

Active learning in sequence labeling

Tomáš Šabata

11. 5. 2017

Czech Technical University in Prague
Faculty of Information technology
Department of Theoretical Computer Science

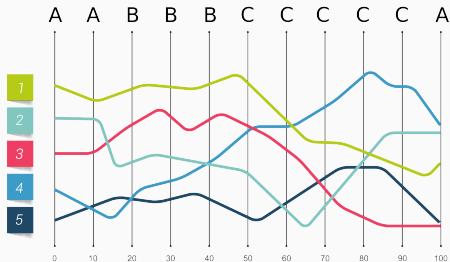
Table of contents

1. Introduction
2. Active learning
3. Graphical models
4. Active learning in sequence labeling
5. Semi-supervised active learning in sequence labeling
6. Experiment
7. Summary

Introduction

Sequence modeling and labeling problem definition

- Sequence of states
- Sequence of observations



Obrázek 1: Sequence representation

1. Sequence modeling

- Given a sequence of states/labels and sequence of observations, find a model that the most likely generates the sequences.

2. Sequence labeling

- Given a sequence of observations, determine an appropriate label/state for each observation.
- Reducing errors by considering relations.

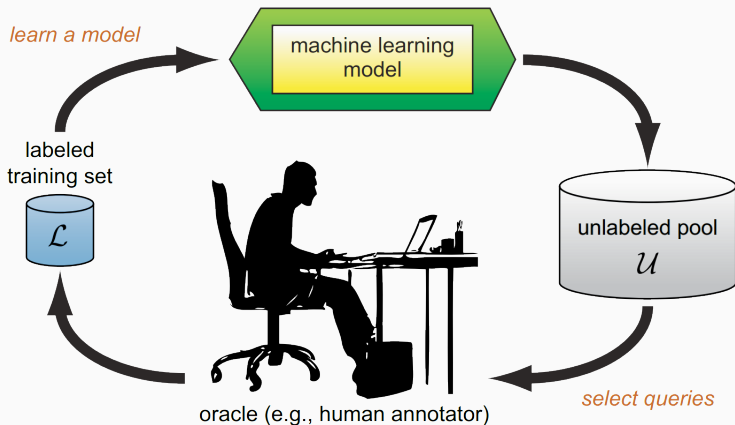
- Handwriting recognition
- Facial expression dynamic modeling
- DNA analysis
- Part-of-speech tagging
- Speech recognition
- Video analysis

Active learning

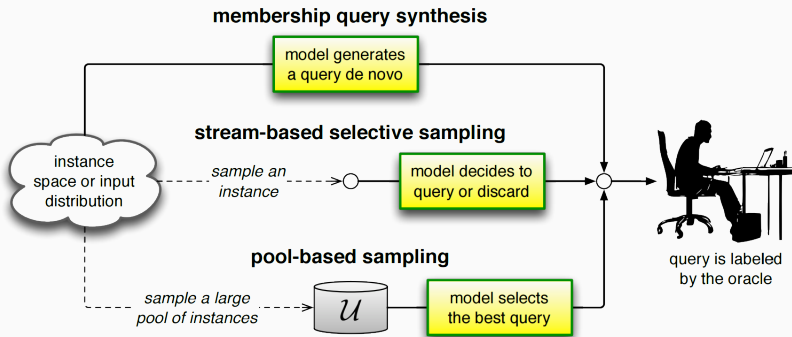
What is active learning?

- *The quality of labels makes a huge difference.
Garbage in, garbage out.*
- *Obtaining "golden" annotation data can be really expensive.*

What is active learning?



Obrázek 2: Active learning cycle



Obrázek 3: Active learning scenarios

AL algorithm

Data:

L - set of labeled examples

U - set of unlabeled examples

θ - utility function

while *stopping criterion is not met* **do**

1. learn model M from L;
2. for all $x_i \in U : u_{x_i} \leftarrow \theta_M(x_i)$;
3. select example $x^* \in U$ with the highest utility function u_i ;
4. query annotator for label of example x^* ;
5. move $\langle y, x^* \rangle$ to L;

end

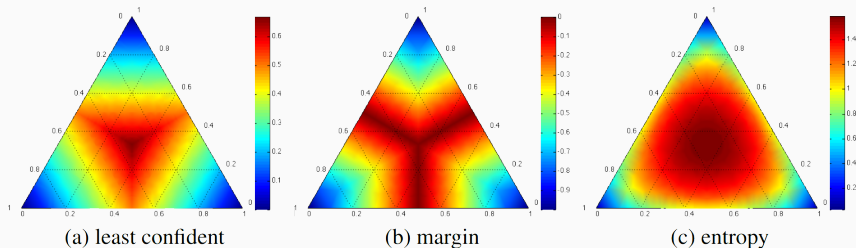
return L

Algorithm 1: General pool-based AL framework

1. Uncertainty Sampling
2. Query-By-Committee
3. Expected Model Change
4. Expected Error Reduction
5. Variance Reduction
6. Density-Weighted Methods

