

# Supplementary Material of “Online Learning with Noisy Labels in Streaming Features”

Jinjie Qiu, Shengda Zhuo, Shuqiang Huang, Changdong Wang, *Member, IEEE* and Philip S. Yu, *Fellow, IEEE*

**Outline.** This document supplements the manuscript in the following aspects. Section 1 introduces compared methods, and Section 2 presents the complete experimental results that supplement.

## 1 Compared Methods

This section presents complete descriptions of methods which used for comparison analysis.

**Compared Methods.** We compared OSLMF with 13 related state-of-the-art online learning (OL), offline learning (OFL), and online noisy labels learning (ONL<sup>2</sup>).

- FOBOS (Singer and Duchi 2009) (OL): is an online learning algorithm that integrates gradient descent with regularization, designed to update models incrementally and promote sparsity by penalizing non-essential features.
- RDA (Singer and Duchi 2009) (OL): offers a base framework for direct training of online learners on observed features, utilizing a projected subgradient approach to enforce sparsity and eliminate features with negligible coefficients.
- FTRL (Singer and Duchi 2009) (OL): is an online learning algorithm that uses regularization techniques to accumulate all past gradients for parameter updates, which balances between aggressive learning and maintaining stability, especially in sparse data environments.
- FTRL-P (Promixal) (Singer and Duchi 2009) (OL): is an online learning algorithm that extends FTRL by incorporating proximal regularization, improving convergence on sparse features and providing robustness against feature correlation, particularly in high-dimensional spaces.
- OVFM (He et al. 2021) (OL): manages mixed data streams by using Gaussian cointegration modeling to link observed and distributed data, enhancing model convergence despite mixed data challenges.
- OLI<sup>2</sup>DS (You et al. 2023) (OL): optimizing feature selection via empirical risk minimization and using dynamic cost strategies for class imbalance to improve model efficacy.
- OSLMF (Wu et al. 2023) (OL): applying Gaussian Copula for imputation and Density Peak Clustering for label reconstruction, thus enriching the model with extensive label information.
- Co-teaching+ (Yu et al. 2019) (OFL): addresses noisy label learning in offline settings, using a dual-model approach to pick samples with small loss and discrepant predictions for mutual updates, focusing on reliable samples to mitigate label noise impact.
- PHuber (Menon et al. 2019) (OFL): employs robustness regularization in offline label noise learning to prevent overfitting to noise by modifying gradients, minimizing noisy label interference in the model.
- BeyondImages (Bey.Ima.) (Zhu, Wang, and Liu 2022) (OFL): tackles label noise in offline learning by correcting loss through an information-theoretic method that extracts key features from low-dimensional representations and estimates a noise transition matrix for model compensation, diminishing noisy label interference.
- Adastream (Zhang et al. 2022) (ONL<sup>2</sup>): innovates in data stream mining under conditions of incomplete and noisy labels, using a landmark assumption to estimate the noise transition matrix and leveraging model reuse ensembles to bolster method stability.
- OPA (Crammer et al. 2006) (ONL<sup>2</sup>): is an online learning algorithm using first-order linear learning with a threshold for sample training to lessen sensitivity to noisy labels and prevent overfitting to noise distributions.
- SCW (Wang, Zhao, and Hoi 2012) (ONL<sup>2</sup>): is an online learning algorithm that employs second-order linear learning and threshold constraints to reduce noise sensitivity and overfitting in model training.

## 2 Complete Experiments for noise-label data streams

This section presents additional experimental results that were precluded of the main paper because of page limitation. These results include:

1. The complete CAR results yielded from all 14 datasets, which is shown in table 1, table2, table3. The results are collected from experiments in the setting of symmetric noisy labeled, asymmetric noise labeled and flip noisy labeled data stream, compared with online learning algorithms.
2. The complete CAR trends yielded from all 14 datasets on symmetric noisy labeled, asymmetric noise labeled and

flip noisy labeled data stream, shown in Figure 1, Figure2, and Figure3. The results are collected from experiments in the setting of symmetric noisy labeled, asymmetric noise labeled and flip noisy labeled data stream, compared with online learning algorithms.

3. The complete CAR results yielded from all 14 datasets, which is shown in table 4, table5, and table6. The results are collected from experiments in the setting of symmetric noisy labeled, asymmetric noise labeled and flip noisy labeled data stream, compared with noisy labels learning algorithms.
4. The complete CAR trends yielded from all 14 datasets on symmetric noisy labeled, asymmetric noise labeled, and flip noisy labeled data stream, shown in Figure 4, Figure5, and Figure6. The results are collected from experiments in the setting of symmetric noisy labeled, asymmetric noise labeled, and flip noisy labeled data stream, compared with noisy labels learning algorithms.

## References

- Crammer, K.; Dekel, O.; Keshet, J.; Shalev-Shwartz, S.; and Singer, Y. 2006. Online passive aggressive algorithms.
- He, Y.; Dong, J.; Hou, B.-J.; Wang, Y.; and Wang, F. 2021. Online Learning in Variable Feature Spaces with Mixed Data. In *ICDM*, 181–190. IEEE.
- Menon, A. K.; Rawat, A. S.; Reddi, S. J.; and Kumar, S. 2019. Can gradient clipping mitigate label noise? In *ICLR*.
- Singer, Y.; and Duchi, J. C. 2009. Efficient learning using forward-backward splitting. *Advances in Neural Information Processing Systems*, 22.
- Wang, J.; Zhao, P.; and Hoi, S. C. 2012. Exact soft confidence-weighted learning. *arXiv preprint arXiv:1206.4612*.
- Wu, D.; Zhuo, S.; Wang, Y.; Chen, Z.; and He, Y. 2023. Online semi-supervised learning with mix-typed streaming features. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 4720–4728.
- You, D.; Xiao, J.; Wang, Y.; Yan, H.; Wu, D.; Chen, Z.; Shen, L.; and Wu, X. 2023. Online Learning From Incomplete and Imbalanced Data Streams. *IEEE Transactions on Knowledge and Data Engineering*.
- Yu, X.; Han, B.; Yao, J.; Niu, G.; Tsang, I.; and Sugiyama, M. 2019. How does disagreement help generalization against label corruption? In *International Conference on Machine Learning*, 7164–7173. PMLR.
- Zhang, Z.-Y.; Qian, Y.-Y.; Zhang, Y.-J.; Jiang, Y.; and Zhou, Z.-H. 2022. Adaptive Learning for Weakly Labeled Streams. In *KDD*, 2556–2564.
- Zhu, Z.; Wang, J.; and Liu, Y. 2022. Beyond images: Label noise transition matrix estimation for tasks with lower-quality features. In *International Conference on Machine Learning*, 27633–27653. PMLR.

Table 1: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR  $\pm$  Standard Variance) for 10 data sets in the case of Symmetric Noisy labeled data. • indicates the cases that our method loses the comparison. OOT: Out Of Time (24 hours).

