

Query Learning for Fish Identification

Gaoang Wang¹, Jenq-Neng Hwang¹, Farron Wallace², Craig S. Rose²

¹University of Washington

²National Oceanic and Atmospheric Administration



ELECTRICAL ENGINEERING

UNIVERSITY *of* WASHINGTON



NOAA

Outline

- Introduction and Overview
 - Fish Identification
 - Query Learning
- Query Learning with Diversity Constraint
- Combine Semi-Supervised Learning
- Results

Fish Identification

- Want to know the species of the fish given input images [1].



Fish Identification Results

- Datasets
 - 2015 chute data (8835 images with 27 classes)
 - 2016 chute data (5032 images with 27 classes)
- Same dataset split into training and testing

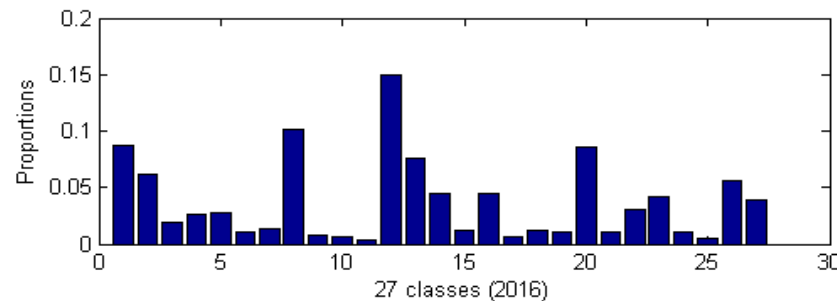
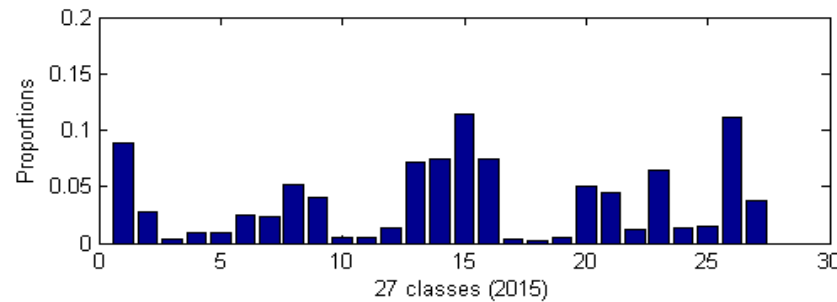
Training Data	Testing Data	Cross Validation	Accuracy (%)
2015	2015	10-fold	96.1
2016	2016	10-fold	98.5
2015+2016	2015+2016	10-fold	96.9

Some Problems with Supervised Learning

- If large difference exists between training and testing datasets
 - Different camera distortions
 - Different colors
 - Different distributions of species

FishID	Fish Name
1	Arrowtooth Flounder
2	Atka Mackerel
3	Bathymaster Signatus
4	Berryteuthis Magister
5	Blackspotted Rockfish
6	Dover Sole
7	Dusky Rockfish
8	Flathead Sole
9	Giant Grenadier
10	Gorgonocephalus Eucnemis
11	Harlequin Rockfish
12	Northern Rock Sole
13	Northern Rockfish
14	Pacific Cod
15	Pacific Halibut
16	Pacific Ocean Perch
17	Pacific Octopus
18	Paragorgia Arborea
19	Prowfish
20	Rex Sole
21	Sablefish
22	Shortraker Rockfish
23	Shortspine Thornyhead
24	Strongylocentrotus sp
25	Sturgeon Poacher
26	Walleye Pollock
27	Yellow Irish Lord