- Simplest, most commonly used
- Intuitive for probabilistic learning model
- Binary problems: choose instance with posterior probability near to 0.5
- Multiclass problems:
 - Least confident - $x_{LC}^* = \arg \max_x 1 - P_\theta(\hat{y}|x)$
 - Margin sampling - $x_M^* = \operatorname{argmin}_x P_\theta(\hat{y}_1|x) - P_\theta(\hat{y}_2|x)$
 - Entropy - $x_H^* = \operatorname{argmax}_x - \sum_i P_\theta(y_i|x) \log P_\theta(y_i|x)$

Frameworks: Uncertainty Sampling



Obrázek 4: Uncertainty sampling for three-class classification problem

- Application dependent
- Entropy - minimizing log-loss
- LC + Margin - minimizing classification error

- We maintain a committee $\mathcal{C} = \{\theta_1, \dots, \theta_C\}$ of models trained on \mathcal{L}
- The most informative query is considered to be the instance about which they most disagree.
- We need to ensure variability of models in the beginning
- Measure of disagreement:
 - Vote entropy - $x_{VE}^* = \operatorname{argmax}_x - \sum_i \frac{V(y_i)}{C} \log \frac{V(y_i)}{C}$
 - Kullback-Leibler divergence - $x_{KL}^* = \operatorname{argmax}_x \frac{1}{C} \sum_{c=1}^C D_{KL}(P_{\theta(c)} || P_C)$

Frameworks: Expected Model Change

- Capable for models using gradient based training.
- Query instance which would cause the largest model change.
- Use a gradient of the objective function $\nabla l_{\theta}(\mathcal{L})$
- $x_{EMC}^* = \operatorname{argmax}_x \sum_i P_{\theta}(y_i|x) \|\nabla l_{\theta}(\mathcal{L} \cup \langle y_i, x \rangle)\|$
- Note: $\|\nabla l_{\theta}(\mathcal{L})\|$ should be close to zero therefore we can use an approximation $\|\nabla l_{\theta}(\mathcal{L} \cup \langle y_i, x \rangle)\| \approx \|\nabla l_{\theta}(\langle y_i, x \rangle)\|$

Frameworks: Expected Error Reduction

- Estimate the expected future error of a model trained on $\mathcal{L} \cup \langle x, y \rangle$
- Methods:
 - Minimizing the expected 0/1-loss
$$x^* = \operatorname{argmin}_x \sum_i P_{\theta}(y_i|x) \left(\sum_{u=1}^U \mathbf{1} - P_{\theta+\langle x, y_i \rangle}(\hat{y}|x^{(u)}) \right)$$
 - Minimizing the expected log-loss
$$x^* = \operatorname{argmin}_x \sum_i P_{\theta}(y_i|x) \left(- \sum_{u=1}^U \sum_j P_{\theta+\langle x, y_i \rangle}(y_j|x^{(u)}) \log P_{\theta+\langle x, y_i \rangle}(y_j|x^{(u)}) \right)$$
- In most cases the most computationally expensive query framework
 - Logistic regression - $\mathcal{O}(ULG)$
 - CRF - $\mathcal{O}(TM^{T+2}ULG)$

- Use the bias-variance decomposition
- $E_T[(\hat{y} - y)^2|x] =$
 $E[(y - E[y|x])^2] + (E_{\mathcal{L}}[\hat{y}] - E[y|x])^2 + E_{\mathcal{L}}[(\hat{y} - E_{\mathcal{L}}[\hat{y}])^2]$
- Model dependent framework

- Informative instances should not only be those which are uncertain, but also those which are "representative" of the underlying distribution
- Uses one of other query strategies as base query strategy (e.g. uncertainty sampling)
- $x^* = \operatorname{argmax}_x \phi_A(x) \times \left(\frac{1}{U} \sum_{u=1}^U \operatorname{sim}(x, x^{(u)}) \right)^\beta$
- The method is more robust to outliers in dataset.