	Symetric Noise $\rho_{-1} = \rho_{+1} = 0.4$								
Datasets	FOBOS	RDA	TG	FTRL	FTRL-P	OSLMF	OVFM	OLI <sup>2</sup> DS	ONLSF
wdbc	.614 $\pm$ .017	.608 $\pm$ .028	.602 $\pm$ .022	.547 $\pm$ .022	.536 $\pm$ .027	.428 $\pm$ .022	.578 $\pm$ .045	.542 $\pm$ .033	<b>.642 <math>\pm</math> .012</b>
breast	.726 $\pm$ .028	.629 $\pm$ .014	.798 $\pm$ .040	.806 $\pm$ .031	.651 $\pm$ .025	.596 $\pm$ .015	.653 $\pm$ .041	.522 $\pm$ .043	<b>.895 <math>\pm</math> .020</b>
dermatology	.592 $\pm$ .027	.607 $\pm$ .008	.594 $\pm$ .010	.682 $\pm$ .027•	.618 $\pm$ .031	.604 $\pm$ .054	.603 $\pm$ .012	.649 $\pm$ .020	.665 $\pm$ .020
wdbc	.586 $\pm$ .013	.587 $\pm$ .005	.582 $\pm$ .017	.587 $\pm$ .030	.572 $\pm$ .023	.598 $\pm$ .031	.579 $\pm$ .029	.544 $\pm$ .009	<b>.691 <math>\pm</math> .009</b>
diabetes	.585 $\pm$ .014	.591 $\pm$ .008	.577 $\pm$ .020	.542 $\pm$ .020	.535 $\pm$ .020	.580 $\pm$ .037	.571 $\pm$ .023	.520 $\pm$ .006	<b>.654 <math>\pm</math> .013</b>
german	.577 $\pm$ .007	.590 $\pm$ .005	.571 $\pm$ .006	.558 $\pm$ .018	.526 $\pm$ .012	.578 $\pm$ .017	.607 $\pm$ .022	.534 $\pm$ .024	<b>.707 <math>\pm</math> .005</b>
contraceptive	.766 $\pm$ .010	.599 $\pm$ .013	.597 $\pm$ .008	.800 $\pm$ .013•	.605 $\pm$ .006	.609 $\pm$ .033	.613 $\pm$ .020	.545 $\pm$ .016	.789 $\pm$ .017
splice	.646 $\pm$ .002	.616 $\pm$ .001	.612 $\pm$ .005	.688 $\pm$ .008	.598 $\pm$ .003	.500 $\pm$ .030	.576 $\pm$ .035	.554 $\pm$ .008	<b>.772 <math>\pm</math> .000</b>
mushroom	.679 $\pm$ .010	.667 $\pm$ .037	.630 $\pm$ .003	.793 $\pm$ .017	.619 $\pm$ .004	.860 $\pm$ .026•	.763 $\pm$ .011	.674 $\pm$ .018	.799 $\pm$ .007
marketing	.652 $\pm$ .003	.628 $\pm$ .011	.635 $\pm$ .007	.766 $\pm$ .009	.597 $\pm$ .006	.717 $\pm$ .038	.696 $\pm$ .006	.565 $\pm$ .010	<b>.818 <math>\pm</math> .000</b>
hapt	.591 $\pm$ .001	.697 $\pm$ .010	.617 $\pm$ .003	.654 $\pm$ .005	.547 $\pm$ .005	.575 $\pm$ .019	.596 $\pm$ .005	.555 $\pm$ .004	<b>.764 <math>\pm</math> .010</b>
ring	.596 $\pm$ .002	.597 $\pm$ .002	.581 $\pm$ .005	.525 $\pm$ .006	.529 $\pm$ .008	.599 $\pm$ .012	.566 $\pm$ .016	.525 $\pm$ .006	<b>.705 <math>\pm</math> .002</b>
magic04	.613 $\pm$ .001	.607 $\pm$ .003	.608 $\pm$ .001	.719 $\pm$ .005	.599 $\pm$ .002	.777 $\pm$ .071	.837 $\pm$ .013•	.553 $\pm$ .002	.810 $\pm$ .008
a8a	.599 $\pm$ .000	.678 $\pm$ .008	.596 $\pm$ .005	.672 $\pm$ .004	.516 $\pm$ .005	.664 $\pm$ .006	.617 $\pm$ .005	.543 $\pm$ .005	.684 $\pm$ .016
loss/win	0/14	0/14	0/14	2/12	0/14	1/13	1/13	0/14	<b>4/108*</b>
p-value	.0001	.0001	.0001	.0009	.0001	.0003	.0001	.0001	—
F-rank	5.3571	5.5357	4.2143	6.2143	2.500	5.2143	5.1429	2.1071	8.7143

Table 2: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR  $\pm$  Standard Variance) for 10 data sets in the case of ASymmetric Noisy labeled data. • indicates the cases that our method loses the comparison. OOT: Out Of Time (24 hours).

