

Active Learning

Nils Murrugarra Llerena
University of Pittsburgh

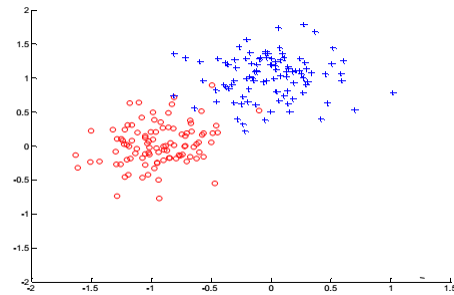
2

Outline

- Introduction
- Why to use active learning?
- Scenarios
- Query Strategies
- Analysis
- Extensions
- Practical Considerations
- Related Research Areas
- Conclusion

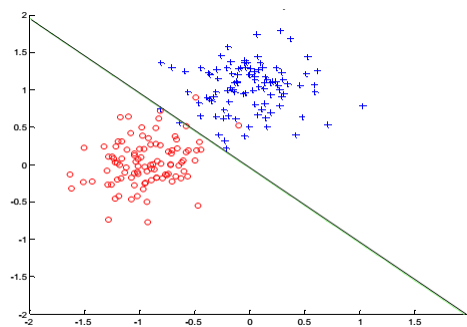
3

Introduction



4

Introduction



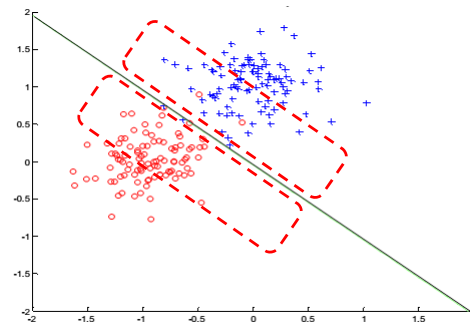
Classifier learned

5

Introduction

Active Learning

- If a learning algorithm is allowed to choose data from which to learn, it will perform better with less training data.
- This means that if the classifier learns the instances that are more “hard” to classify that will be a good classifier using less data.



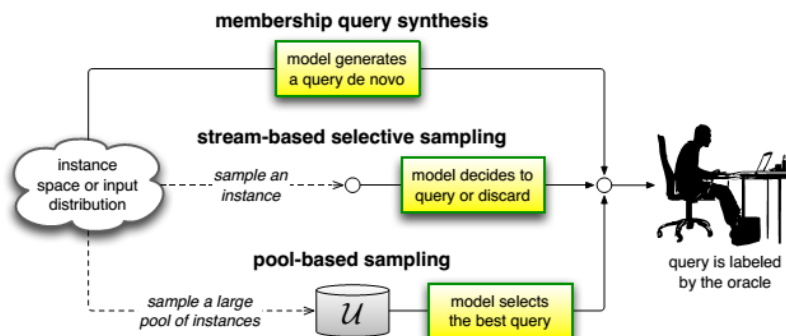
6

Why active learning?

- There are many tasks where labels are: **time-consuming** and/or **expensive** to obtain.
 - **Speech Recognition**
 - Trained Linguistics needed
 - Annotation at word level takes longer time than the audio length
 - **Information Extraction**
 - Finding entities and relations in a news text can take half-hour or more
 - Need some expertise in medical domains
 - **Classification and Filtering**
 - Annotating thousands of data examples can be tedious and redundant

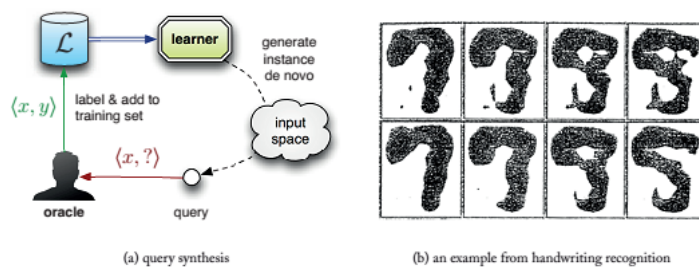
7

Scenarios – Active Learning



8

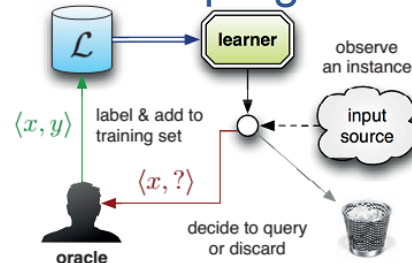
Scenarios: Membership Query Synthesis



- The learner has a definition of the input space
 - Feature dimensions and ranges are known.
- **Problem:**
 - Many generated images doesn't contain recognizable numbers (neural network).

