (LS-SVM) that does not require old data

State of the Art (Details)

	Model	Exact?	Requires old data?	Decremental learning?
[1], [2], [3]	SVM	No	Yes	No
[4]	SVM	Yes	Yes	No
[5]	SVM	Yes	Yes	Yes
[6], [7]	SVM	No	No	No
[8], [9]	LS-SVM	Yes	Yes	No
[10], [11]	LS-SVM	No	No	No
Our work	LS-SVM	Yes	No	Yes

[1] Syed, N., Liu, H., and Sung, K. Incremental learning with support vector machines. In International Joint Conference on Artificial Intelligence, 1999.

[2] Ruping, S. Incremental learning with support vector machines. In IEEE International Conference on Data Mining, 2001.

[3] Carlotta Domeniconi and Dimitrios Gunopulos. Incremental support vector machine construction, in IEEE International Conference on Data Mining, 2001.

[4] Diehl, C. P. and Cauwenberghs, G. SVM incremental learning, adaptation and optimization. In International Joint Conference on Neural Networks, 2003

[5] Cauwenberghs, G. and Poggio, T. Incremental and decremental support vector machine learning. In Advances in neural information processing systems, 2001.

[6] Yi-Min Wen and Bao-Liang Lu. Incremental learning of support vector machines by classifier combining, in Pacific-Asia Conference on Knowledge Discovery and Data Mining, 2007.

[7] Youlu Xing, Furao Shen, Chaomin Luo, and Jinxi Zhao. L3-SVM: a lifelong learning method for SVM. In IEEE IJCNN, 2015

[8] Hoi-Ming Chi and Okan K Ersoy. Recursive update algorithm for least squares support vector machines. Neural Processing Letters, 2003

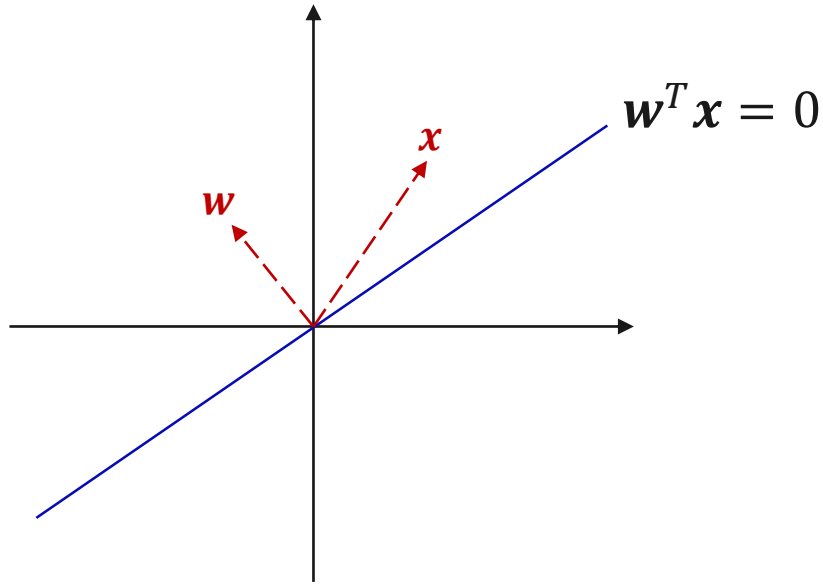
[9] Zhifeng Hao, Shu Yu, Xiaowei Yang, Feng Zhao, Rong Hu, and Yanchun Liang. Online LS-SVM learning for classification problems based on incremental chunk. In International Symposium on Neural Networks, 2004.

[10] Yaakov Engel, Shie Mannor, and Ron Meir. The kernel recursive least squares algorithm. IEEE Trans. On Signal Processing, 2004.

[11] Ling Jian, Shuqian Shen, Jundong Li, Xijun Liang, and Lei Li. Budget online learning algorithm for least squares SVM. IEEE Trans. on Neural Networks and Learning Systems, 2017.



Support Vector Machine (SVM)



- If $w^T x > 0$, classify as label +1
- If $w^T x < 0$, classify as label -1

- Binary classifier
- Can be converted into multi-class classifier using one-versus-all or one-versus-one ensemble approaches

Least-Squares Support Vector Machine (LS-SVM)

- Optimal solution for model parameter

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^J} \underbrace{\rho \|\mathbf{w}\|^2 + \sum_{n=1}^N (\mathbf{w}^T \vec{\phi}(\mathbf{x}_n) - y_n)^2}_{\text{Learning loss function}}$$

Regularization constant

Data sample transformed into feature space

Ground-truth label

$$= \mathbf{\Phi}[\mathbf{\Phi}^T \mathbf{\Phi} + \rho \mathbf{I}_N]^{-1} \mathbf{y} = [\rho \mathbf{I}_J + \mathbf{\Phi} \mathbf{\Phi}^T]^{-1} \mathbf{\Phi} \mathbf{y} \quad (\text{Analytical optimal solution})$$