	ASymmetric Noise $\rho_{-1} = 0.3, \rho_{+1} = 0.5$								
Datasets	FOBOS	RDA	TG	FTRL	FTRL-P	OSLMF	OVFM	OLI <sup>2</sup> DS	ONLSF
wdbc	.651 $\pm$ .010	.652 $\pm$ .017	.640 $\pm$ .015	.567 $\pm$ .024	.584 $\pm$ .013	.525 $\pm$ .059	.632 $\pm$ .038	.547 $\pm$ .031	<b>.676 <math>\pm</math> .006</b>
breast	.795 $\pm$ .029•	.665 $\pm$ .009	.804 $\pm$ .010•	.795 $\pm$ .018•	.672 $\pm$ .018	.591 $\pm$ .037	.656 $\pm$ .013	.533 $\pm$ .033	.778 $\pm$ .013
dermatology	.643 $\pm$ .013	.639 $\pm$ .013	.639 $\pm$ .018	.746 $\pm$ .026•	.666 $\pm$ .043	.610 $\pm$ .027	.673 $\pm$ .023	.657 $\pm$ .024	.686 $\pm$ .006
wdbc	.620 $\pm$ .017	.608 $\pm$ .011	.609 $\pm$ .008	.606 $\pm$ .029	.599 $\pm$ .018	.588 $\pm$ .047	.616 $\pm$ .016	.556 $\pm$ .011	<b>.709 <math>\pm</math> .006</b>
diabetes	.548 $\pm$ .004	.547 $\pm$ .011	.545 $\pm$ .007	.498 $\pm$ .011	.492 $\pm$ .010	.510 $\pm$ .035	.492 $\pm$ .016	.510 $\pm$ .015	<b>.563 <math>\pm</math> .026</b>
german	.517 $\pm$ .007	.543 $\pm$ .003	.532 $\pm$ .012	.481 $\pm$ .021	.491 $\pm$ .025	.476 $\pm$ .033	.487 $\pm$ .039	.536 $\pm$ .019	<b>.557 <math>\pm</math> .062</b>
contraceptive	.750 $\pm$ .022•	.616 $\pm$ .012	.619 $\pm$ .012	.759 $\pm$ .016•	.618 $\pm$ .005	.622 $\pm$ .047	.645 $\pm$ .038	.569 $\pm$ .017	.728 $\pm$ .039
splice	.742 $\pm$ .009	.676 $\pm$ .002	.668 $\pm$ .005	.762 $\pm$ .003•	.651 $\pm$ .001	.607 $\pm$ .025	.639 $\pm$ .021	.547 $\pm$ .010	.759 $\pm$ .000
mushroom	.716 $\pm$ .007	.656 $\pm$ .015	.651 $\pm$ .005	.732 $\pm$ .013•	.637 $\pm$ .006	.736 $\pm$ .071•	.743 $\pm$ .006•	.679 $\pm$ .016	.725 $\pm$ .007
marketing	.795 $\pm$ .003	.758 $\pm$ .006	.738 $\pm$ .003	.820 $\pm$ .003	.680 $\pm$ .006	.767 $\pm$ .017	.783 $\pm$ .006	.574 $\pm$ .017	.818 $\pm$ .000
hapt	.662 $\pm$ .001	.774 $\pm$ .007	.690 $\pm$ .004	.741 $\pm$ .005	.568 $\pm$ .012	.653 $\pm$ .009	.653 $\pm$ .017	.561 $\pm$ .005	<b>.820 <math>\pm</math> .002</b>
ring	.592 $\pm$ .003	.592 $\pm$ .001	.572 $\pm$ .004	.503 $\pm$ .006	.501 $\pm$ .006	.502 $\pm$ .012	.514 $\pm$ .006	.528 $\pm$ .003	<b>.607 <math>\pm</math> .020</b>
magic04	.643 $\pm$ .002	.639 $\pm$ .003	.637 $\pm$ .002	.733 $\pm$ .002	.630 $\pm$ .003	.631 $\pm$ .066	.796 $\pm$ .007•	.566 $\pm$ .004	.787 $\pm$ .017
a8a	.649 $\pm$ .000	.685 $\pm$ .013	.621 $\pm$ .006	.740 $\pm$ .013•	.520 $\pm$ .007	.740 $\pm$ .005•	.684 $\pm$ .005	.547 $\pm$ .004	.736 $\pm$ .001
loss/win	2/12	0/14	1/13	6/8	0/14	2/12	2/12	0/14	<b>13/99*</b>
p-value	.0034	.0001	.0003	.0765	.0001	.0003	.0004	.0001	—
F-rank	6.3571	5.4286	4.8929	6.2857	2.8214	3.4643	5.1429	2.6071	8.0000

Table 3: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR  $\pm$  Standard Variance) for 8 data sets in the case of Flip Noisy labeled data. • indicates the cases that our method loses the comparison. OOT: Out Of Time (24 hours)

	Filp Noise Rate = 0.4								
Datasets	FOBOS	RDA	TG	FTRL	FTRL-P	OSLMF	OVFM	OLI2DS	ONLSF
wdbc	.608 $\pm$ .011	.605 $\pm$ .007	.584 $\pm$ .031	.508 $\pm$ .014	.503 $\pm$ .039	.571 $\pm$ .079	.528 $\pm$ .046	.509 $\pm$ .026	<b>.626 <math>\pm</math> .012</b>
breast	.710 $\pm$ .017	.630 $\pm$ .027	.744 $\pm$ .051	.788 $\pm$ .042	.660 $\pm$ .027	.582 $\pm$ .070	.656 $\pm$ .013	.561 $\pm$ .056	<b>.917 <math>\pm</math> .034</b>
dermatology	.598 $\pm$ .017	.596 $\pm$ .012	.584 $\pm$ .018	.697 $\pm$ .027•	.616 $\pm$ .051	.534 $\pm$ .053	.601 $\pm$ .013	.640 $\pm$ .032	.658 $\pm$ .015
wdbc	.589 $\pm$ .005	.583 $\pm$ .008	.577 $\pm$ .007	.572 $\pm$ .019	.583 $\pm$ .027	.574 $\pm$ .018	.576 $\pm$ .020	.543 $\pm$ .012	<b>.686 <math>\pm</math> .005</b>
diabetes	.596 $\pm$ .007	.592 $\pm$ .008	.577 $\pm$ .008	.530 $\pm$ .012	.516 $\pm$ .012	.573 $\pm$ .053	.577 $\pm$ .012	.523 $\pm$ .004	<b>.652 <math>\pm</math> .005</b>
german	.583 $\pm$ .014	.593 $\pm$ .012	.579 $\pm$ .009	.555 $\pm$ .011	.540 $\pm$ .020	.557 $\pm$ .039	.570 $\pm$ .028	.521 $\pm$ .016	<b>.690 <math>\pm</math> .010</b>
contraceptive	.761 $\pm$ .017	.604 $\pm$ .009	.600 $\pm$ .007	.794 $\pm$ .020•	.606 $\pm$ .004	.560 $\pm$ .044	.637 $\pm$ .028	.562 $\pm$ .020	.794 $\pm$ .023
splice	.642 $\pm$ .004	.615 $\pm$ .004	.611 $\pm$ .005	.692 $\pm$ .009•	.597 $\pm$ .005	.526 $\pm$ .031	.573 $\pm$ .027	.551 $\pm$ .009	<b>.772 <math>\pm</math> .000</b>
mushroom	.683 $\pm$ .005	.682 $\pm$ .015	.630 $\pm$ .010	.799 $\pm$ .008	.623 $\pm$ .006	.906 $\pm$ .009•	.765 $\pm$ .014	.677 $\pm$ .015	.811 $\pm$ .008
marketing	.650 $\pm$ .009	.646 $\pm$ .022	.638 $\pm$ .002	.771 $\pm$ .007	.592 $\pm$ .008	.736 $\pm$ .019	.694 $\pm$ .012	.567 $\pm$ .008	<b>.817 <math>\pm</math> .000</b>
hapt	.592 $\pm$ .002	.709 $\pm$ .017	.615 $\pm$ .004	.654 $\pm$ .003	.569 $\pm$ .009	.568 $\pm$ .011	.587 $\pm$ .012	.552 $\pm$ .003	<b>.764 <math>\pm</math> .011</b>
ring	.597 $\pm$ .002	.598 $\pm$ .002	.584 $\pm$ .004	.532 $\pm$ .002	.530 $\pm$ .007	.605 $\pm$ .008	.562 $\pm$ .005	.526 $\pm$ .005	<b>.705 <math>\pm</math> .002</b>
magic04	.613 $\pm$ .002	.609 $\pm$ .003	.610 $\pm$ .003	.719 $\pm$ .003	.602 $\pm$ .002	.745 $\pm$ .039	.830 $\pm$ .007•	.554 $\pm$ .003	.809 $\pm$ .013
a8a	.599 $\pm$ .000	.678 $\pm$ .017•	.598 $\pm$ .003	.681 $\pm$ .008•	.523 $\pm$ .009	.676 $\pm$ .008•	.614 $\pm$ .004	.544 $\pm$ .003	.663 $\pm$ .007
loss/win	0/14	1/13	0/14	4/10	0/14	2/12	1/13	0/14	<b>8/104*</b>
p-value	.0001	.0001	.0001	.0026	.0001	.0012	.0001	.0001	—
F-rank	6.0714	5.6071	4.4643	5.9643	2.9643	4.3571	5.0357	2.0000	8.5357

