

Efficient Sequence Regression by Learning Linear Models in All-Subsequence Space

Severin Gsponer, Barry Smyth, Georgiana Ifrim

19.09.2017

Outline

- Sequence Regression Problem
- Related Work and Background
- Our Approach: Linear Model for Sequence Regression (SqLoss)
- Evaluation

Sequence Regression Problem

Problem Setting

Score	Sequence
290.5	AGTCCACAAGGCTAGGATAGCTATCCGGATCGA
315.1	TATCCTGCAGTACAAGTCCGTAATTCTCAATCCA
805.6	AGTCCGCTAGGCTAGGATAGCTAGCCCGATCGA
799.7	AGCCAAGACCTGAAATAGGCTCCTGAGATACAG
???	CGGGTCGTATCCGCACTGAATATCCAGAGATACG

$$\Sigma = \{A, C, G, T\}$$

Sequence Regression Problem

Problem Setting

Score	Sequence
290.5	AGTCC*C*AGGCTAGGATAGCTATCCGGATCGA
315.1	TATCCTGCAGTACAAGTCCGTAATT*C*AATCCA
805.6	AGTCCGCTAGGCTAGGATAGCTAGCCCGATCGA
799.7	AGCCAAGACCTGAAATAGGCTCCTGAGATACAG
???	CGGGTCGTATCCGCACTGAATATCTAGGCTTACG

$$\Sigma = \{A, C, G, T\}$$

Weight	k-mer
796.6	TAGGCT
402.5	C*C*A
-125.3	TCCG

Related Work

Bag of Words

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

Kernel SVM

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

Neural Networks

- Hard to train
- Time consuming
- Not easily interpretable (Blackbox)

Design Goals

Interpretable and Simple

Linear models with k -mers as features

Accurate

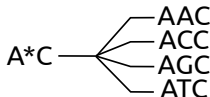
All-length k -mer feature space and wildcards

Efficient

Fast search for most promising features

All-Subsequence Feature Space

k-mers	A	C	G	T	AA	AC	...	CA	CC	...	GTCCTA	GTCCTC	...	TTTTTT
Binary vector	1	1	1	1	1	0	...	0	1	...	1	0	...	0



Sample sequence: ...GTCCTAATCCTA ...

1-mer: A, C, G, T (4 possible)

2-mer: GT, TC, CC, CT, TA, ... ($4^2 = 16$ possible)

3-mer: GTC, TCC, CCT, CTA, TAA, ... ($4^3 = 64$ possible)

⋮

8-mer: GTCCTAAT, TCCTAATC, ... ($4^8 = 65536$ possible)

Regression with Square Loss

Given:

Training set of labeled examples:

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

Goal:

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

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^d} L(\beta) = \arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^N (y_i - \beta^T \mathbf{x}_i)^2 + CR(\beta)$$

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

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

Our Approach: SqLoss Algorithm

Algorithm 1 Coordinate Descent with Gauss Southwell Selection

- 1: Set $\beta^{(0)} = 0$
 - 2: **while** not termination condition **do**
 - 3: Adjust intercept
 - 4: Calculate objective function $L(\beta^{(t)})$
 - 5: **Find coordinate j_t with maximum gradient value**
 - 6: Find optimal step size η_{j_t} by line search or exact optimization
 - 7: Update $\beta^{(t)} = \beta^{(t-1)} - \eta_{j_t} \frac{\partial L}{\partial \beta_{j_t}}(\beta^{(t-1)}) e_{j_t}$
 - 8: Add corresponding feature to feature set
 - 9: **end while**
-

Our Approach: SqLoss Algorithm

Algorithm 1 Coordinate Descent with Gauss Southwell Selection

- 1: Set $\beta^{(0)} = 0$
 - 2: **while** not termination condition **do**
 - 3: Adjust intercept
 - 4: Calculate objective function $L(\beta^{(t)})$
 - 5: **Find coordinate j_t with maximum gradient value**
 - 6: Find optimal step size η_{j_t} by line search or exact optimization
 - 7: Update $\beta^{(t)} = \beta^{(t-1)} - \eta_{j_t} \frac{\partial L}{\partial \beta_{j_t}}(\beta^{(t-1)}) e_{j_t}$
 - 8: Add corresponding feature to feature set
 - 9: **end while**
-

How do we find coordinate j_t efficiently?