Main Result for Incremental/Decremental Learning

Theorem 1. *For a given model*

$$\mathbf{w} = \Phi[\Phi^T \Phi + \rho \mathbf{I}_N]^{-1} \mathbf{y} \quad (3)$$

and auxiliary matrix

$$\mathbf{C} = \Phi[\Phi^T \Phi + \rho \mathbf{I}_N]^{-1} \Phi^T, \quad (4)$$

when adding new training data $(\mathbf{X}_a, \mathbf{y}_a)$ and removing existing training data $(\mathbf{X}_r, \mathbf{y}_r)$, we can compute the new values of \mathbf{w} and \mathbf{C} using

Model updating equations, only involving added/removed data samples and auxiliary matrix \mathbf{C}

$$\left\{ \begin{array}{l} \mathbf{w}_{\text{new}} = \mathbf{w} + (\mathbf{C} - \mathbf{I}_J) \Phi_c \left(\rho \mathbf{I} - \Phi_c'^T (\mathbf{C} - \mathbf{I}_J) \Phi_c \right)^{-1} \left(\Phi_c'^T \mathbf{w} - \mathbf{y}_c \right) \end{array} \right. \quad (5)$$

$$\left\{ \begin{array}{l} \mathbf{C}_{\text{new}} = \mathbf{C} + (\mathbf{C} - \mathbf{I}_J) \Phi_c \left(\rho \mathbf{I} - \Phi_c'^T (\mathbf{C} - \mathbf{I}_J) \Phi_c \right)^{-1} \Phi_c'^T (\mathbf{C} - \mathbf{I}_J) \end{array} \right. \quad (6)$$

where we define $\Phi_c = (\Phi_a, \Phi_r)$, $\Phi_c' = (\Phi_a, -\Phi_r)$, and $\mathbf{y}_c = (\mathbf{y}_a, -\mathbf{y}_r)$.

Algorithm

- Initial model training

- Compute $\mathbf{w} = \Phi[\Phi^T \Phi + \rho \mathbf{I}_N]^{-1} \mathbf{y}$ and $\mathbf{C} = \Phi[\Phi^T \Phi + \rho \mathbf{I}_N]^{-1} \Phi^T$

- Incremental/decremental learning using new/removed data

- Compute

- $$\mathbf{w}_{\text{new}} = \mathbf{w} + (\mathbf{C} - \mathbf{I}_J) \Phi_c \left(\rho \mathbf{I} - \Phi_c'^T (\mathbf{C} - \mathbf{I}_J) \Phi_c \right)^{-1} \left(\Phi_c'^T \mathbf{w} - \mathbf{y}_c \right)$$

- $$\mathbf{C}_{\text{new}} = \mathbf{C} + (\mathbf{C} - \mathbf{I}_J) \Phi_c \left(\rho \mathbf{I} - \Phi_c'^T (\mathbf{C} - \mathbf{I}_J) \Phi_c \right)^{-1} \Phi_c'^T (\mathbf{C} - \mathbf{I}_J)$$

- **Complexity:** $O(L^3 + JL^2 + J^2L + J^3)$

Number of added/removed
data samples in a “batch”

Dimension of
feature vector

Algorithm

- Initial model training

- Compute $\mathbf{w} = \Phi[\Phi^T\Phi + \rho\mathbf{I}_N]^{-1}\mathbf{y}$ and $\mathbf{C} = \Phi[\Phi^T\Phi + \rho\mathbf{I}_N]^{-1}\Phi^T$

- Incremental/decremental learning using new/removed data

- Compute

$$\mathbf{w}_{\text{new}} = \mathbf{w} + (\mathbf{C} - \mathbf{I}_J)\Phi_c \left(\rho\mathbf{I} - \Phi_c'^T(\mathbf{C} - \mathbf{I}_J)\Phi_c \right)^{-1} \left(\Phi_c'^T\mathbf{w} - \mathbf{y}_c \right)$$

$$\mathbf{C}_{\text{new}} = \mathbf{C} + (\mathbf{C} - \mathbf{I}_J)\Phi_c \left(\rho\mathbf{I} - \Phi_c'^T(\mathbf{C} - \mathbf{I}_J)\Phi_c \right)^{-1} \Phi_c'^T(\mathbf{C} - \mathbf{I}_J)$$

- **Complexity:** $O(L^3 + JL^2 + J^2L + J^3)$

Number of added/removed data samples in a “batch”

Dimension of feature vector



Multiple data sample updates can be done either in one or multiple batches

- When using one sample per batch
 - Complexity is **linear** in the number updated samples
 - May not be best in practice though due to efficient implementations of matrix multiplication (details in experiments)