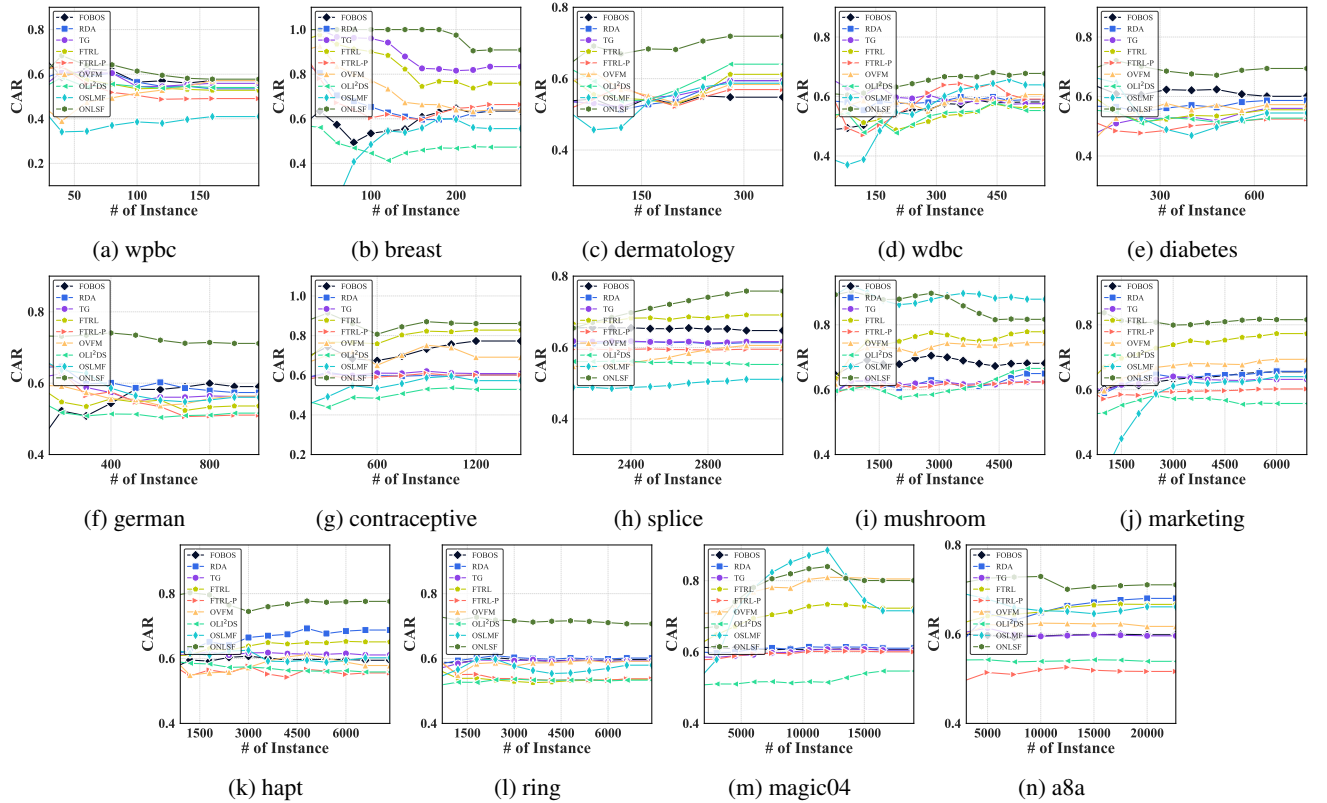


Figure 1: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM, OL<sup>2</sup>DS and ONLSF in all 14 symmetric noisy labeled data streams.