9

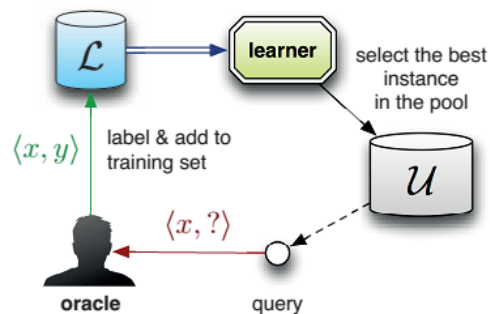
Scenarios: Stream-Based Selective Sampling



- Obtain an unlabeled instance is free.
- **How to query?**
 - Use an "informativeness measure" or "query strategy" (threshold).
 - Compute a region of uncertainty and pick instances that fall in that region.
- **Applications:**
 - news on the web, computer network traffic, phone conversations, ATM transactions, web searches, sensor data, ...
- **Advantage:**
 - Suitable for mobile and embedded devices (memory and power is limited)

10

Scenarios: Pool-based active learning

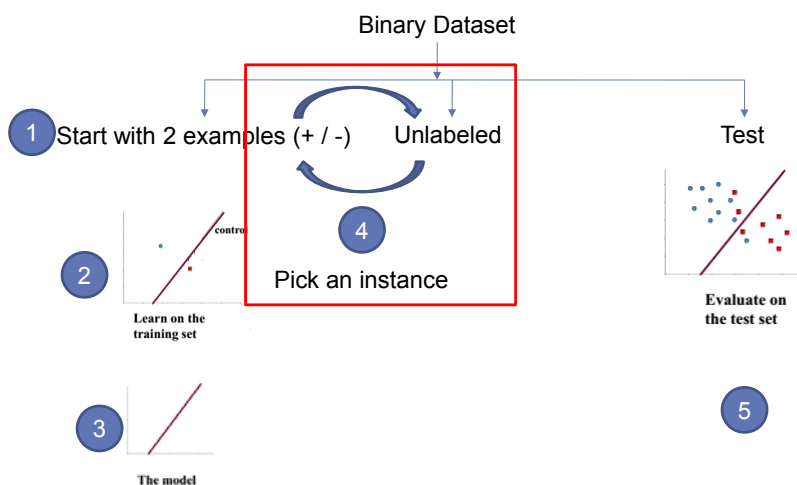


- Large amount of data can be collected at once.
- **How to select?**
 - Use an informativeness measure to evaluate all instances.
- **Comparison:**
 - **Pool-based:** ranks the entire collection.
 - **Stream-based:** scans the data sequentially and make individual decisions.

11

Scenarios: Example

Steps in an Active Learning Approach



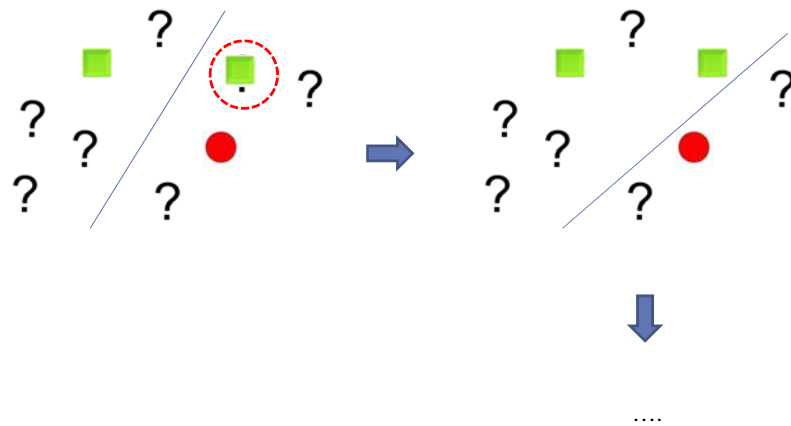
12

Query Strategies

How we evaluate the informativeness of unlabeled instances?