Special Cases

- One data sample in a batch, either incremental or decremental

Corollary 1 (Incremental learning of a single data sample).
We can update \mathbf{w} and \mathbf{C} to include the influence of a new training data sample $(\mathbf{x}_{N+1}, y_{N+1})$ using

$$\mathbf{w}_{\text{new}} = \mathbf{w} + \frac{(\mathbf{C} - \mathbf{I}_J)\vec{\phi}(\mathbf{x}_{N+1})(\vec{\phi}(\mathbf{x}_{N+1})^T \mathbf{w} - y_{N+1})}{\vec{\phi}(\mathbf{x}_{N+1})^T \vec{\phi}(\mathbf{x}_{N+1}) + \rho - \vec{\phi}(\mathbf{x}_{N+1})^T \mathbf{C} \vec{\phi}(\mathbf{x}_{N+1})}$$
$$\mathbf{C}_{\text{new}} = \mathbf{C} + \frac{(\mathbf{C} - \mathbf{I}_J)\vec{\phi}(\mathbf{x}_{N+1})\vec{\phi}(\mathbf{x}_{N+1})^T(\mathbf{C} - \mathbf{I}_J)}{\vec{\phi}(\mathbf{x}_{N+1})^T \vec{\phi}(\mathbf{x}_{N+1}) + \rho - \vec{\phi}(\mathbf{x}_{N+1})^T \mathbf{C} \vec{\phi}(\mathbf{x}_{N+1})}.$$

Corollary 2 (Decremental learning of a single data sample).
We can update \mathbf{w} and \mathbf{C} to remove the influence of an existing training data sample (\mathbf{x}_r, y_r) using

$$\mathbf{w}_{\text{new}} = \mathbf{w} - \frac{(\mathbf{C} - \mathbf{I}_J)\vec{\phi}(\mathbf{x}_r)(\vec{\phi}(\mathbf{x}_r)^T \mathbf{w} - y_r)}{-\vec{\phi}(\mathbf{x}_r)^T \vec{\phi}(\mathbf{x}_r) + \rho + \vec{\phi}(\mathbf{x}_r)^T \mathbf{C} \vec{\phi}(\mathbf{x}_r)}$$
$$\mathbf{C}_{\text{new}} = \mathbf{C} - \frac{(\mathbf{C} - \mathbf{I}_J)\vec{\phi}(\mathbf{x}_r)\vec{\phi}(\mathbf{x}_r)^T(\mathbf{C} - \mathbf{I}_J)}{-\vec{\phi}(\mathbf{x}_r)^T \vec{\phi}(\mathbf{x}_r) + \rho + \vec{\phi}(\mathbf{x}_r)^T \mathbf{C} \vec{\phi}(\mathbf{x}_r)}.$$

Special Cases

- Multiple data samples in a batch, either incremental or decremental

Corollary 3 (Incremental learning of a batch of data samples). *We can update w and C to include the influence of a batch of new training data samples (X_a, y_a) using*

$$w_{\text{new}} = w + (C - I_J)\Phi_a \left(\rho I - \Phi_a^T (C - I_J)\Phi_a \right)^{-1} \left(\Phi_a^T w - y_a \right)$$

$$C_{\text{new}} = C + (C - I_J)\Phi_a \left(\rho I - \Phi_a^T (C - I_J)\Phi_a \right)^{-1} \Phi_a^T (C - I_J)$$

where $\Phi_a = \Phi(X_a)$.

Corollary 4 (Decremental learning of a batch of data samples). *We can update w and C to remove the influence of a batch of existing training data samples (X_r, y_r) using*

$$w_{\text{new}} = w - (C - I_J)\Phi_r \left(\rho I + \Phi_r^T (C - I_J)\Phi_r \right)^{-1} \left(\Phi_r^T w - y_r \right)$$

$$C_{\text{new}} = C - (C - I_J)\Phi_r \left(\rho I + \Phi_r^T (C - I_J)\Phi_r \right)^{-1} \Phi_r^T (C - I_J)$$

where $\Phi_r = \Phi(X_r)$.