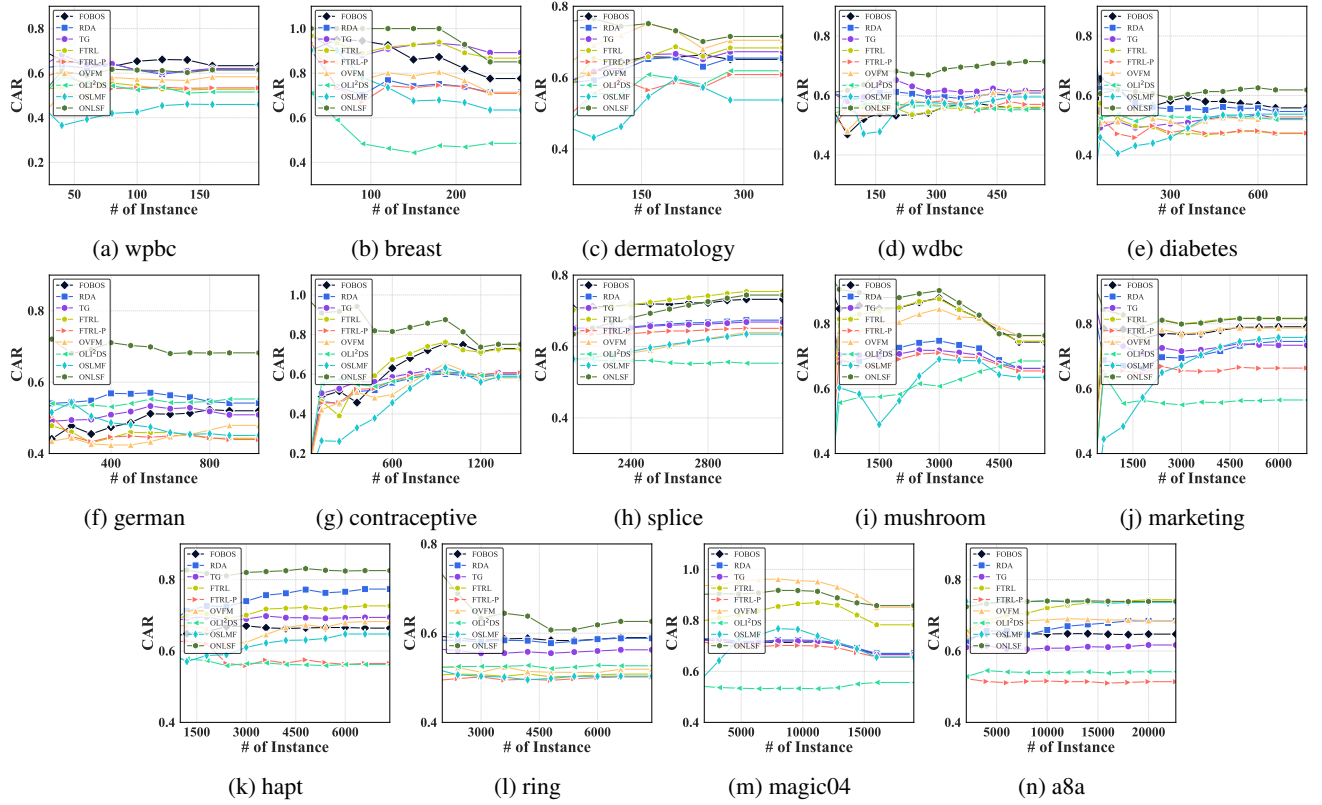


Figure 2: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM, OL<sup>2</sup>DS and ONLSF in all 14 asymmetric noisy labeled data streams.

Table 4: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR  $\pm$  Standard Variance) for 10 data sets in the case of Symmetric Noisy labeled data. • indicates the cases that our method loses the comparison. OOT: Out Of Time (24 hours)

	Symetric Noise $\rho_{-1} = \rho_{+1} = 0.4$						
Datasets	Co-teach+	Phuber	Bey.Ima.	OPA	SCW	AdaStream	ONLSF
wpbc	.448 $\pm$ .259	.658 $\pm$ .210•	.551 $\pm$ .256	.606 $\pm$ .000	.606 $\pm$ .000	.585 $\pm$ .146	.642 $\pm$ .012
breast	.439 $\pm$ .119	.768 $\pm$ .077	.545 $\pm$ .149	.603 $\pm$ .000	.628 $\pm$ .020	.518 $\pm$ .136	<b>.895 <math>\pm</math> .020</b>
dermatology	.460 $\pm$ .186	.616 $\pm$ .153	.464 $\pm$ .185	.603 $\pm$ .000	.606 $\pm$ .002	.499 $\pm$ .171	<b>.665 <math>\pm</math> .012</b>
wdbc	.473 $\pm$ .126	.475 $\pm$ .125	.424 $\pm$ .102	.603 $\pm$ .000	.610 $\pm$ .006	.530 $\pm$ .080	<b>.691 <math>\pm</math> .009</b>
diabetes	.469 $\pm$ .145	.470 $\pm$ .148	.430 $\pm$ .101	.600 $\pm$ .000	.601 $\pm$ .008	.480 $\pm$ .066	<b>.654 <math>\pm</math> .013</b>
german	.458 $\pm$ .183	.540 $\pm$ .196	.556 $\pm$ .155	.600 $\pm$ .000	.599 $\pm$ .006	.377 $\pm$ .094	<b>.707 <math>\pm</math> .005</b>
contraceptive	.456 $\pm$ .059	.427 $\pm$ .000	.438 $\pm$ .077	.601 $\pm$ .000	.534 $\pm$ .039	.507 $\pm$ .021	<b>.789 <math>\pm</math> .017</b>
splice	.666 $\pm$ .062	.656 $\pm$ .207	.630 $\pm$ .194	.600 $\pm$ .000	.601 $\pm$ .000	.548 $\pm$ .071	<b>.772 <math>\pm</math> .000</b>
mushroom	.523 $\pm$ .115	.817 $\pm$ .024•	.523 $\pm$ .144	.600 $\pm$ .000	.802 $\pm$ .081•	.533 $\pm$ .079	.799 $\pm$ .007
marketing	.691 $\pm$ .254	.796 $\pm$ .017	.683 $\pm$ .247	.600 $\pm$ .000	.795 $\pm$ .010	.517 $\pm$ .121	<b>.818 <math>\pm</math> .000</b>
hapt	.684 $\pm$ .047	.537 $\pm$ .201	.461 $\pm$ .208	.600 $\pm$ .000	.601 $\pm$ .000	.649 $\pm$ .058	<b>.764 <math>\pm</math> .010</b>
ring	.476 $\pm$ .063	.567 $\pm$ .089	.548 $\pm$ .083	.600 $\pm$ .000	.619 $\pm$ .014	.490 $\pm$ .052	<b>.705 <math>\pm</math> .002</b>
magic04	.435 $\pm$ .129	.529 $\pm$ .146	.374 $\pm$ .203	.600 $\pm$ .000	.627 $\pm$ .015	.642 $\pm$ .008	<b>.810 <math>\pm</math> .008</b>
a8a	.657 $\pm$ .064	.697 $\pm$ .009•	.534 $\pm$ .062	.600 $\pm$ .000	.724 $\pm$ .010•	.687 $\pm$ .004•	.684 $\pm$ .016
loss/win	14/0	11/3	14/0	14/0	12/2	13/1	6/78
p-value	.0001	.0009	.0001	.0001	.0006	.0001	—
F-rank	2.5357	4.4286	2.1071	4.0357	5.1786	3.1429	6.5714

