

Subject Section

GkmExplain: Fast and Accurate Interpretation of Nonlinear Gapped k -mer SVMs

Avanti Shrikumar^{1,*†}, Eva Prakash^{2,†} and Anshul Kundaje^{1,3,*}

[†]Co-first authors

¹Department of Computer Science, Stanford University, Stanford, CA 94305, USA,

²BASIS Independent Silicon Valley, San Jose, CA 95126, USA and

³Department of Genetics, Stanford University, Stanford, CA 94305, USA

*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Summary: Support Vector Machines with gapped k -mer kernels (gkm-SVMs) have been used to learn predictive models of regulatory DNA sequence. However, interpreting predictive sequence patterns learned by gkm-SVMs can be challenging. Existing interpretation methods such as deltaSVM, in-silico mutagenesis (ISM), or SHAP either do not scale well or make limiting assumptions about the model that can produce misleading results when the gkm kernel is combined with nonlinear kernels. Here, we propose GkmExplain: a computationally efficient feature attribution method for interpreting predictive sequence patterns from gkm-SVM models that has theoretical connections to the method of Integrated Gradients. Using simulated regulatory DNA sequences, we show that GkmExplain identifies predictive patterns with high accuracy while avoiding pitfalls of deltaSVM and ISM and being orders of magnitude more computationally efficient than SHAP. By applying GkmExplain and a recently developed motif discovery method called TF-MoDISco to gkm-SVM models trained on *in vivo* TF binding data, we recover consolidated, non-redundant transcription factor (TF) motifs. Mutation impact scores derived using GkmExplain consistently outperform deltaSVM and ISM at identifying regulatory genetic variants from gkm-SVM models of chromatin accessibility in lymphoblastoid cell-lines.

Availability: Code and example notebooks to reproduce results are at <https://github.com/kundajelab/gkmexplain>. **Contact:** avanti@stanford.edu, evaprakash2@gmail.com, akundaje@stanford.edu

1 Introduction

Deciphering the combinatorial regulatory DNA sequence patterns that determine transcription factor (TF) binding and chromatin state is critical to understand gene regulation and interpret the molecular impact of regulatory genetic variation. High-throughput *in vivo* and *in vitro* functional genomics experiments provide large datasets to train predictive models using machine learning approaches that can learn the relationship between regulatory DNA sequences and their associated molecular phenotypes. A Support Vector Machine (SVM) is a popular type of supervised classification model that learns an optimal linear separating hyperplane in a high-dimensional feature space. SVMs rely on a function

called a kernel that measures the similarity between all pairs of datapoints in the high-dimensional feature space. SVMs are appealing because they are stable to train and, when used with an appropriate kernel, can model complex input-output relationships. The gapped k -mer (gkm) string kernel [12, 7] was developed to enable training SVMs on string inputs such as DNA or protein sequences. The gkm kernel computes a similarity between pairs of sequences based on the biologically motivated notion of shared approximate occurrences of short subsequences allowing for gaps and mismatches. The gkm kernel can be further combined with other nonlinear kernels such as the radial basis function (RBF) kernel to capture complex nonlinear relationships between the input string features. Gapped k -mer SVMs and their extensions have been successfully applied to several prediction tasks in regulatory genomics such as TF binding and chromatin

A Appendices

A.1 AuROC and AuPRC for identifying TAL motif

Similar to Fig. 4, the auROC and auPRC curves for identifying the TAL motif using importance scores are shown in fig A.1.