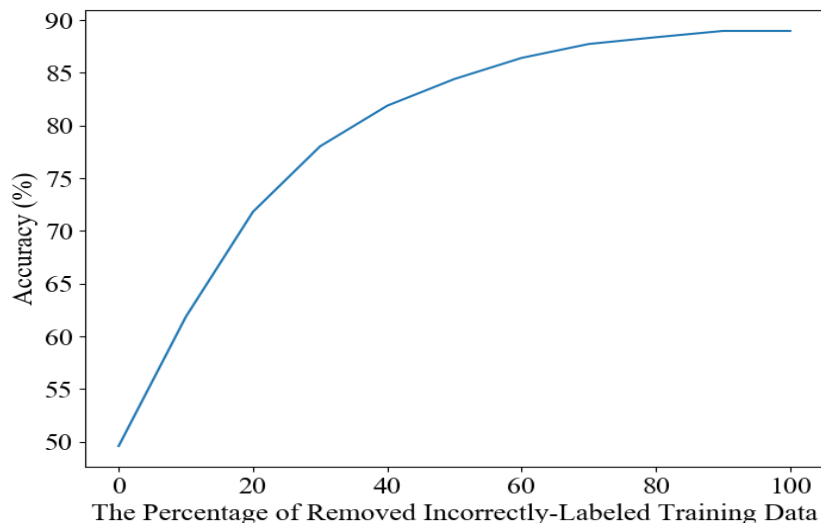
Experiments

- MNIST dataset of handwritten digits
- Classify even/odd digits using LS-SVM
- 60,000 training data samples
- 10,000 test data samples

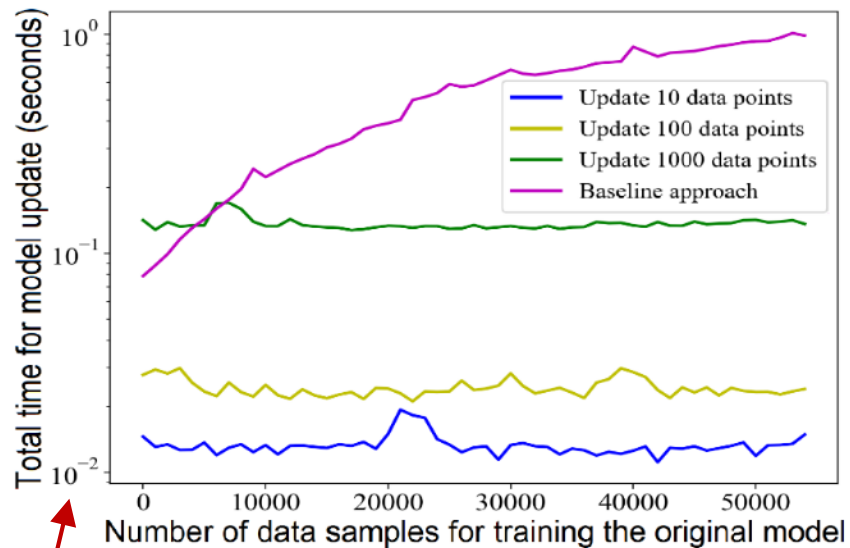


Training with Incorrectly Labeled Data

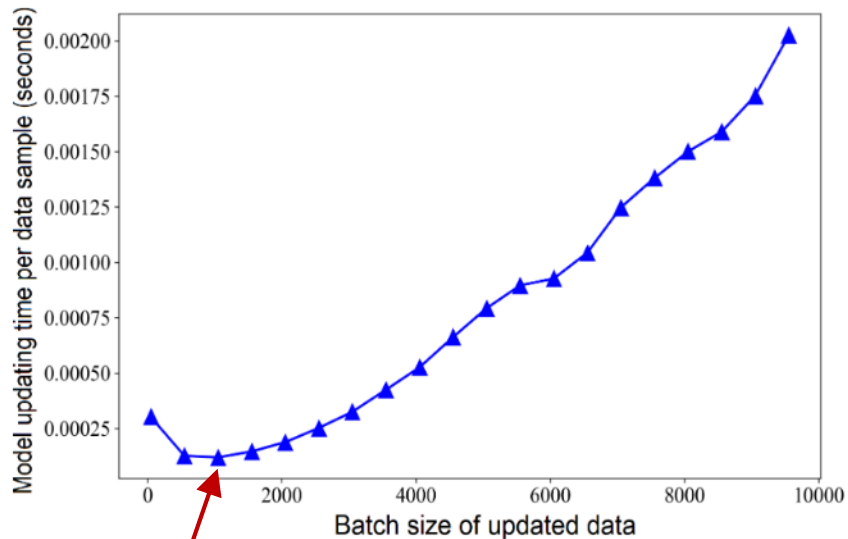
- Initially, 50% of data samples are mislabeled
- Removing 10% of mislabeled data increases accuracy by 12%
 - Retraining the entire model takes **0.81 seconds** (on a personal laptop)
 - Updating using our proposed approach takes **0.26 seconds**