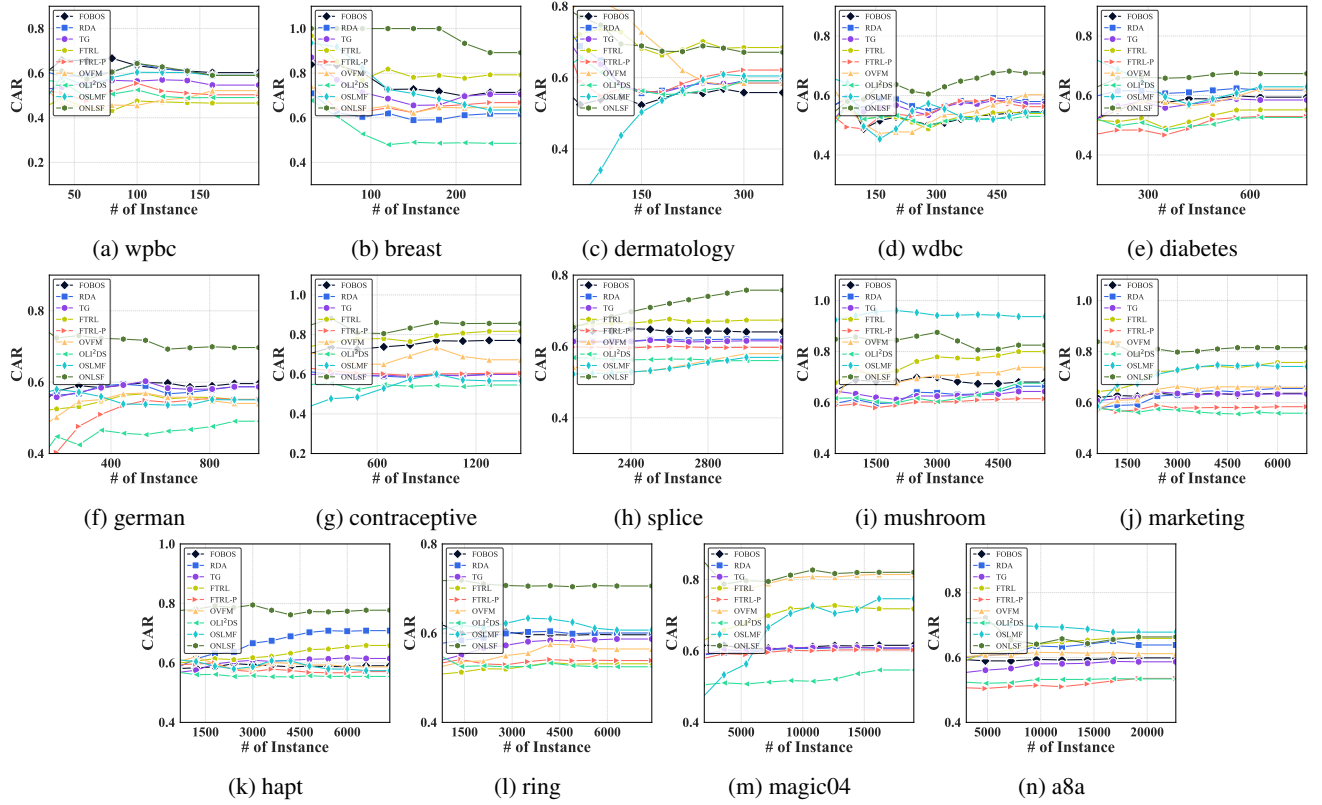


Figure 3: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM, OL<sup>2</sup>DS and ONLSF in all 14 flipped noisy labeled data streams.

Table 5: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR  $\pm$  Standard Variance) for 10 data sets in the case of Symmetric Noisy labeled data. • indicates the cases that our method loses the comparison. OOT: Out Of Time (24 hours)

	ASymmetric Noise $\rho_{-1} = 0.3, \rho_{+1} = 0.5$						
Datasets	Co-teach+	Phuber	Bey.Ima.	OPA	SCW	AdaStream	ONLSF
wpbc	.448 $\pm$ .257	.658 $\pm$ .210	.655 $\pm$ .209	.657 $\pm$ .000	.657 $\pm$ .000	.579 $\pm$ .151	<b>.676 <math>\pm</math> .006</b>
breast	.479 $\pm$ .154	.705 $\pm$ .004	.535 $\pm$ .186	.646 $\pm$ .000	.642 $\pm$ .021	.657 $\pm$ .098	<b>.778 <math>\pm</math> .013</b>
dermatology	.606 $\pm$ .173	.616 $\pm$ .153	.542 $\pm$ .188	.640 $\pm$ .000	.640 $\pm$ .001	.614 $\pm$ .152	<b>.686 <math>\pm</math> .006</b>
wdbc	.473 $\pm$ .126	.373 $\pm$ .000	.536 $\pm$ .134	.626 $\pm$ .000	.638 $\pm$ .004	.568 $\pm$ .073	<b>.709 <math>\pm</math> .006</b>
diabetes	.406 $\pm$ .123	.470 $\pm$ .148	.576 $\pm$ .116•	.570 $\pm$ .000•	.555 $\pm$ .015	.349 $\pm$ .000	.563 $\pm$ .026
german	.619 $\pm$ .158•	.540 $\pm$ .196	.524 $\pm$ .175	.560 $\pm$ .000•	.529 $\pm$ .015	.417 $\pm$ .096	.557 $\pm$ .062
contraceptive	.514 $\pm$ .073	.485 $\pm$ .072	.507 $\pm$ .067	.615 $\pm$ .000	.652 $\pm$ .038	.513 $\pm$ .012	<b>.728 <math>\pm</math> .039</b>
splice	.599 $\pm$ .189	.759 $\pm$ .000	.555 $\pm$ .242	.652 $\pm$ .000	.653 $\pm$ .000	.480 $\pm$ .092	.759 $\pm$ .000
mushroom	.521 $\pm$ .104	.670 $\pm$ .059•	.469 $\pm$ .085	.624 $\pm$ .000	.665 $\pm$ .082	.513 $\pm$ .068	<b>.725 <math>\pm</math> .007</b>
marketing	.623 $\pm$ .239	.816 $\pm$ .002	.574 $\pm$ .296	.664 $\pm$ .000	.805 $\pm$ .005	.488 $\pm$ .085	<b>.818 <math>\pm</math> .000</b>
hapt	.674 $\pm$ .212	.735 $\pm$ .036	.542 $\pm$ .209	.667 $\pm$ .000	.668 $\pm$ .000	.762 $\pm$ .028	<b>.820 <math>\pm</math> .002</b>
ring	.482 $\pm$ .058	.522 $\pm$ .109	.498 $\pm$ .058	.599 $\pm$ .000	.566 $\pm$ .020	.462 $\pm$ .021	<b>.607 <math>\pm</math> .020</b>
magic04	.426 $\pm$ .227	.679 $\pm$ .045	.410 $\pm$ .119	.630 $\pm$ .000	.636 $\pm$ .001	.642 $\pm$ .007	<b>.787 <math>\pm</math> .017</b>
a8a	.682 $\pm$ .045	.712 $\pm$ .008	.492 $\pm$ .055	.651 $\pm$ .010	.733 $\pm$ .003	.716 $\pm$ .005	<b>.736 <math>\pm</math> .001</b>
loss/win	13/1	13/1	13/1	12/2	14/0	14/0	5/79
p-value	.0002	.0008	.0001	.0003	.0001	.0001	—
F-rank	2.7857	4.4643	2.2143	4.3571	4.6429	2.8571	6.6786

Table 6: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR  $\pm$  Standard Variance) for 10 data sets in the case of Symmetric Noisy labeled data. • indicates the cases that our method loses the comparison. OOT: Out Of Time (24 hours)