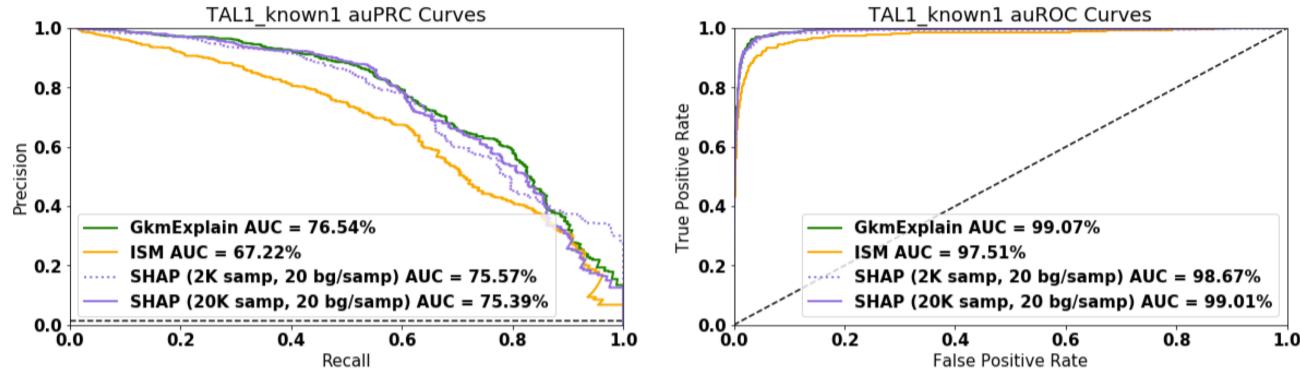


Fig. A.1: **GkmExplain outperforms ISM at identifying TAL1 motifs.** Analogous to Fig. 4, but using 16bp windows (the length of the TAL1_known1 motif) rather than 10bp windows.

$$\text{wgkmrbf_SVM_IG}_i^{S_x} = \int_{\beta=0}^{\beta=1} \sum_{j=1}^m \alpha_j y^j \gamma \frac{\partial K_{\text{wgkm}}^{(S_x, S_z)}(\eta^{S_x})}{\partial \pi_i(\beta)} \exp(\gamma \beta K_{\text{wgkm}}(S_x, S_z) - \gamma) d\beta$$

Recall that $\frac{\partial K_{\text{wgkm}}^{(S_x, S_z)}(\eta^{S_x})}{\partial \pi_i(\beta)}$ is constant with respect to $\pi_i(\beta)$. After performing the integral with respect to β , we obtain **Eqn. 13**, reproduced below for convenience:

$$\text{wgkmrbf_SVM_IG}_i^{S_x} = \phi_{i, S_x}^{\text{wgkmrbfsvm}} = \sum_{j=1}^m \alpha_j y^j \phi_{i, S_x}^{\text{wgkmrbf}} \quad (26)$$

A.3 Normalization of scores for TF-MoDISco

Empirically, we found that the following normalization of the importance scores produces improved results with TF-MoDISco. Let $f_h(S_x, i, B^*)$ be a function that returns the hypothetical importance of base B^* at position i in sequence S_x , and let S_x^i be the original base at position i in sequence S_x . To normalize the hypothetical importance scores at position i , we divide by the sum of all hypothetical scores at position i with the same sign as $f_h(S_x, i, S_x^i)$. The rationale is that if a different base at position i could produce a score of higher magnitude than S_x^i , then S_x^i is relatively less important. Let $1\{x > 0\}$ be an indicator function that returns 1 if x is positive and 0 otherwise. Formally, the normalized hypothetical importance for base B^* at position i is defined as:

$$\bar{f}_h(S_x, i, B^*) := \frac{f_h(S_x, i, B^*) f_h(S_x, i, S_x^i)}{\sum_{B'} f_h(S_x, i, B') 1\{(f_h(S_x, i, B') f_h(S_x, i, S_x^i)) > 0\}} \quad (27)$$

Similarly, we define the normalized importance score as the value of the normalized hypothetical importance score for the base that is actually present in the original sequence:

$$\bar{f}_i(S_x, i) := \bar{f}_h(S_x, i, S_x^i) \quad (28)$$

We find that the normalized importance scores appear to be less noisy relative to the unnormalized importance scores, as illustrated in Fig. A.2.