Update Time with Different Batch Sizes



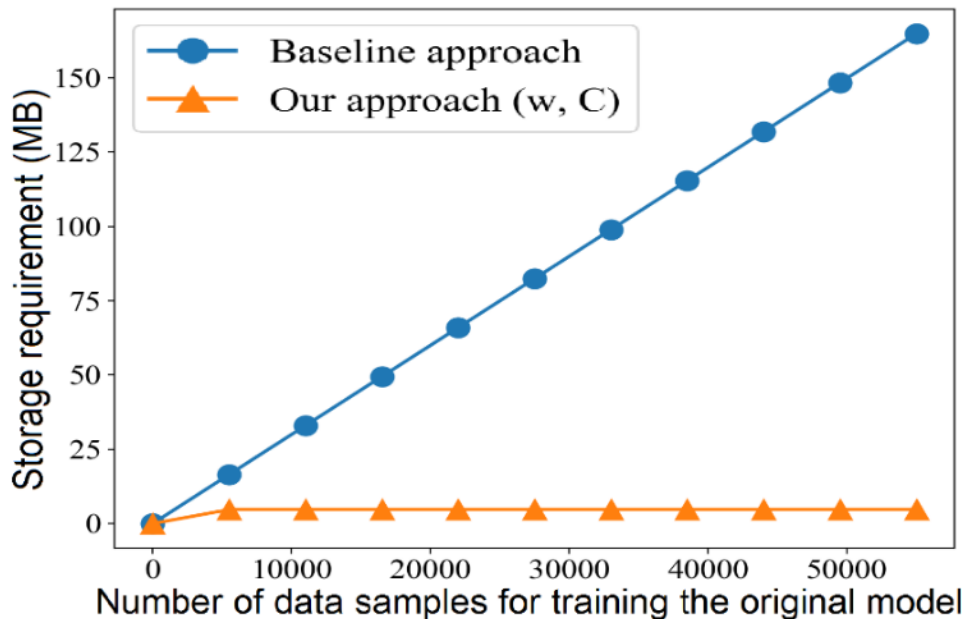
Log-scale



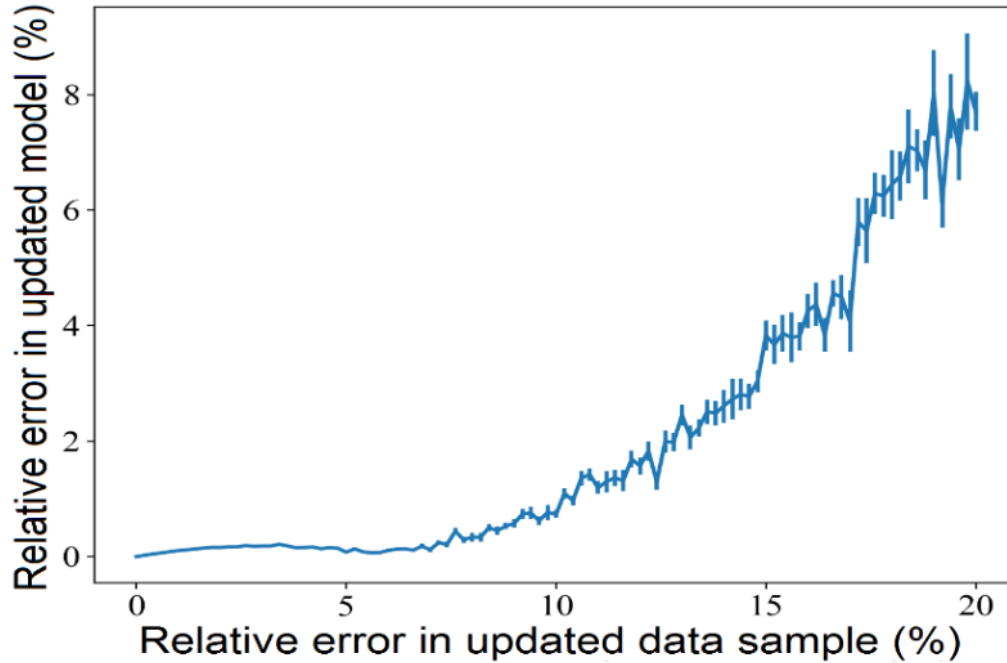
Optimal batch size from experimental (empirical) evaluation

Storage Saving

- We do not need to store original data



When we don't exactly know the data samples to be removed...



(Removing one data sample)

Summary

- What we have done...
 - Provably optimal incremental and decremental learning for LS-SVM
 - Does not require old data
 - Only requires an auxiliary matrix and data samples that are added/removed
 - **Reduces** model update **time** and **storage**, compared to retraining
 - **Preserves privacy** of training data when sharing updatable models
- What remains to be done...
 - Incremental and decremental learning (particularly decremental learning) for generic models such as deep neural networks
 - How to properly define decremental learning?
 - Algorithms for joint incremental and decremental learning
 - Anything provable?
 - Any intuitive insights?