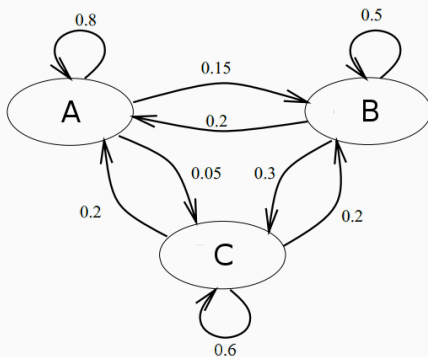
Active learning problem variants

- Active Learning for Structured Outputs
 - Instance is not represented by a single feature vector, but rather a structure.
 - e.g.: Sequences, trees, grammars.
- Active Feature Acquisition
 - Selection of salient unused features
 - e.g.: Medical tests, sensitive information
- Active Class Selection
 - Learner is allowed to query a known class label, and obtaining each instance incurs a cost.
- Active Clustering
 - Generate (or subsample) instances in such a way that they self-organize into groupings
 - Try to get less overlap or noise than with random sampling

Active learning problem variants

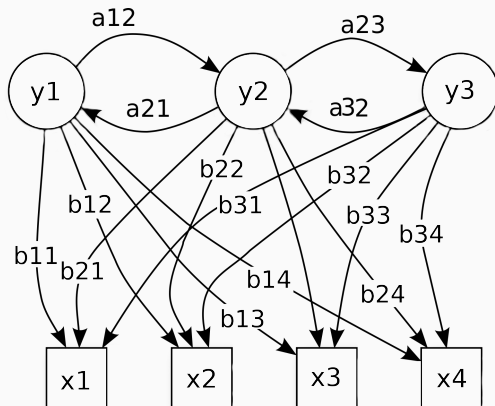
- Active Learning for Structured Outputs
 - Instance is not represented by a single feature vector, but rather a structure.
 - e.g.: Sequences, trees, grammars.
- Active Feature Acquisition
 - Selection of salient unused features
 - e.g.: Medical tests, sensitive information
- Active Class Selection
 - Learner is allowed to query a known class label, and obtaining each instance incurs a cost.
- Active Clustering
 - Generate (or subsample) instances in such a way that they self-organize into groupings
 - Try to get less overlap or noise than with random sampling

Graphical models



Obrázek 5: Markov model/chain

Models: Hidden Markov model



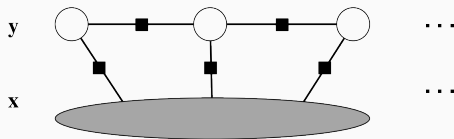
Obrázek 6: Hidden Markov model

Models: Hidden Markov model

$$\lambda = (A, B, \pi)$$

- Set of hidden states $Y = \{y_1, y_2, \dots, y_N\}$, set of observable values $X = \{x_1, x_2, \dots, x_M\}$
- Sequence of states $Q = q_1 q_2 q_3 \dots q_T$ sequence of outputs $O = o_1 o_2 o_3 \dots o_T$
- Transition probability matrix $A = \{a_{ij}\}$
 $a_{ij} = P(q_t = y_j | q_{t-1} = y_i), \quad 1 \leq i, j \leq N$
- Emission probability distribution $B = \{b_{i,j}\}$
 $b_{i,j} = P(o_t = x_j | q_t = y_i), \quad 1 \leq i \leq N, \quad 1 \leq j \leq M$
- Initial probability distribution $\pi = \{\pi_i\}$
 $\pi_i = P(q_1 = y_i), \quad 1 \leq i \leq N$

Models: Conditional random field (linear chain)



Obrázek 7: Linear chain conditional random field

- Discriminative model $P(Y|X)$, we do not explicitly model $P(X)$.
- Perform better than HMMs when the true data distribution has higher-order dependencies than the model.
- $P(Y|X) = \frac{1}{Z(X)} \prod_{t=1}^T \exp\left(\sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, x_t)\right)$
- $Z(X) = \sum_Y \prod_{t=1}^T \exp\left(\sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, x_t)\right)$

Models: Summary

1. HMM

- $P(Q, O) \propto \prod_{t=1}^T P(q_t|q_{t-1}) P(o_t|q_t)$

2. CRF

- $P(Q|O) \propto \frac{1}{Z_O} \prod_{t=1}^T \exp \left(\sum_j \lambda_j f_j(q_t, q_{t-1}) + \sum_k \mu_k g_k(q_t, o_t) \right)$