	Flip Noise rate = 0.4						
Datasets	Co-teach+	Phuber	Bey.Ima.	OPA	SCW	AdaStream	ONLSF
wdbc	.449 $\pm$ .233	.553 $\pm$ .257	.387 $\pm$ .207	.601 $\pm$ .000	.601 $\pm$ .000	.645 $\pm$ .146•	.626 $\pm$ .012
breast	.494 $\pm$ .110	.695 $\pm$ .074	.550 $\pm$ .193	.603 $\pm$ .000	.609 $\pm$ .025	.508 $\pm$ .127	<b>.917 <math>\pm</math> .034</b>
dermatology	.524 $\pm$ .179	.463 $\pm$ .188	.576 $\pm$ .149	.601 $\pm$ .000	.599 $\pm$ .001	.538 $\pm$ .186	<b>.658 <math>\pm</math> .015</b>
wdbc	.576 $\pm$ .102	.525 $\pm$ .125	.475 $\pm$ .126	.601 $\pm$ .000	.608 $\pm$ .009	.576 $\pm$ .063	<b>.686 <math>\pm</math> .005</b>
diabetes	.410 $\pm$ .122	.470 $\pm$ .148	.501 $\pm$ .129	.600 $\pm$ .000	.609 $\pm$ .023	.406 $\pm$ .071	<b>.652 <math>\pm</math> .005</b>
german	.530 $\pm$ .189	.380 $\pm$ .160	.561 $\pm$ .169	.600 $\pm$ .000	.588 $\pm$ .009	.335 $\pm$ .070	<b>.690 <math>\pm</math> .010</b>
contraceptive	.458 $\pm$ .105	.456 $\pm$ .058	.485 $\pm$ .072	.600 $\pm$ .000	.589 $\pm$ .046	.500 $\pm$ .023	<b>.794 <math>\pm</math> .023</b>
splice	.657 $\pm$ .084	.344 $\pm$ .207	.641 $\pm$ .200	.600 $\pm$ .000	.601 $\pm$ .000	.488 $\pm$ .086	<b>.772 <math>\pm</math> .000</b>
mushroom	.458 $\pm$ .097	.668 $\pm$ .178•	.526 $\pm$ .112	.600 $\pm$ .000	.790 $\pm$ .070	.499 $\pm$ .060	<b>.811 <math>\pm</math> .008</b>
marketing	.331 $\pm$ .246	.805 $\pm$ .012	.564 $\pm$ .308	.600 $\pm$ .000	.788 $\pm$ .028	.447 $\pm$ .105	<b>.817 <math>\pm</math> .000</b>
hapt	.548 $\pm$ .204	.626 $\pm$ .162	.445 $\pm$ .185	.600 $\pm$ .000	.601 $\pm$ .000	.713 $\pm$ .044	<b>.764 <math>\pm</math> .011</b>
ring	.485 $\pm$ .029	.389 $\pm$ .000	.498 $\pm$ .095	.600 $\pm$ .000	.632 $\pm$ .014	.508 $\pm$ .042	<b>.705 <math>\pm</math> .002</b>
magic04	.580 $\pm$ .109	.689 $\pm$ .053	.565 $\pm$ .140	.600 $\pm$ .000	.633 $\pm$ .004	.642 $\pm$ .008	<b>.809 <math>\pm</math> .013</b>
a8a	.567 $\pm$ .140	.685 $\pm$ .013	.532 $\pm$ .068	.600 $\pm$ .000	.728 $\pm$ .012•	.684 $\pm$ .004•	.663 $\pm$ .007
loss/win	14/0	13/1	14/0	14/0	13/1	12/2	4/80
p-value	.0001	.0002	.0001	.0001	.0009	.0003	—
F-rank	2.2500	3.4286	2.6429	4.3929	5.1786	3.3929	6.7143

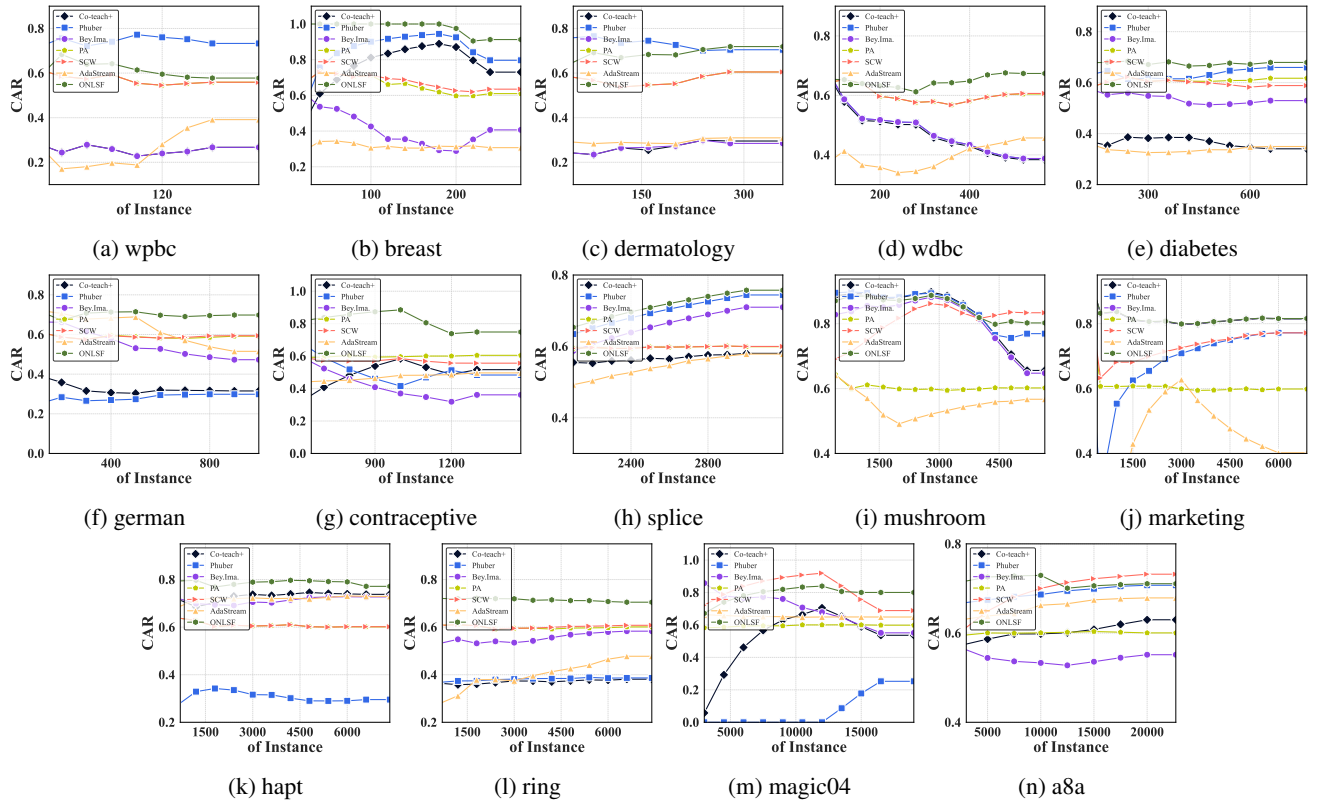


Figure 4: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM, OL<sup>2</sup>DS and ONLSF in all 14 symmetric noisy labeled data streams.