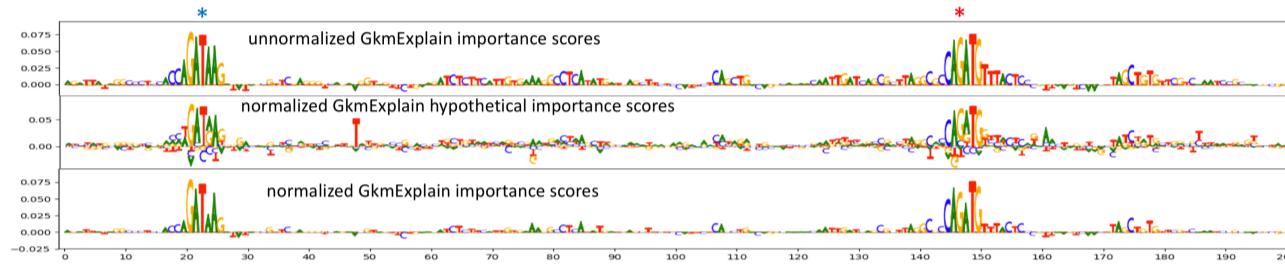


Fig. A.2: **Normalization of GkmExplain importance scores for TF-MoDISco.** A gkmrbf SVM was trained as described in Sec. 6.1. Shown are unnormalized GkmExplain importance scores (top row), normalized hypothetical importance scores (Eqn. 27), and normalized importance scores (Eqn. 28) on a single sequence. The blue star indicates the location of an embedded GATA1 motif (GATAAG), and the red star indicates the location of an embedded TAL1 motif (CAGATG). The normalized importance scores appear less noisy than the unnormalized importance scores.

A.4 AuROC on dsQTL data

Performance in terms of auROC on the dsQTL dataset is shown in **Fig. A.3**. Although GkmExplain gives a weaker auROC than ISM in 4 out of 5 cases, the difference is not statistically significant by a Wilcoxon test (possibly owing to the small number of samples). We note that when negatives greatly outnumber positives, auPRC is generally considered a more useful performance metric than auROC. We also note that GkmExplain consistently outperforms deltaSVM.

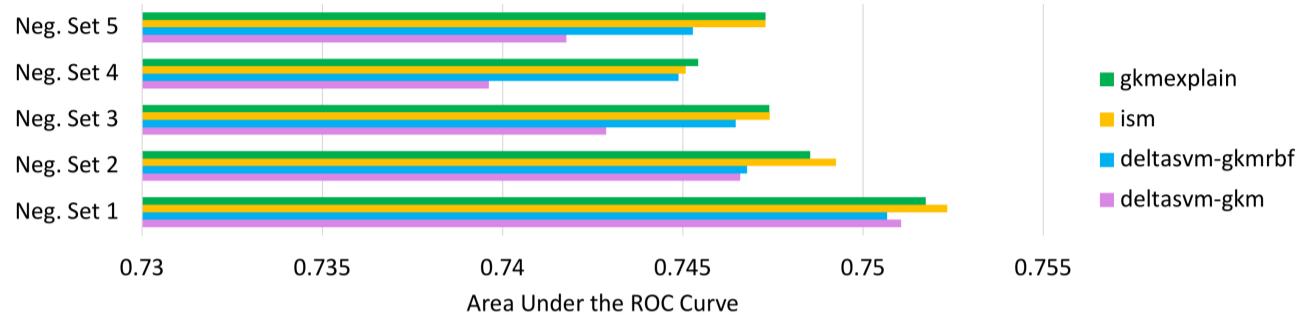


Fig. A.3: AuROC for dsQTL identification. Analogous to **Fig. 8**, but showing auROC instead of auPRC. The GkmExplain and ISM methods consistently outperform deltaSVM.

A.5 Integrated Gradients for interpreting other non-linear SVMs

To explore the general applicability of Integrated Gradients (IG) to SVM interpretation beyond genomics, we used the MADELON dataset [1]. This dataset has a training set of 2000 examples and a validation set of 600 examples, all with 500 features, of which 20 features are informative and the remaining are “distractor” features with no predictive power. Because the Gaussian SVM had a strong tendency to overfit to the full dataset, we created 10 reduced feature sets, each with the same set of examples in the training and validation sets, but with 20 likely informative features (identified based on the challenge submission results) and a different set of 20 distractor features in each set. Each of these 10 sets served as a distinct data point on which to evaluate the accuracy of the importance scoring algorithm in question (Integrated Gradients or SHAP). Each feature in each set was normalized by subtracting the average feature value and dividing by the standard deviation. We trained Gaussian kernel SVM on each of the datasets and achieved approximately 0.8 accuracy on the validation set. These classifiers were then provided to the importance scoring algorithms.