Efficient GS Selection via Gradient Bounding

Example

$$s_j = \text{"ACTC"} \text{ and } s_p = \text{"ACT"} \\ \text{gradient}(s_j) \leq \mu(s_p)$$

Key Ideas

Bound gradient of each feature using only information about its prefix.

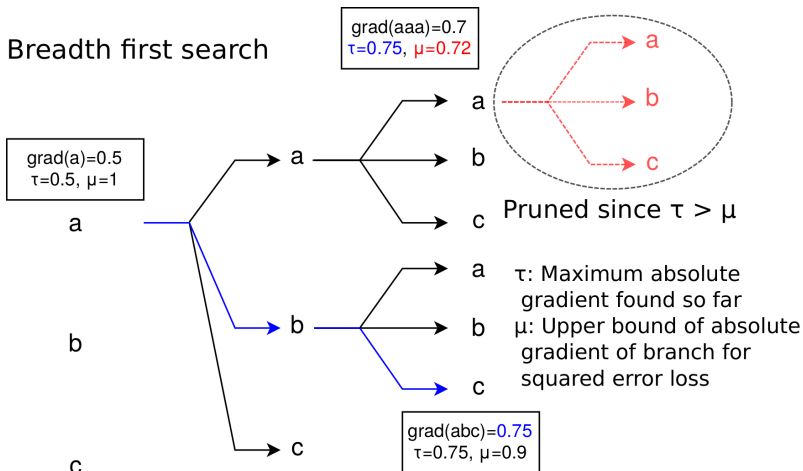
Separate positive and negative terms of the gradient.

Theorem For any subsequence $s_p \subseteq s_j$ it holds:

$$\left| \frac{\partial L}{\partial \beta_j}(\beta) \right| \leq \mu(s_p) = \max \left\{ \sum_{\{i | s_p \in x_i, y_i - \beta^T x_i \geq 0\}} -2 \cdot x_{ip} \cdot (y_i - \beta^T \cdot x_i), \right. \\ \left. \sum_{\{i | s_p \in x_i, y_i - \beta^T x_i \leq 0\}} -2 \cdot x_{ip} \cdot (y_i - \beta^T \cdot x_i) \right\}$$

Search and Pruning

Breadth first search



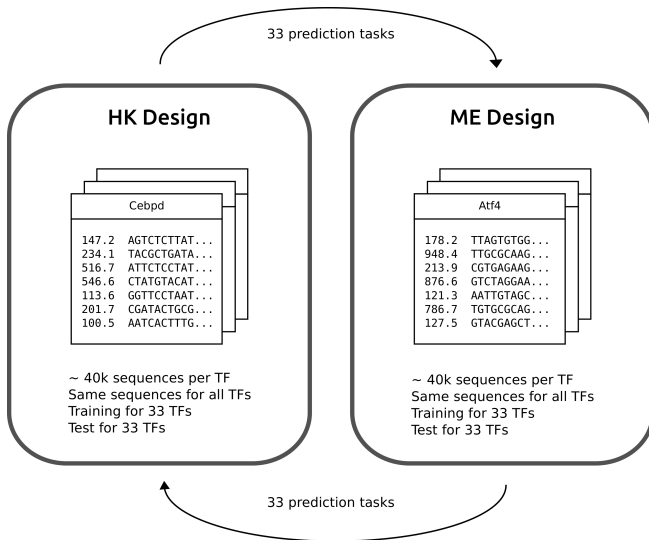
Evaluation

DREAM5 Transcription-Factor Binding Affinity Prediction
DNA-Motif Recognition Challenge in 2011
Collection of 66 regression tasks in a biological domain.