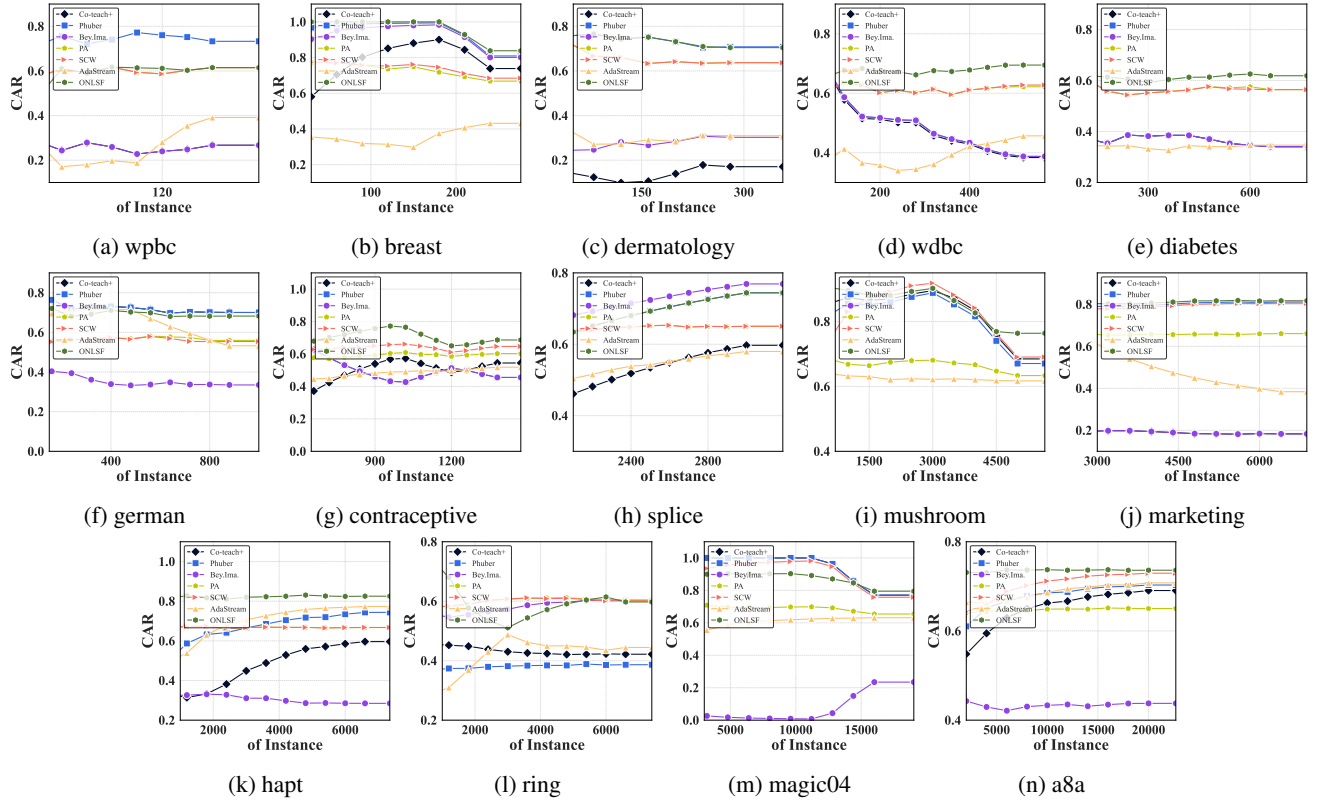


Figure 5: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM,  $OL^2DS$  and ONLSF in all 14 asymmetric noisy labeled data streams.



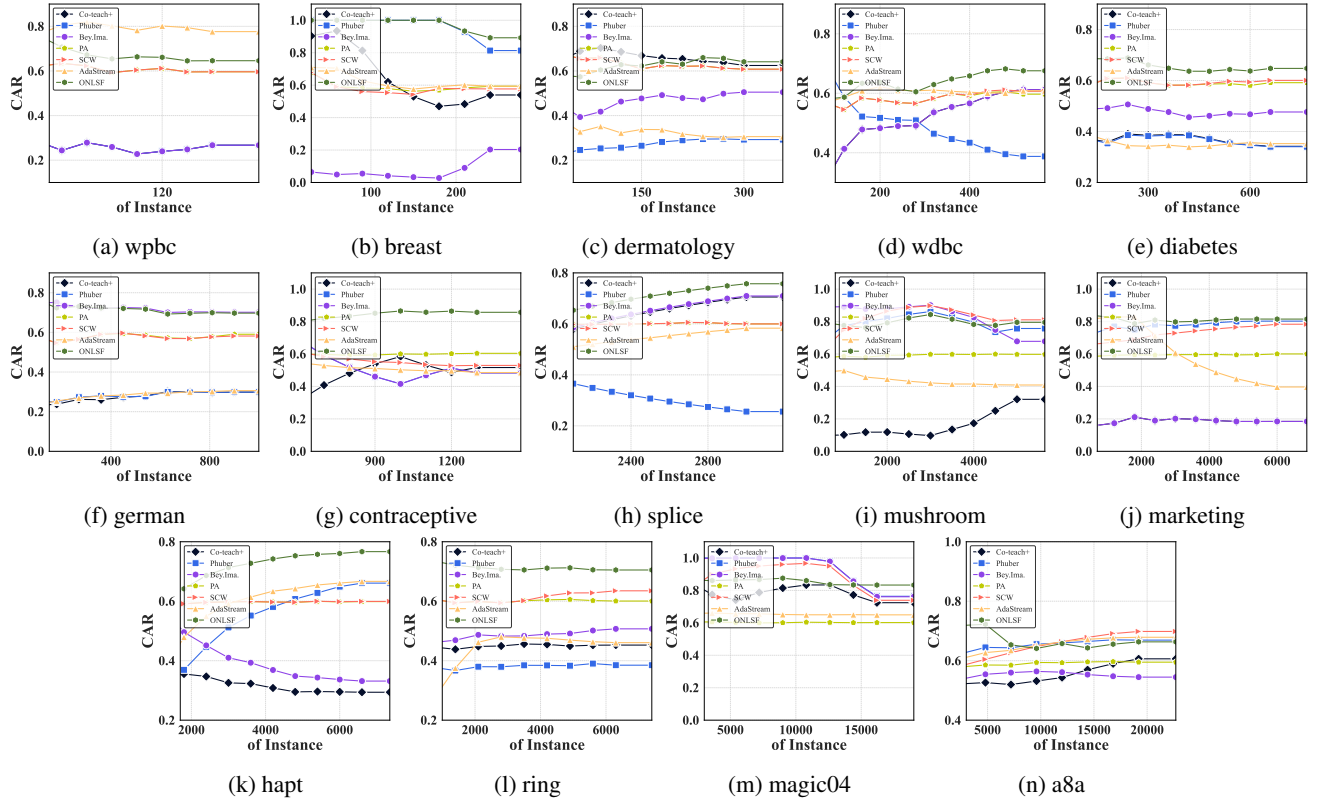


Figure 6: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM,  $OL^2DS$  and ONLSF in all 14 flipped noisy labeled data streams.