Active learning in sequence labeling

- Least confident

- $x_{LC}^* = \operatorname{argmax}_x 1 - P_\theta(\hat{y}|x)$
- Viterbi path

- Margin sampling

- $x_M^* = \operatorname{argmin}_x P_\theta(\hat{y}_1|x) - P_\theta(\hat{y}_2|x)$
- N-best algorithm

- Entropy

- Token entropy - $x_{TE}^* = \operatorname{argmax}_x -\frac{1}{T} \sum_{t=1}^T \sum_{m=1}^M P_\theta(y_t = m) \log P_\theta(y_t = m)$
- Total token entropy
- Sequence entropy $x_{SE}^* = \operatorname{argmax}_x - \sum_{\hat{y}} P_\theta(\hat{y}|x) \log P_\theta(\hat{y}|x)$
- N-best sequence entropy $x_{SE}^* = \operatorname{argmax}_x - \sum_{\hat{y} \in \mathcal{N}} P_\theta(\hat{y}|x) \log P_\theta(\hat{y}|x)$

Query by Committee

- Query-by-bagging (each model has unique modified set $\mathcal{L}^{(c)}$)
- Vote entropy - disagreement over Viterbi's paths

$$x_{VE}^* = \operatorname{argmax}_x - \frac{1}{T} \sum_{t=1}^T \sum_{m=1}^M \frac{V(y_t, m)}{C} \log \frac{V(y_t, m)}{C}$$

- Kullback Leibler

$$x_{KL}^* = \operatorname{argmax}_x \frac{1}{T} \sum_{t=1}^T \frac{1}{C} \sum_{c=1}^C D_{KL}(\theta^{(c)} || \mathcal{C})$$

$$D(\theta^{(c)} || \mathcal{C}) = \sum_{m=1}^M P_{\theta^{(c)}}(y_t = m) \log \frac{P_{\theta^{(c)}}(y_t = m)}{P_{\mathcal{C}}(y_t = m)}$$

- Non-normalized variants
- Sequence vote entropy

$$x_{SVE}^* = \operatorname{argmax}_x - \sum_{\hat{y} \in \mathcal{N}^C} P(\hat{y} | x, \mathcal{C}) \log P(\hat{y} | x, \mathcal{C})$$

- Sequence Kullback-Leibler $x_{SKL}^* = \operatorname{argmax}_x \frac{1}{C} \sum_{c=1}^C \sum_{\hat{y} \in \mathcal{N}^C} \log \frac{P_{\theta^{(c)}}(\hat{y} | x)}{P_{\mathcal{C}}(\hat{y} | x)}$

Expected Gradient Length

- Expectation over the N-best labelings.
- $x_{EMC}^* = \operatorname{argmax}_x \sum_{\hat{y} \in \mathcal{N}^c} P_{\theta}(\hat{y}|x) \|\nabla l_{\theta}(\mathcal{L} \cup \langle y_i, x \rangle)\|$

- $\mathcal{O}(TM^{T+2}ULG)$ too expensive :(

- Information density
 - solves problem that US and QBC are prone to querying outliers
- We need a distance measure for sequences.
 - Kullback-Leibler
 - Euclidean distance
 - Cosine distance
- Drawback: number of required similarity calculations grows quadratically with the number of instances in \mathcal{U} .
- Solution: Precompute them.

Semi-supervised active learning in sequence labeling

- FuSAL:
 - Sequence is handled as a whole unit.
 - Sequence-wise vs. token-wise utility functions
- SeSAL
 - Some subsequences can be easily labelled automatically.
 - Decrease labelling effort.
 - Usage of self-training principle.

Fully supervised general algorithm

Data:

B - number of examples to be selected

L - set of labeled examples

U - set of unlabeled examples

θ - utility function

while *stopping criterion is not met* **do**

1. learn model M from L;
2. for all $x_i \in U$: $u_{x_i} \leftarrow \theta_M(x_i)$;
3. select B examples $x_i \in U$ with highest utility function u_i ;
4. annotate sequences using M;
5. query for labels of non-confidential tokens;
6. move newly annotated examples to L;

end

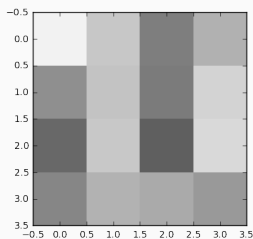
return L

Algorithm 2: General AL framework

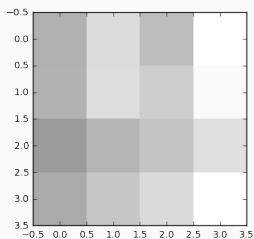
Experiment

Problem definition

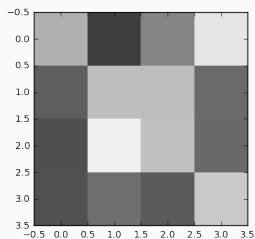
- "Handwritten" letters recognition.
- Each letter is randomly written in one of 6 fonts
- Downscaled to 4x4 pixels
- Letters are organized in real sentences.