Species Distributions

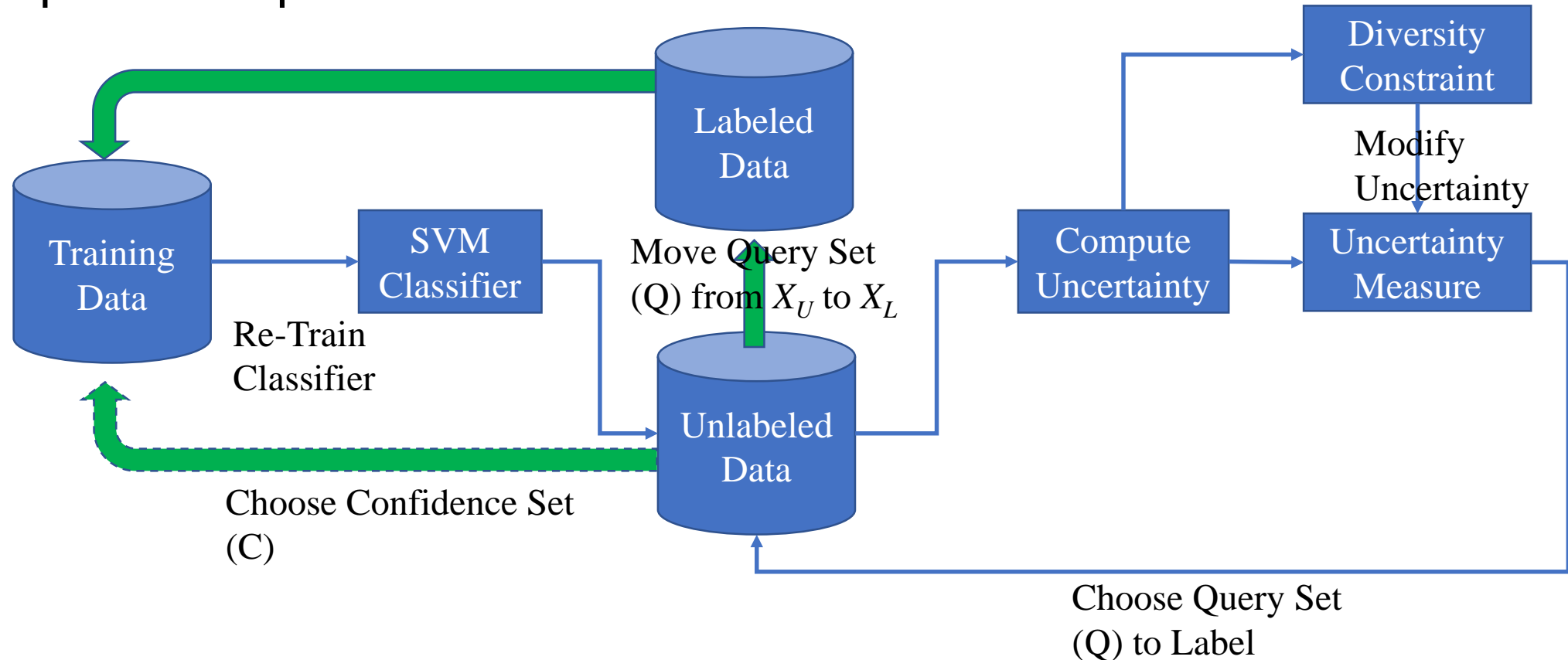


- Comparison with across-year datasets

Training Data	Testing Data	Acc (%)
2015 dataset (5%)	2015 dataset (95%)	83.9
2016 dataset (5%)	2016 dataset (95%)	86.6
2015 dataset (100%)	2016 dataset (100%)	69.5
2015 dataset+2016 dataset (5%)	2016 dataset (95%)	88.1

Query Learning (Active Learning)

- Goal: **iteratively** select **informative** samples for human labeling to improve the performance of the classifier.



Uncertainty Measure

- **Best vs. Second Best (BvSB) Strategy:**

- Predictions from multi-class SVM classifier

$$p = [p_1, p_2, \dots, p_K] \text{ for } K \text{ classes.}$$

p_k is the likelihood of the sample x belongs to k -th class.

- Define the uncertainty by

$$f(x) = \max(1 + p_{k_2} - p_{k_1}, 0) \in [0, 1].$$

k_1, k_2 : the first two most likely predicted classes.

x : arbitrary unlabeled sample.

f : uncertainty score

Query Learning with Diversity Constraint

- Motivation: choose large uncertain samples with **large diversity**.

Fast greedy search:

1. **Choose the top 1 most uncertain** sample x_k to the query set.

$$k^t = \arg \max_i f_i^{t-1}$$

2. **Decrease the uncertainty** of all samples x_i using a Gaussian kernel, like

$$f_i^t = f_i^{t-1} - f_k^{t-1} \exp\left(-\frac{\|x_i - x_k\|^2}{\sigma^2}\right), \forall i$$

where f_i and f_k are the uncertainty scores of x_i and x_k .

3. Go to step 1 until we select enough samples.

x : unlabeled samples.
 f : uncertainty score.
Red: large uncertainty sample.
Green: small uncertainty sample.
Black: selected sample.



Unlabeled samples

Query Learning with Diversity Constraint

- Motivation: choose large uncertain samples with **large diversity**.

Fast greedy search:

1. **Choose the top 1 most uncertain** sample x_k to the query set.

$$k^t = \arg \max_i f_i^{t-1}$$

2. **Decrease the uncertainty** of all samples x_i using a Gaussian kernel, like

$$f_i^t = f_i^{t-1} - f_k^{t-1} \exp\left(-\frac{\|x_i - x_k\|^2}{\sigma^2}\right), \forall i$$

where f_i and f_k are the uncertainty scores of x_i and x_k .