Goal Predict binding affinity of protein (TF) to a given DNA sequence

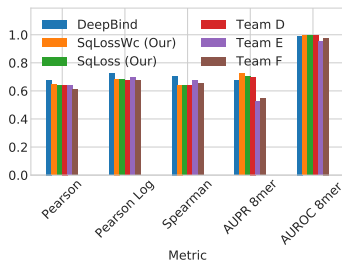
Value	Data points
290.507	AGGGCATCATGGAGCTGTCCAG
679.305	ATCACAATTTTGCCGAGAGCGA
1998.715	GTACACCCCGTTTCGGCGGCCCA
447.803	CCTTTAGCCCATCGTTGGCCAA

TF Binding Affinity Prediction



TF Binding Affinity Results I

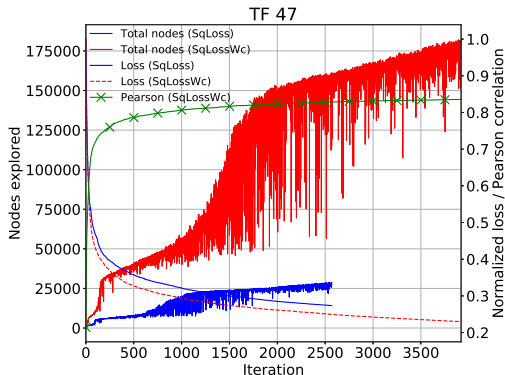
Team	Pearson	Pearson_Log	Spearman	AUPR_8mer	AUROC_8mer
DeepBind	0.6780	0.7260	0.7060	0.6760	0.9910
SqLossWc	0.6483	0.6846	0.6423	0.7236	0.9967
SqLoss	0.6399	0.6791	0.6390	0.7049	0.9953
Team_D	0.6413	0.6742	0.6394	0.6997	0.9942
Team_E	0.6375	0.6936	0.6735	0.5223	0.9524
Team_F	0.6103	0.6732	0.6555	0.5456	0.9766



- Good results across all metrics
- No domain knowledge needed
- Comparable to other linear models

TF Binding Affinity Results II

Motif	Weight
TAAT*A	0.733985
TAATG*G	0.706344
ATG*AAA	0.674507
:	:
GGATA	-0.188202
TCAAT	-0.214858
G*ATAG	-0.218132



- Final model is a list of weighted *k*-mers
- Wildcards improve results but also imply a computational burden

Conclusion & Future Work

- Linear model for sequence regression
- All-length subsequence feature space
- Mean squared error as loss function
- Coordinate descent with Gauss-Southwell selection
- Branch-and-Bound strategy for efficient GS selection

Future Work

- Expand to other data structures
- Ensemble multiple models
- Stochastic variant of our approach

Further Information

Paper, Code and Data available

Efficient Sequence Regression by Learning Linear Models in
All-Subsequence Space

Severin Gsponer, Barry Smyth, Georgiana Ifrim

Code and data available at:
github.com/svgspartner/SqLoss

Email:
severin.gsponer@insight-centre.org

Appendix: Search and Pruning

grad(a)=0.5
 $\delta=0.5, \mu=1$

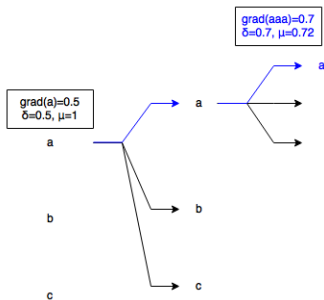
a

b

c

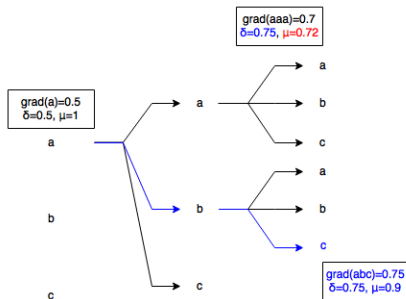
Calculate gradient and bound μ for each node. Save global maximal gradient in δ .

Appendix: Search and Pruning



Continue in a breadth first manner.

Appendix: Search and Pruning



Continue in a breadth first manner.