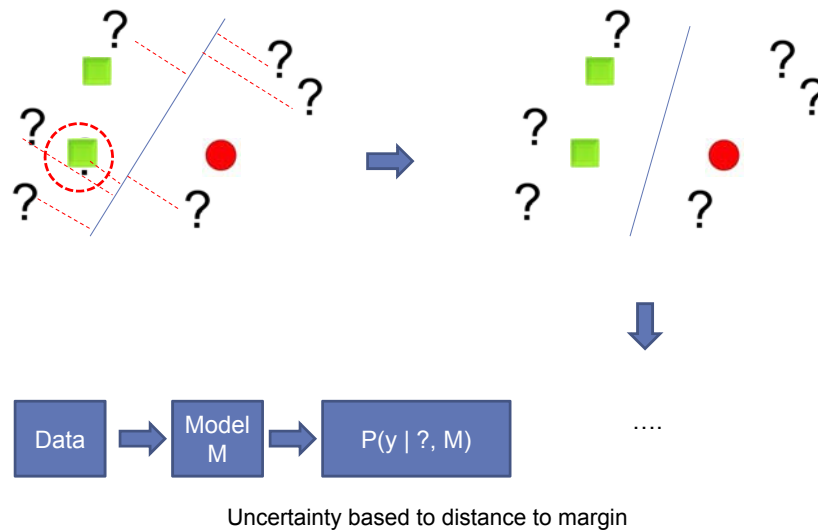
13

Query Strategy: Random



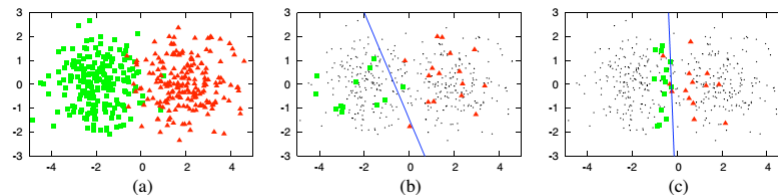
14

Query Strategy: Uncertainty



15

Query Strategies - Potential of AL: Random versus Uncertainty Sampling



Two class Gaussians

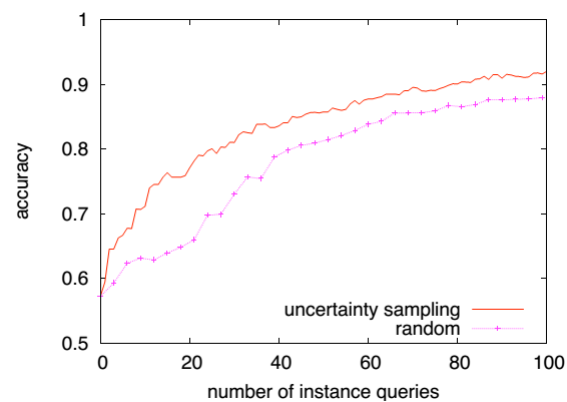
30 instances
selected randomly
[70% accuracy]

30 instances
selected using
uncertainty
sampling [90%
accuracy]

How we compare Query Strategies?

16

Query Strategies – AL Evaluation



Learning Curves

Learning curves are employed to compare query strategies

17

Query Strategy: Uncertainty Sampling

Least Confident

- Query an instance for which the learner is least certain how to label it.
 - **Two classes:** Select the instance whose positive posterior probability is near 0.5
 - **Three or more:** Select the instance whose prediction is the least confident.

$$x_{LC}^* = \arg \max_x (1 - P_{\theta}(\hat{y} | x))$$

\hat{y} : class label with the highest probability

$$\hat{y} = \arg \max_y P_{\theta}(y | x)$$

18

Query Strategy: Uncertainty Sampling

$$x_{LC}^* = \arg \max_x (1 - P_{\theta}(\hat{y} | x))$$

- $P \approx 0$, produce a **higher value** (1) => Pick least certain classifier
- $P \approx 1$, produce a **lower value** (0)

The model's belief that it will mislabel x.

Drawback

- It only considers information about the most probable label.
 - Throws away information about the remaining label distribution.

19

Query Strategy: Uncertainty Sampling

Margin Sampling

$$x_M^* = \arg \min_x (P_\theta(\hat{y}_1 | x) - P_\theta(\hat{y}_2 | x))$$

\hat{y}_1 and \hat{y}_2 : first and second most probable class labels under the model θ

- **Large margin**, instances easy to differentiate
- **Small margin**, more **ambiguous** to differentiate

Drawback

- For very large label sets, the margin approach still ignores the output distribution of the remaining classes.

How to incorporate all labels distribution?

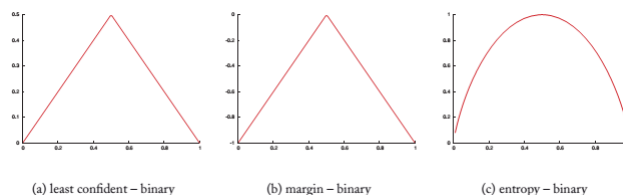
20

Query Strategy: Uncertainty Sampling

Entropy

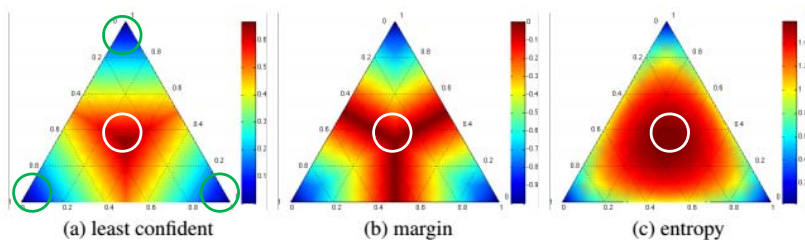
$$\begin{aligned} x_H^* &= \arg \max_x (H_\theta(Y | x)) \\ &= \arg \max_x \left(- \sum_y P_\theta(y | x) \log P_\theta(y | x) \right) \end{aligned}$$