3. Go to step 1 until we select enough samples.

x : unlabeled samples.
 f : uncertainty score.
Red: large uncertainty sample.
Green: small uncertainty sample.
Black: selected sample.



Unlabeled samples

Query Learning with Diversity Constraint

- Motivation: choose large uncertain samples with **large diversity**.

Fast greedy search:

1. **Choose the top 1 most uncertain** sample x_k to the query set.

$$k^t = \arg \max_i f_i^{t-1}$$

2. **Decrease the uncertainty** of all samples x_i using a Gaussian kernel, like

$$f_i^t = f_i^{t-1} - f_k^{t-1} \exp\left(-\frac{\|x_i - x_k\|^2}{\sigma^2}\right), \forall i$$

where f_i and f_k are the uncertainty scores of x_i and x_k .

3. Go to step 1 until we select enough samples.

x : unlabeled samples.
 f : uncertainty score.
Red: large uncertainty sample.
Green: small uncertainty sample.
Black: selected sample.



Unlabeled samples

Query Learning with Diversity Constraint

- Motivation: choose large uncertain samples with **large diversity**.

Fast greedy search:

1. **Choose the top 1 most uncertain** sample x_k to the query set.

$$k^t = \arg \max_i f_i^{t-1}$$

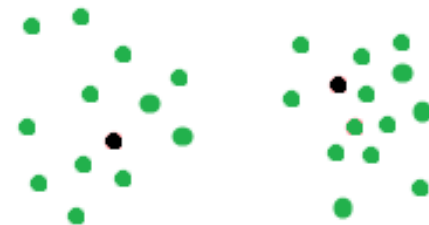
2. **Decrease the uncertainty** of all samples x_i using a Gaussian kernel, like

$$f_i^t = f_i^{t-1} - f_k^{t-1} \exp\left(-\frac{\|x_i - x_k\|^2}{\sigma^2}\right), \forall i$$

where f_i and f_k are the uncertainty scores of x_i and x_k .

3. Go to step 1 until we select enough samples.

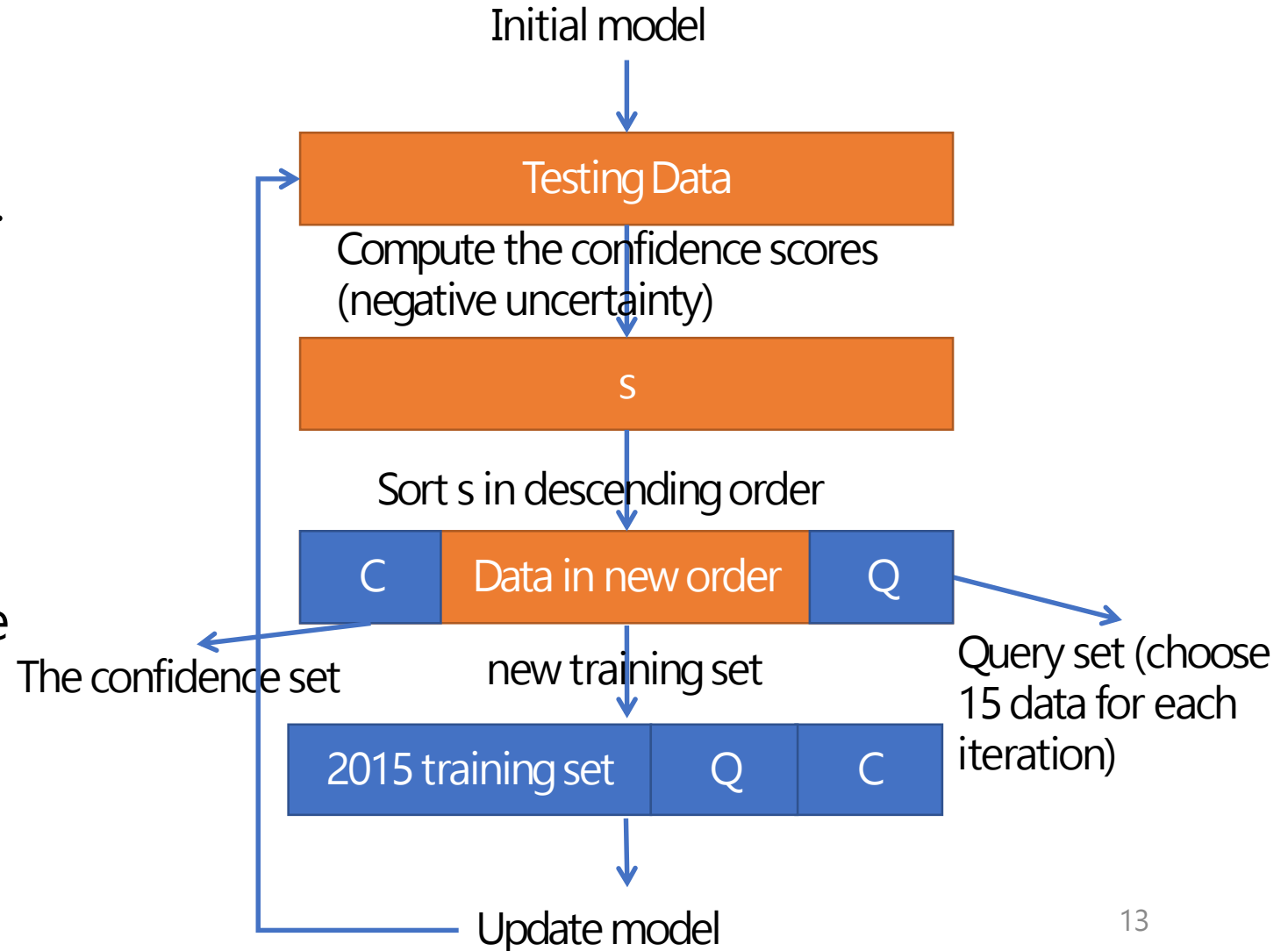
x : unlabeled samples.
 f : uncertainty score.
Red: large uncertainty sample.
Green: small uncertainty sample.
Black: selected sample.