Thank you

Q & A



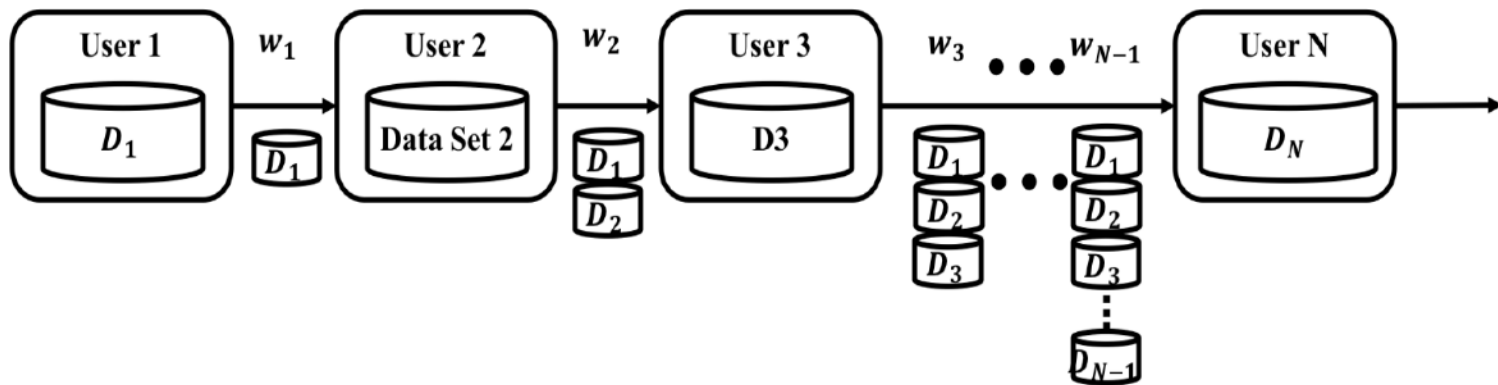
Backup slides



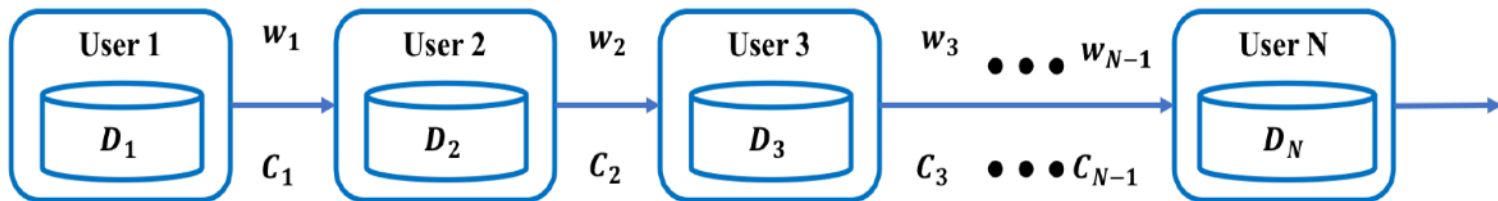
Application

- Decentralized learning

Data Sharing

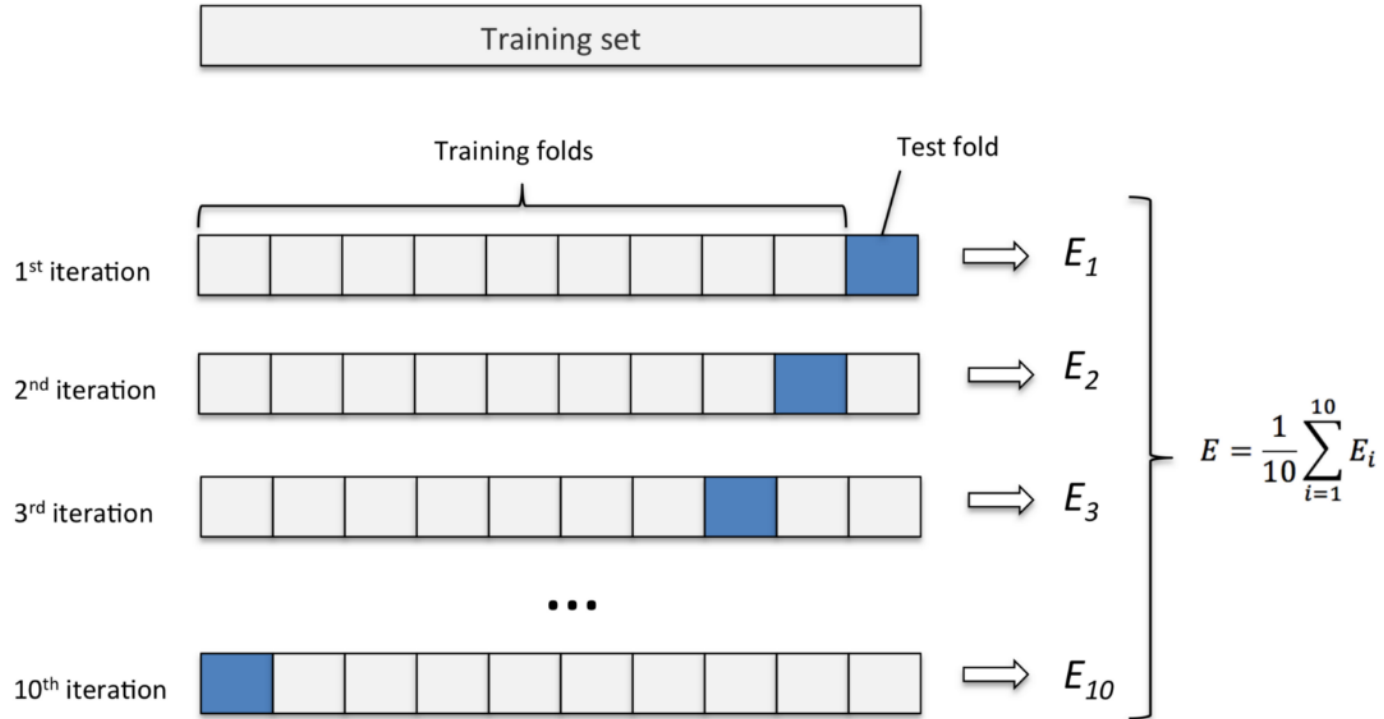


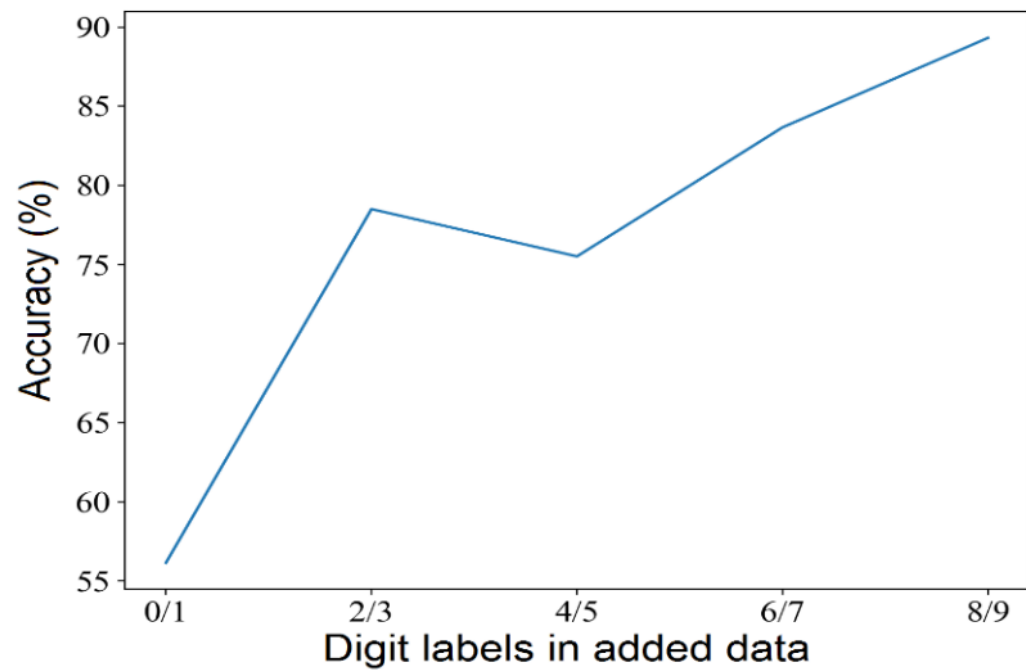
Our Method



Application

- k-fold cross validation (incrementally replace one fold)





Proof. For the updated training dataset, we get $\Phi_{\text{new}} = (\Phi, \Phi_a | \Phi_r)$ and $\mathbf{y}_{\text{new}} = (\mathbf{y}, \mathbf{y}_a | \mathbf{y}_r)$, where $(\mathbf{Z}_1 | \mathbf{Z}_2)$ denotes removing the columns (of a matrix) or elements (of a vector) in \mathbf{Z}_2 from \mathbf{Z}_1 . Using (2) and (3), we can compute the new model \mathbf{w}_{new} as

$$\begin{aligned} \mathbf{w}_{\text{new}} &= \Phi_{\text{new}} [\Phi_{\text{new}}^T \Phi_{\text{new}} + \rho \mathbf{I}_{N+N_a-N_r}]^{-1} \mathbf{y}_{\text{new}} \\ &= [\rho \mathbf{I}_J + \Phi_{\text{new}} \Phi_{\text{new}}^T]^{-1} \Phi_{\text{new}} \mathbf{y}_{\text{new}} \\ &= [\rho \mathbf{I}_J + \Phi \Phi^T + \Phi_a \Phi_a^T - \Phi_r \Phi_r^T]^{-1} [\Phi \mathbf{y} + \Phi_a \mathbf{y}_a - \Phi_r \mathbf{y}_r] \\ &= [\rho \mathbf{I}_J + \Phi \Phi^T + \Phi_c \Phi_c'^T]^{-1} [\Phi \mathbf{y} + \Phi_c \mathbf{y}_c] \end{aligned} \quad (7)$$

where we note that $\Phi_c = (\Phi_a, \Phi_r)$, $\Phi_c' = (\Phi_a, -\Phi_r)$ and $\mathbf{y}_c = (\mathbf{y}_a, -\mathbf{y}_r)$ by definition.