- Is a measure of variable's average information content.
 - **Impurity measure**
 - **Worst case**, (2 classes), probability 0.5
- Measure if all labels have very similar classification probabilities



21

Query Strategies: Uncertainty Sampling



3 Classifiers: c_1 vs (c_2, c_3) ; c_2 vs (c_1, c_3) and c_3 vs (c_1, c_2)

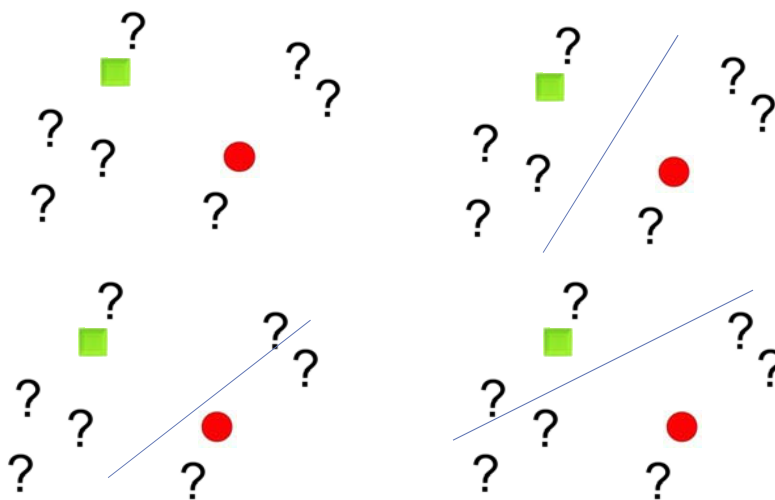
Advantage (entropy)

- Doesn't favor instances where only one of the labels is highly unlikely.
 - Entropy, minimize log-loss
 - Least confident and margin, reduce classification error.

Why if we incorporate more than one model ?

22

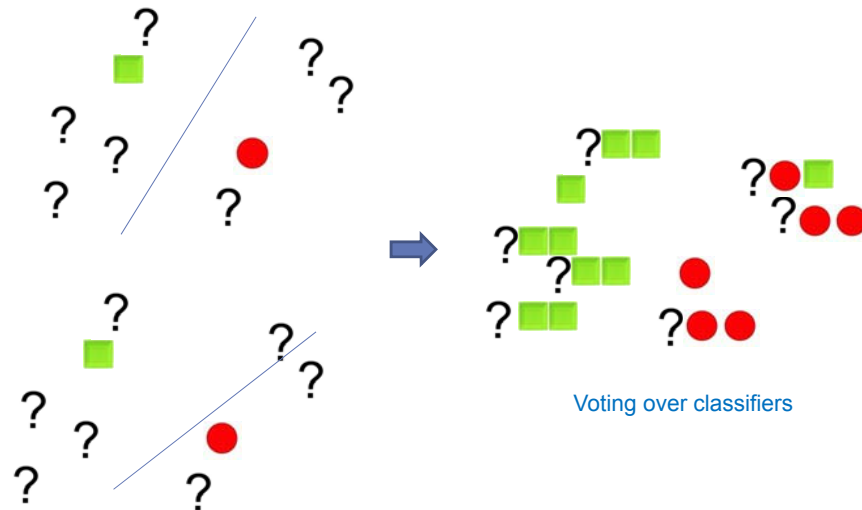
Query Strategy: Query-By-Committee



Different model initialized randomly by a Perceptron approach

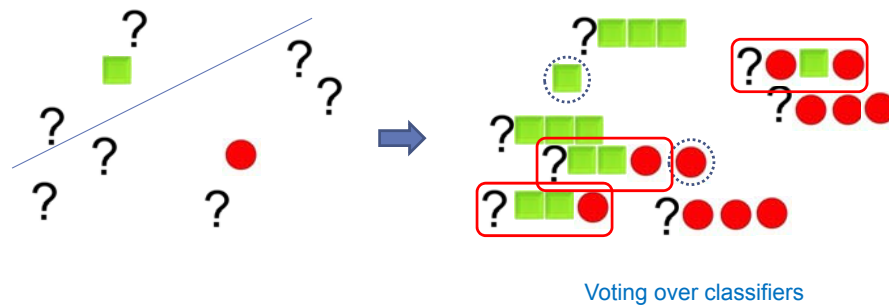
23

Query Strategy: Query-By-Committee



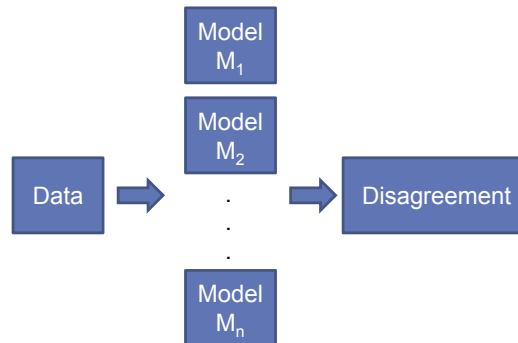
24

Query Strategy: Query-By-Committee



25

Query Strategy: Query-By-Committee



How to generate different models?