When applying Integrated Gradients, we used a reference that was the average feature value of the negative labeled training points and approximated the integration of the partial derivatives using 10 linearly-spaced intermediate points between each test point and the reference. With SHAP, we experimented with 10 and 50 samples around the test points using the same reference as was used for IG. We judged the quality of each importance scoring algorithm by feeding the top 10 most important features reported by the algorithm to a Gaussian kernel SVM and measuring the SVM’s accuracy on the validation set. The results are shown in Figure A.4. Integrated Gradients is significantly faster than SHAP and tends to be more accurate than when SHAP is used with a small number of perturbation samples. With a larger number of perturbation samples, SHAP is more accurate, though it is still slower. A notebook demonstrating how to use Integrated Gradients on a non-genomic dataset is at:

<https://colab.research.google.com/drive/1LUGMIwOHLKddK3deg7IIZ71kb2Nu1qv2>.

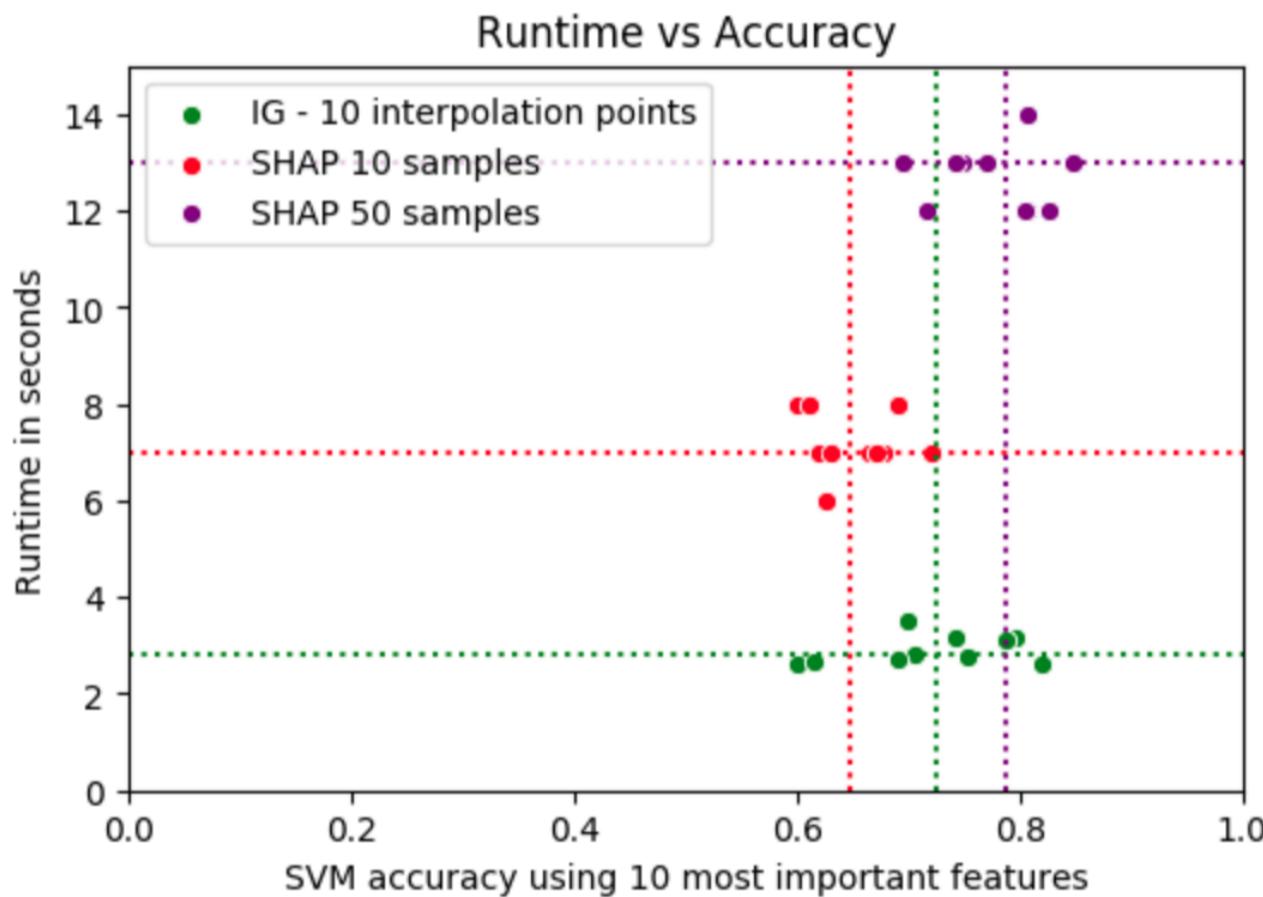


Fig. A.4: Comparison of Integrated Gradients with SHAP for Gaussian kernel SVM interpretation on the Madelon dataset. Dashed lines indicate means. IG appears to be pareto optimal in the sense that it can provide more accurate results faster than compared to SHAP used with a small number of perturbation samples.

References

- [1] Uci machine learning repository: Madelon data set. <https://archive.ics.uci.edu/ml/datasets/Madelon>. (Accessed on 08/01/2018).
- [2] www.beerlab.org/deltasvm/. <http://www.beerlab.org/deltasvm/>. (Accessed on 10/26/2018).

- [3]Timothy L Bailey, Mikael Boden, Fabian A Buske, Martin Frith, Charles E Grant, Luca Clementi, Jingyuan Ren, Wilfred W Li, and William S Noble. MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res.*, 37(Web Server issue):W202–8, July 2009.
- [4]Yana Bromberg and Burkhard Rost. Comprehensive in silico mutagenesis highlights functionally important residues in proteins. *Bioinformatics*, 24(16):i207–12, August 2008.