According to the Woodbury matrix identity [24], we have $(\mathbf{A} + \mathbf{U}\mathbf{B}\mathbf{V})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{B}\mathbf{V}\mathbf{A}^{-1}\mathbf{U})\mathbf{B}\mathbf{V}\mathbf{A}^{-1}$. We define $\mathbf{A} = \rho \mathbf{I}_J + \Phi \Phi^T$, $\mathbf{B} = 1$, $\mathbf{U} = \Phi_c$, $\mathbf{V} = \Phi_c'^T$, and $\Psi = \mathbf{I} + \Phi_c'^T \mathbf{A}^{-1} \Phi_c$. From (7), we can obtain the following updating process for \mathbf{w} :

$$\begin{aligned} \mathbf{w}_{\text{new}} &= (\mathbf{A}^{-1} - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T \mathbf{A}^{-1}) [\Phi \mathbf{y} + \Phi_c \mathbf{y}_c] \\ &= \mathbf{w} + \mathbf{A}^{-1} \Phi_c \mathbf{y}_c - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T \mathbf{w} \\ &\quad - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T \mathbf{A}^{-1} \Phi_c \mathbf{y}_c \\ &= \mathbf{w} + \mathbf{A}^{-1} \Phi_c \Psi^{-1} \mathbf{y}_c - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T \mathbf{w} \\ &= \mathbf{w} - \mathbf{A}^{-1} \Phi_c \Psi^{-1} (\Phi_c'^T \mathbf{w} - \mathbf{y}_c) \\ &= \mathbf{w} + (\mathbf{C} - \mathbf{I}_J) \Phi_c (\rho \mathbf{I} + \Phi_c'^T (\mathbf{C} - \mathbf{I}_J) \Phi_c)^{-1} (\Phi_c'^T \mathbf{w} - \mathbf{y}_c). \end{aligned}$$

The last equality holds because

$$\begin{aligned} \rho \mathbf{A}^{-1} + \mathbf{C} &= \rho [\rho \mathbf{I}_J + \Phi \Phi^T]^{-1} + \Phi [\rho \mathbf{I}_N + \Phi^T \Phi]^{-1} \Phi^T \\ &= \rho [\rho \mathbf{I}_J + \Phi \Phi^T]^{-1} + [\rho \mathbf{I}_J + \Phi \Phi^T]^{-1} \Phi \Phi^T \\ &= [\rho \mathbf{I}_J + \Phi \Phi^T]^{-1} [\rho \mathbf{I}_J + \Phi \Phi^T] = \mathbf{I}. \end{aligned}$$

Similarly, we can compute \mathbf{C}_{new} as

$$\begin{aligned} \mathbf{C}_{\text{new}} &= \Phi_{\text{new}} [\Phi_{\text{new}}^T \Phi_{\text{new}} + \rho \mathbf{I}_{N+N_a-N_r}]^{-1} \Phi_{\text{new}}^T \\ &= [\rho \mathbf{I}_J + \Phi_{\text{new}} \Phi_{\text{new}}^T]^{-1} \Phi_{\text{new}} \Phi_{\text{new}}^T \\ &= [\mathbf{A} + \Phi_a \Phi_a^T - \Phi_r \Phi_r^T]^{-1} [\Phi \Phi^T + \Phi_a \Phi_a^T - \Phi_r \Phi_r^T] \\ &= [\mathbf{A} + \Phi_c \Phi_c'^T]^{-1} [\Phi \Phi^T + \Phi_c \Phi_c'^T] \\ &= (\mathbf{A}^{-1} - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T \mathbf{A}^{-1}) [\Phi \Phi^T + \Phi_c \Phi_c'^T] \\ &= \mathbf{C} + \mathbf{A}^{-1} \Phi_c \Phi_c'^T - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T \mathbf{C} \\ &\quad - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T \mathbf{A}^{-1} \Phi_c \Phi_c'^T \\ &= \mathbf{C} + \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T \mathbf{C} \\ &= \mathbf{C} - \mathbf{A}^{-1} \Phi_c \Psi^{-1} \Phi_c'^T (\mathbf{C} - \mathbf{I}_J) \\ &= \mathbf{C} + (\mathbf{C} - \mathbf{I}_J) \Phi_c (\rho \mathbf{I} + \Phi_c'^T (\mathbf{C} - \mathbf{I}_J) \Phi_c)^{-1} \Phi_c'^T (\mathbf{C} - \mathbf{I}_J). \end{aligned}$$