Unlabeled samples

Combine Semi-Supervised Learning

- Train the classifier based on 2015 dataset.
- Test the classifier on 2016 dataset. For each testing data, define a confidence score based on the uncertainty measure.
- Select 15 samples using query learning with diversity constraint for human labeling at each iteration. (Stop human labeling if the labeled set achieves 5% of the testing data.)
- Adaptively re-train the classifier.

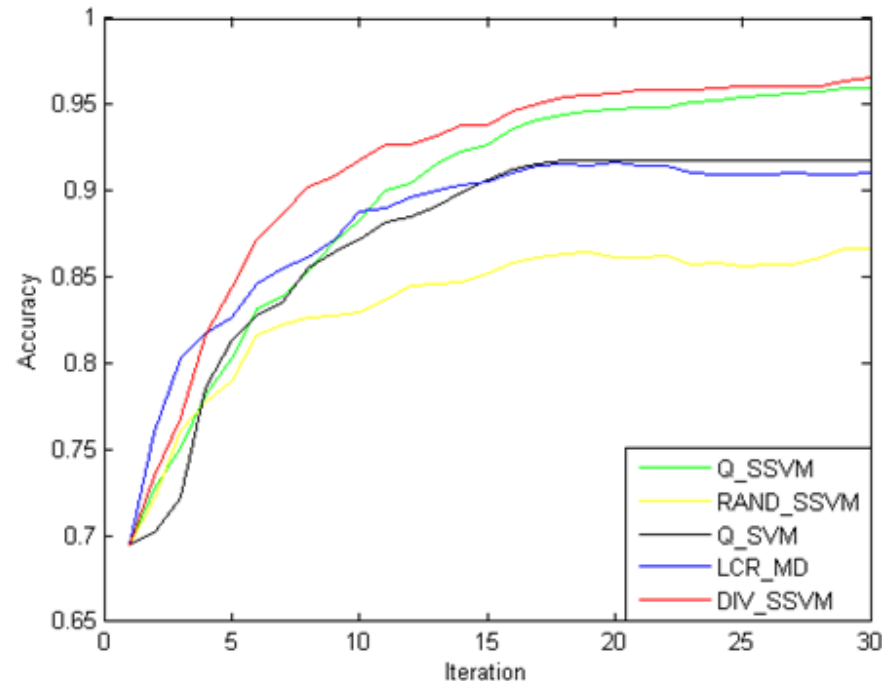


Summary of Query Learning

- Some Properties:
 - 1. **Iteratively** select **informative** samples for human labeling.
 - 2. Consider **uncertainty** and **diversity** in the selection.
 - **Uncertainty**: ambiguity of unlabeled data.
 - **Diversity**: data with large **dissimilarity**.
 - 3. Combine **semi-supervised** learning.
 - Add **very confident** unlabeled data to the training set.
 - 4. Update training samples and **re-train** the classifier for each iteration.

Results

- Training set: 2015 dataset+2016 dataset (5%)
- Accuracy: 96.8% (88.1% to 96.8%)



Q_SSVM: Query learning with semi-supervised learning without diversity constraint.
RAND_SSVM: Query learning based on random sample selection.
Q_SVM: Query learning without semi-supervised learning.
LCR_MD: (Leng et al. 2013).
DIV_SSVM: Query learning with both semi-supervised learning and diversity constraint.

Thank you!

Q & A