



Learning to learn from noisy web videos Serena Yeung¹, Vignesh Ramanathan¹, Olga Russakovsky², Liyue Shen¹, Greg Mori³, Li Fei-Fei¹ ¹Stanford University ²Carnegie Mellon University ³Simon Fraser University

Introduction

- Manually labeling training videos for action recognition is impractical to scale to the long-tail of action categories: e.g. fine-gained, rare, or niche classes.
- We can leverage noisy data from web queries to learn new actions, using semi-supervised or "webly-supervised" approaches. However, existing methods typically do not learn and leverage domain-specific knowledge, or rely on iterative hand-tuned data labeling policies.
- Our insight is that good labeling policies can be learned from existing annotated datasets. A good policy should label noisy data such that a classifier trained on the labels would achieve high classification accuracy on the existing datasets.
- We propose a reinforcement learning-based formulation for learning data labeling policies from noisy web data. Concretely, we introduce a joint formulation of a Q-learning agent and a class recognition model. The agent selects web search examples to label as positives, which are then used to train the recognition model.

Sports-1M action recognition

- Training classes (used to learn policy): 300 Sports-1M classes
- Test classes: 105 Sports-1M classes
- Policy labels noisy YouTube data, using videos returned by the YouTube query suggestion feature for 30 different query expansions per class.
- At training time of learning the policy, rewards are based on classification accuracy, where classifiers are trained on the policy-labeled noisy data and evaluated on the annotated reward dataset (Sports-1M test videos for the 300 classes).
- To evaluate the learned policy, classifiers are trained on policy-labeled noisy data for the 105 previously unseen Sports-1M test classes, and evaluated on annotated Sports-1M test videos for these classes.

Method	Budget-60	Budget-80	Budget-100
Seed	64.3	64.3	64.3
Label propagation	65.4	65.4	67.2
Label spreading	65.4	66.6	67.3
TSVM	70.7	71.7	72.5
Greedy	71.7	73.8	74.8
Greedy-clustering	72.3	73.2	74.3
Greedy-KL	74.1	74.7	74.7
Ours	75.4	76.2	77.0

mAP on Sports-1M with different budgets for the number of selected positive examples.





 $s = \{H_{pos}, H_{pos}, \{H_{D_1}, \dots, H_{D_K}\}, P\}, \text{ where } \{H_{pos}, M_{D_1}, \dots, H_{D_K}\}, P\}$ classifier scores for the positive set, the negative





method is robust to semantic drift and selects useful subcategories of bobsleigh videos such as crashes and pov.





Noisy MNIST digit classification

7777771111211111111 7777 11212 222 21 21 21 21 21 8678996989888899999 67787878688888999999

Ten sample query subsets in Noisy MNIST for the digit 7. *Top row.* Different translation and rotation transformations. *Bottom row.* The two leftmost queries have different amounts of noise, the center one is a

Label propagation	Label spreading	TSVM	Iterative classifier	Ours
37.9	41.1	39.5	43.1	60.9
40.8	45.6	44.4	43.2	61.3
42.2	46.7	46.2	42.4	71.4
51.1	48.6	46.1	49.7	55.1
48.8	48.5	42.6	48.5	57.6
48.1	46.6	39.7	47.4	55.7
35.0	35.2	41.2	38.3	56.2
40.0	34.1	39.6	39.6	55.6
42.0	30.2	40.8	38.0	55.5
37.5	36.5	41.4	52.4	52.4
37.9	37.4	38.9	53.5	53.5
38.0	37.6	39.5	55.7	55.7
40.4	40.3	42.1	43.4	56.1
41.9	41.4	41.4	43.4	57.0
42.6	40.3	41.5	42.3	59.5

AP on Noisy MNIST, with budgets of 60, 80 and 100 for the numbers of positive examples selected from D_{cand} .

		· •	\sim
6 6	86666	77666	66666
66	6666	66666	66666
0	6600	6666	6655
ها و	66000	6666	6666
9 9	97799	77777	77777
99	999999	77717	77777
2 9	8888	99899	8 8 8 8
5	99999	<u>ه</u> ه 9 ه 9	> 9 9 ? 9

Comparison of positive query subsets selected using our method versus the greedy classifier baseline. Subsets chosen by each method are shown from left to right, for the digit 6 (top example) and the digit 9 (bottom example). Our model is better able to select useful positives with visual diversity, while avoiding semantic drift.

Long-tail action labeling

Taking	a selfie
	1 2 3 3 4 5 5 Taking a selfie
	61 62 63 64 65
	Taking a selfie every day
	96 97 98 98 99 99 99 99 99 99 99 99 90 90 90 90 90
Olympic gy	mastics
	1 2 3 3 4 5 Olympic gymnastics
5	
	Olympic gymnastics hoop
	Olympic acrobatic gymnastics
A	86 87 88 OURS 89 90 Olympic gymnastics dance

OURS

Comparison of positive query subsets selected using our method versus the greedy classifier baseline, for two long-tail classes. See Sports-1M figure for explanation of figure structure. *Top example*. The greedy classifier selects many similar-looking examples of taking a selfie, while our method learns domain-specific knowledge that positives in different environments are more useful, e.g. with a tornado or underwater. Bottom example. The greedy classifier selects similar examples of gymnastics,

References

[5] L.-J. Li and L. Fei-Fei. OPTIMOL: automatic online picture collection via incremental model learning. International Journal of Computer Vision, 88(2):147–168, 2010.