26

Query Strategy: Query-By-Committee

How to generate different models?

- Use a bootstrap procedure (e.g. bagging) to subsample the labeled dataset.
- Try different parameters in the classifier
 - Radial SVM, change gamma and cost parameters
 - Decision trees, try different pruning algorithms.

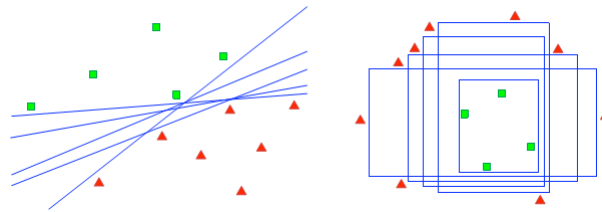


27

Query Strategy: Query-By-Committee

Maintain a committee of learners $C = \{\theta^{(1)}, \dots, \theta^{(C)}\}$, which are all trained in the labeled set L (or subsets).

- Each learner **vote** on the label of the query candidate
- Pick the instance where they most **disagree**.



Considerations

- Consider learners that represent different regions
- Have a measure of disagreement among the learners

How measure disagreement for more than 2 classes?

28

Query Strategy: Query-By-Committee

Disagreement measures

- **Vote entropy**

$$x_{VE}^* = \arg \max_x \left(- \sum_i \frac{V(y_i)}{C} * \log \frac{V(y_i)}{C} \right)$$

Measure Impurity

- $V(y_i)$, number of votes the label y_i receives
- C , committee size

- **KL - Divergence**

$$x_{KL}^* = \arg \max_x \left(\frac{1}{C} * \sum_{c=1}^C D(P_{\theta^{(c)}} \| P_C) \right)$$

$$D(P_{\theta^{(c)}} \| P_C) = \sum_i P_{\theta^{(c)}}(y_i | x) * \log \frac{P_{\theta^{(c)}}(y_i | x)}{P_C(y_i | x)}$$

- $\theta^{(c)}$, a model in the committee
- C , all the committee

$$P_C(y_i | x) = \frac{1}{C} \sum_{c=1}^C P_{\theta^{(c)}}(y_i | x)$$

29

Query Strategy: Query-By-Committee

Disagreement measures

• KL – Divergence

- It measures the difference between two probabilities
- **Most informative query**: Instance that has the largest average difference between:
 - any one committee member
 - and the consensus (all learners)

• KL - Divergence

$$x_{KL}^* = \arg \max_x \left(\frac{1}{C} * \sum_{c=1}^C D(P_{\theta(c)} \parallel P_C) \right)$$

$$D(P_{\theta(c)} \parallel P_C) = \sum_i P_{\theta(c)}(y_i | x) * \log \frac{P_{\theta(c)}(y_i | x)}{P_C(y_i | x)}$$

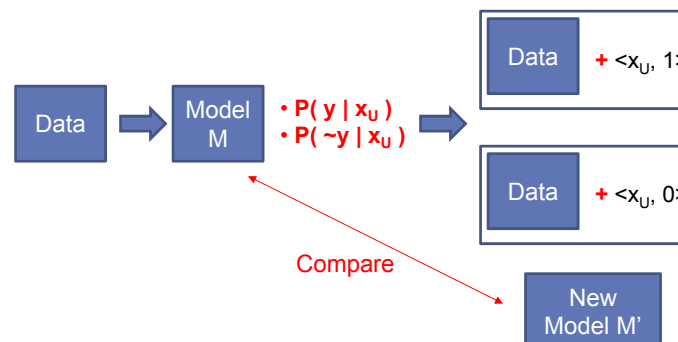
- $\theta^{(c)}$, a model in the committee
- C , all the committee

$$P_C(y_i | x) = \frac{1}{C} \sum_{c=1}^C P_{\theta(c)}(y_i | x)$$

30

Query Strategy: Expected Model Change

- Select the instance that would impact the greatest change to the current model



Compare and quantify the change due to the point inclusion in the labeled set

31

Query Strategy: Expected Model Change

•Expected Gradient Length (EGL)

- Can be applied to any learning algorithm that uses gradient based parameter training
- It determines the importance of the data point with respect to its influence on the model parameters (their change)

$\nabla E(\theta)$: Gradient of error E with respect to the current model θ (M)

$$\nabla E(\theta) = \left[\frac{\partial E}{\partial \theta_1}, \frac{\partial E}{\partial \theta_2}, \dots, \frac{\partial E}{\partial \theta_m} \right]$$

