



# Machine Learning for Sequence Learning

## Learning in an All-Subsequence Space

Severin Gsponer, Georgiana Ifrim, Barry Smyth

January 20, 2016

# Outline

- Background
- Linear Classifiers for Sequences
- SEQL Approach
- Contribution
- Future Work

# Background for Sequence Learning

## Definition of a sequence

A sequence consists of symbols of a given finite alphabet  $\Sigma$  in a given order:  $s_0, s_1, \dots, s_n$

## Examples

- Genetic sequence: **AGCTGTTTCGT** ,  $|\Sigma| = 4$ ,  $\Sigma = \{A, C, G, T\}$
- Protein sequence: **KVKTGCKATLR** ,  $|\Sigma| = 20$
- Text: **The house is blue** ,  $|\Sigma| = 4$ , (# distinct words in corpus)

# Sequence Classification

Class	Data points
+1	C70124045F0*EE*ADC00E9D64A000C6689CCF1C70
+1	7413BAEF01000668951488B7000F0*EE*AD00081CA
-1	08F9C81A80C18B484000895110B8040000C20C00CCC
-1	CCCFF8CC84C8B5C8BC18B484C8B505C8340240481

Find subsequences that can be used to identify the class.

?? CC8CC84C8BC8B458B4CC0F82B505FB4C83B4B0481

# Related Work

## Bag of Words

- Loss of structural order ( e.g., Mary is faster than John)
- Often not accurate enough

## Kernel SVM

- Lift into implicit high-dimensional feature space through kernel trick
- Restrict features for scale (e.g., max 5-gram)
- Not easily interpretable (Blackbox)

## SEQL (Our Approach)

- Works in explicit high-dimensional feature space
- Unrestricted features (i.e. all-length subsequences)
- Interpretable classifier (Whitebox)

# All-Subsequence Feature Space

Sample sequence:	... F09EE1AD ...
Uni-gram (all):	0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F (16 possible)
Bi-gram:	F0, 09, 9E, EE, 1A,... (16 <sup>2</sup> = 256 possible)
Tri-gram:	F09, 09E, EE1, E1A, 1AD,... (16 <sup>3</sup> = 4096 possible)
⋮	⋮
8-gram:	F09EE1AD,... (16 <sup>8</sup> = 4294967296 possible)

Representation of sequence in explicit vectorspace of all subsequences:

$$\mathbf{x}_i = (0, 1, 2, 3, 4, \dots, F, 00, 01, 02, 03, \dots, FF, 000, 0001, \dots)$$

# Linear Sequence Classifier

Given:

Training set of labeled examples:

$$\{x_i, y_i\} \text{ for } i = 1, \dots, N \text{ where } y_i \in \{-1, 1\}$$
$$\mathbf{x}_i \in \mathbb{R}^d \text{ with } d = \text{number of features}$$

Goal:

Find  $\beta = (\beta_1, \beta_2, \dots, \beta_d)$ ,  $\beta_i \in \mathbb{R}$  by optimizing:

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^d} L(\beta) = \arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^N \xi(y_i, x_i, \beta) + CR(\beta)$$

Classical **gradient** descent is computationally infeasible for a large feature space

$$\beta^{(t)} = \beta^{(t-1)} - \eta_t \nabla L(\beta^{(t-1)})$$

---

## Algorithm 1 SEQL workflow

---

Set  $\beta^{(0)} = 0$

**while** !termination condition **do**

    Calculate objective function  $L(\beta^{(t)})$

    Find feature with maximum gradient value

    Find step length  $\eta_t$  by line search

    Update  $\beta^{(t)} = \beta^{(t-1)} - \eta_t \frac{\partial L}{\partial \beta_{j_t}}(\beta^{(t-1)})$

    Add corresponding feature to feature set

**end while**

---



# Contribution

1. Study influence of problem characteristics on classification performance (simulation)
2. Extend SEQL approach to regression (gradient bound for squared error loss)
3. Real-World Applications

# Contribution 1: Simulation Dimensions

- **Alphabet size**  $|\Sigma|$
- **Sequence length**  $L$
- **Data set size**  $N$
- Motif length  $m$
- Sparsity of the feature space
- Noise in the motifs

# Contribution 1: Analysis

## Accuracy

- Classification performance (ACC, AUC, F1, ...)

## Speed

- Number of iterations
- Quality of gradient bound (pruning ration)
- Run time

## Interpretability

- Number of produced features

# Contribution 1: Simulation Framework

Systematic experiments on generated sequences:

Generation of  $N$  sequences of length  $L$

$$l_1, l_2, \dots, l_L$$

where  $l_i \sim U(\text{Alphabet})$

Insert **motifs** of length  $m$  in positive sequences.

Ratio of positive to negative sequences is 1:10

# Contribution 1: Data Generation

1. Random generation of a motif
2. Determine motif insertion **position** randomly for each sequence
3. Random generation of sequence and insertion of motif at **position**

# Contribution 1: Data Generation

---

## Algorithm 2 Positive sequences generation

---

```
Generate motif by drawing  $m$  symbols from  $\sim U(\text{Alphabet})$ 
for  $i < N \cdot 0.1$  do
   $pos \sim U(L - m)$ 
  for  $l < (L - m)$  do
    if  $l = pos$  then
      add motif to sequence
    else
      add symbol  $l \sim U(\text{Alphabet})$  to sequence
    end if
  end for
  add sequence to data set
end for
```

---

## Contribution 2: Extension to Regression

Value	Data points
+0.2	C70124045C00E9D64A000CCCF1C70
+1.4	7413BAEF0100051488B700000081CA
-3.2	08F9C81A80000895110B8040000C20
-0.1	CCF8CC84C8B5C8BC8B505C834024

Implementation of squared error loss and new gradient bound

$$\xi(y_i, x_i, \beta) = \sum_{i=1}^N (y_i - \beta^t x_i)^2$$

With L1 regularization known as LASSO.

**Questions** Influence of loss function and quality of the bound

# Contribution 3: Real World Application

## Classification Task

Microsoft Malware Challenge (BIG 2015)

Kaggle Competition in early 2015

**Goal** Classification of Malware into 9 families

**Data** ~500GB of hexadecimal sequences

## Regression Task

We are still looking for problem domains for sequence regression?



# Future Work

## Regression applications

Test on real world application.

## Rescaling of features

TF-IDF style rescaling of feature instead of binary indicator [1]  
and analysis of influence for the gradient bound quality.

# References

## Bibliography



L. Miratrix and R. Ackerman.

Conducting sparse feature selection on arbitrarily long phrases in text corpora with a focus on interpretability.

pages 1--41, 2015.