Obrázek 8: A



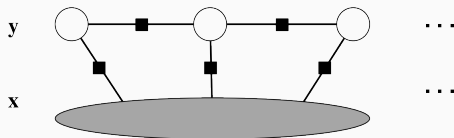
Obrázek 9: B



Obrázek 10: C

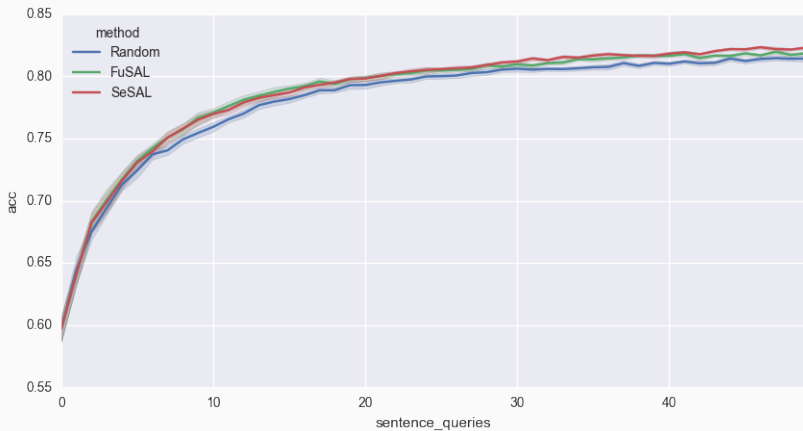
Model and training

- Linear chain conditional random fields.
- 3 labeled sentences in train dataset in the beginning .
- Labeled sentences are added to the dataset iteratively (50 times).
 - Random choice
 - FuSAL (Least Confident)
 - SeSAL (Least Confident + marginal probability)



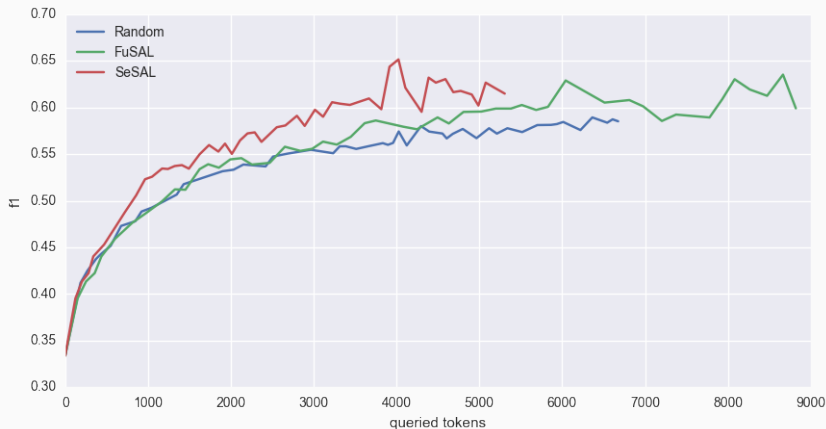
Results

Does it even worth it?

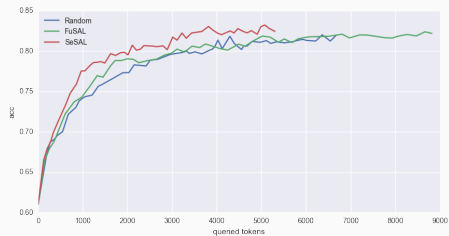
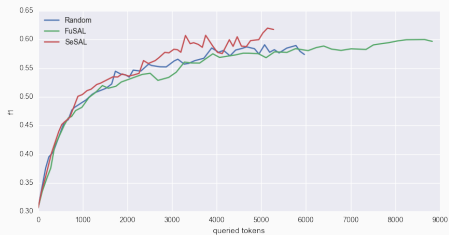


Results

Let's look from another point of view.



Results



Summary

- It does not work in all cases.
- It can be implementation overhead.
- In the most of cases the active learning helps.
- It can be applied to different structures like sequences or trees
- The combination with semi-supervised learning can lead to rapid save of costs.
- Future work: Different query costs, automatic threshold finding, CT-HMM.

**Thank you for your attention.
Questions?**