- instance $\langle x_i, y \rangle$ is selected

$\nabla E_i^+(\theta)$: new gradient by adding $\langle x_i, 1 \rangle$

$\nabla E_i^-(\theta)$: new gradient by adding $\langle x_i, 0 \rangle$

Combine

$$= o_i \|\nabla E_i^+(\theta)\| + (1 - o_i) \|\nabla E_i^-(\theta)\| \quad \|x\| = \sqrt{x_1^2 + \dots + x_n^2}$$

32

Query Strategy: Expected Model Change

•Expected Gradient Length (EGL)

- **How to measure the impact/change?**: Consider the norm of the training gradient (i.e. vector used to re-estimate parameter values).

$$x_{EGL}^* = \arg \max_x \left(\sum_i P_\theta(y_i | x)^* \|\nabla \ell_\theta(T \cup \langle x, y_i \rangle)\| \right)$$

$\hookrightarrow \nabla \ell_\theta(T \cup \langle x, y_i \rangle) \approx \nabla \ell_\theta(\langle x, y_i \rangle)$

We don't know the correct label y , for that we consider an expectation over all possible labels.

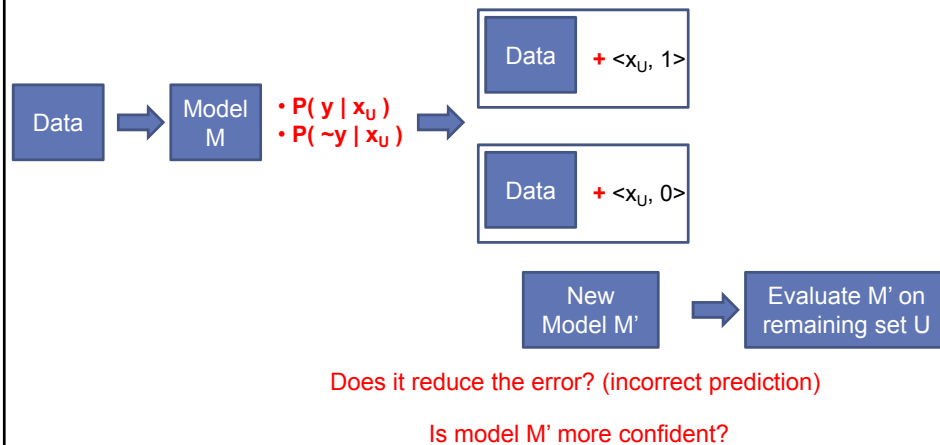
Drawback

- **Computational expensive**, if both the feature space and set of labels are very large

33

Query Strategy: Expected Error Reduction

- Select the instance that **reduce** the **generalization error**.



34

Query Strategy: Expected Error Reduction

- Select the instance that **reduce** the **generalization error**.

- Minimize the Expected 0/1-loss function

$$x_{0/1}^* = \arg \min_x \left(\sum_i P_{\theta}(y_i | x) * \left(\sum_{u=1}^U \underbrace{(1 - P_{\theta+\langle x, y_i \rangle}(\hat{y} | x^{(u)}))}_{\text{Loss function on unlabeled data (\# of incorrect predictions)}} \right) \right)$$

New model after train with $\langle x, y_i \rangle$

Goal: Reduce the expected total number of incorrect predictions.

35

Query Strategy: Expected Error Reduction

- Reduce expected entropy over U

$$x_{\log}^* = \arg \min_x \left(\sum_i P_{\theta}(y_i | x) * \overbrace{\left(\sum_{u=1}^U - \sum_j \underbrace{P_{\theta+\langle x, y_i \rangle}(y_j | x^{(u)}) * \log P_{\theta+\langle x, y_i \rangle}(y_j | x^{(u)})}_{\text{Entropy}} \right)}^{\text{Entropy over U}} \right)$$

Goal: Increase confidence in prediction (minimize entropy).

Drawback

- Most computational expensive framework,
 - require estimate the future error over U for each query
 - a new model is retrained for each query (iterate over all the pool)
- Usually employed in binary classification tasks.

36

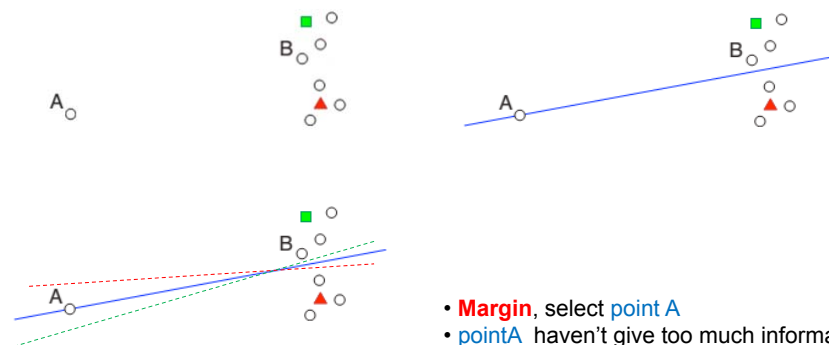
Query Strategy: Density-Weighted Methods

Previous Approaches

- Uncertainty, QBC and EGL are more likely to pick outliers
 - Uncertainty: See example
 - QBC and EGL could pick possible outliers
 - Controversial
 - Generate significant change in the model
- Expected error avoid the previous problems (less probable to pick outliers)
 - Because they focus on the entire input space than individual instances.