- [5]Jacob F Degner, Athma A Pai, Roger Pique-Regi, Jean-Baptiste Veyrieras, Daniel J Gaffney, Joseph K Pickrell, Sherryl De Leon, Katelyn Michelini, Noah Lewellen, Gregory E Crawford, Matthew Stephens, Yoav Gilad, and Jonathan K Pritchard. DNase I sensitivity QTLs are a major determinant of human expression variation. *Nature*, 482(7385):390–394, February 2012.
- [6]ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489(7414):57–74, September 2012.
- [7]Mahmoud Ghandi, Dongwon Lee, Morteza Mohammad-Noori, and Michael A. Beer. Enhanced regulatory sequence prediction using gapped k-mer features. *PLoS Comput. Biol.*, 10(7):e1003711, jul 2014.
- [8]Sven Heinz, Christopher Benner, Nathanael Spann, Eric Bertolino, Yin C Lin, Peter Laslo, Jason X Cheng, Cornelis Murre, Harinder Singh, and Christopher K Glass. Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol. Cell*, 38(4):576–589, May 2010.
- [9]Pouya Kheradpour and Manolis Kellis. Systematic discovery and characterization of regulatory motifs in ENCODE TF binding experiments. *Nucleic Acids Res.*, 42(5):2976–2987, March 2014.
- [10]Dongwon Lee. LS-GKM: a new gkm-SVM for large-scale datasets. *Bioinformatics*, 32(14):2196–2198, July 2016.
- [11]Dongwon Lee, David U Gorkin, Maggie Baker, Benjamin J Strober, Alessandro L Asoni, Andrew S McCallion, and Michael A Beer. A method to predict the impact of regulatory variants from DNA sequence. *Nat. Genet.*, 47(8):955–961, August 2015.
- [12]Christina Leslie and Rui Kuang. Fast String Kernels using Inexact Matching for Protein Sequences. *J. Mach. Learn. Res.*, 5(Nov):1435–1455, 2004.
- [13]Scott Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *NIPS*, 2017.
- [14]Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. *CoRR*, abs/1704.02685, 2017.
- [15]Avanti Shrikumar, Katherine Tian, Anna Shcherbina, Ágáta Avsec, Abhimanyu Banerjee, Mahfuz Sharmin, Surag Nair, and Anshul Kundaje. Tf-modisco v0.4.2.2-alpha: Technical note. *CoRR*, abs/1811.00416.
- [16]Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. *CoRR*, abs/1703.01365, 2017.
- [17]Jian Zhou and Olga G Troyanskaya. Predicting effects of noncoding variants with deep learning-based sequence model. *Nat. Methods*, 12(10):931–934, August 2015.