37

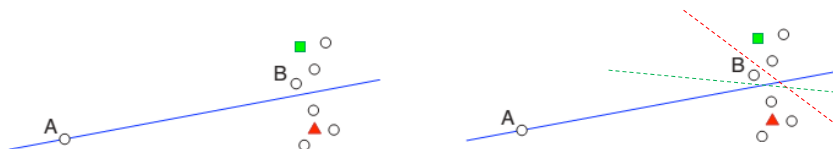
Query Strategy: Density-Weighted Methods



- **Margin**, select **point A**
- **point A** haven't give too much information.
 - Decision boundary almost the same.
 - **data distribution** as a **whole**

38

Query Strategy: Density-Weighted Methods



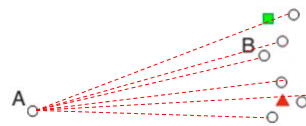
- **point B** is more **informative** about the data distribution

How to combine informativeness and data distribution information?

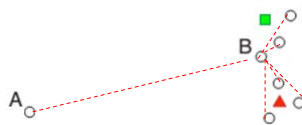
39

Query Strategy: Density-Weighted Methods

Data Distribution



Get average distance



- Distance ≈ 0 , similar examples (Dissimilarity measure)

Similarity measure

- Similar examples, value ≈ 1

40

Query Strategy: Density-Weighted Methods

Model the input distribution during the query selection

- Define informative instances as:
 - uncertain
 - are “representative” of the data distribution

Average similarity to all other instances

- ≈ 1 , more similar with all data

$$x_{ID}^* = \arg \max_x \underbrace{\phi_A(x)}_{\text{Informativeness of query}} * \underbrace{\left(\frac{1}{U} * \sum_{u=1}^U sim(x, x^{(u)}) \right)}_{\text{Average similarity to all other instances}} \underbrace{\beta}_{\text{Control parameter}}$$

Informativeness of query (e.g. uncertainty sampling)

- ≈ 1 , more informative

41

Analysis of Active Learning

Empirical Analysis

- AL helps to reduce the number of labeled instances required to achieve a certain accuracy in the majority of reported results.

Theoretical Analysis

- **Would be Nice!!**
 - Sort of bound in the number of queries to learn a sufficient accurate model
 - This number should be less than passive learning.
- Let's consider instances in one-dimensional line and our model is:

$$g(x; \theta) = \begin{cases} 1 & \text{if } (x > \theta), \text{ and} \\ 0 & \text{otherwise} \end{cases}$$

42

Analysis of Active Learning

Theoretical Analysis

- Let's consider instances in one-dimensional line and our model is:

$$g(x; \theta) = \begin{cases} 1 & \text{if } (x > \theta), \text{ and} \\ 0 & \text{otherwise} \end{cases}$$

According to **PAC model**

- The data distribution can be perfectly classified with $O(1/e)$ random labeled instances.

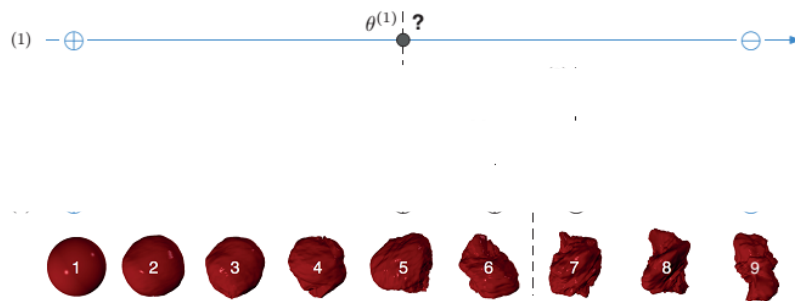
Pool-based AL

- Consider the point on a real line: their labels are a sequence of 0's and 1's.
- **Goal:** Discover the location where the transition occurs

43

Analysis of Active Learning

Theoretical Analysis



Pool-based AL

- Using a binary search. A classifier with error less than ϵ can be obtained with $O(\log 1/\epsilon)$

44

Analysis of Active Learning

Theoretical Analysis

According to **Bayesian Assumption**

- It is possible to achieve generalization error ϵ after seeing $O(d/\epsilon)$ unlabeled instances (d is the VC dimension).

Stream-based and Pool-based AL (QBC)

- It is possible to achieve generalization error ϵ , requesting only $O(d \log 1/\epsilon)$
- Exponential improvement

45

Extensions of Active Learning

AL for Structured Outputs

- Sequential models can produce a probability distribution for every possible label sequence \mathbf{y} , the number of which can grow exponentially in the sequence length T .
- Least confident approach is famous in this setting, because the most likely output sequence $\hat{\mathbf{y}}$ and the associated $P_{\theta}(\hat{\mathbf{y}} | \mathbf{x})$ can be efficiently computed with dynamic programming (Viterbi algorithm).



46

Extensions of Active Learning

Active Feature Acquisition

Instances may have incomplete feature descriptions

- Credit card company can have access to their clients information but not the transactions for other credit companies
- For medical diagnosis, can have access to some basic symptoms, but not all (complex, expensive or risky procedures)

Goal: Select most informative feature to obtain (request) [train time]

Solution:

- Impute the missing values and then acquire the ones that the model is less certain

47

Extensions of Active Learning

Active Classification

Missing feature values can be acquired at test time.

Active Class Selection

Query an instance of a given class label

Active Clustering

Subsample unlabeled instances in a way that they self-organize into groups:

- less overlap or noise

48

Practical Considerations

Batch-Mode Active Learning

Majority of active learning techniques consider that queries are selected one at a time.

- time to induce a model is **expensive**
- All process is **inefficient**

Goal: Query instances in groups.

How to select the optimal query set?

- k-best queries doesn't work properly
 - it fails to consider overlap information in k-best instances
- Most approaches use greedy heuristics that instances in the query are diverse and informative.
- e.g. query centroids of clusters that lie closes to the decision boundary

49

Practical Considerations

Noisy Oracles

Even if labels come from human experts, they might **not be reliable**:

- Some instances are really difficult to annotate
- People can be distracted or fatigued over time

How to use non-experts as oracles?

- Averaging labels of multiple non-experts

50

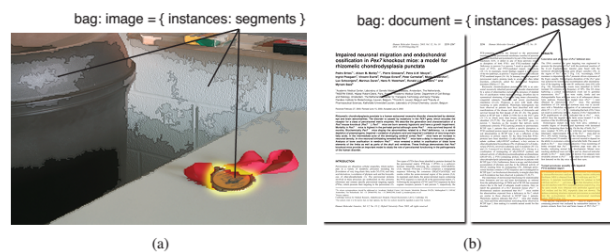
Practical Considerations

Alternative Query Types

• Multiple-instance Active Learning

Instances are grouped in bags:

- labeled **negative**, if all of its instances are negative
- labeled **positive**, if at least one instance is positive



Advantages

- Coarse labels sometimes are available at low cost.
- Allowed to query for labels are finer granularity.
- Could consider approaches of mixed-granularity.

51

Practical Considerations

Alternative Query Types

• Tandem Learning

- Interleave instance-label queries with feature-salient queries.
- e.g. is the word “ball” a discriminative feature for sport documents ?

Multi-Task Active Learning

Same instances may be labeled in multiple ways for different subtasks.

- [parsing and NER](#)
 - Alternating
 - Rank-combination, each task rank the queries and select the highest combined rank
- [Images for binary classification tasks](#).

Stopping Criteria

When accuracy has reached a non-change state?

- Use intrinsic measure of stability within the learner.
 - If the measure degrades, STOP active learning
- Real Stop, based on economic factors (before intrinsic measures)

52

Related Research Areas

Semi-supervised Learning (SSL)

In conjunction with AL, they try to get the most out of the unlabeled data

- [Self training](#) pick the [most confident](#) unlabeled instance. In contrast, AL uncertainty sampling pick the least confident instance.
- [Co-training](#) consider ensemble methods as QBC consider them for AL.

AL and SSL attack the problem from opposite directions

Reinforcement Learning

- In order to improve
 - the learner must take risks and try actions for which it is uncertain about the final result (as AL)

Equivalence Query Learning

- Similar to membership query learning
- It generates an hypothesis of the target concept class
 - The oracle confirm or deny the hypothesis

53

Related Research Areas

Model Parroting and Compression

- Neural Networks achieve better generalization accuracy than decision trees in many applications.
- Decision trees are more comprehensible by humans.

Proposal: Extract high accurate decision trees from neural networks.

AL

- Consider an “oracle model”, trained using a small set of the available labeled data
- Consider a “parrot model”, that can query using the “oracle model”
 - label of any unlabeled data (pool-based)
 - Synthesize new instances (membership-query)

54

Conclusions

- AL is a growing research area
 - Data is **easy** to **obtain**
 - **Difficult/costly** to **label**
- AL has been studied related to:
 - scenarios
 - query strategies
 - Extensions
 - Practical Considerations
 - Related Areas
- However there are still much work to do and open questions ...



Questions

