Large Scale Learning

TSS recognition

Discussion

Large Scale Machine Learning for Genomic Sequence Analysis (Support Vector Machine Based Signal Detectors)

> Sören Sonnenburg Friedrich Miescher Laboratory, Tübingen

joint work with Alexander Zien, Jonas Behr, Gabriele Schweikert, Petra Philips and Gunnar Rätsch



Introd	

Large Scale Learning

TSS recognition

Outline



- 2 Large Scale Learning
- 3 TSS recognition



 Introduction
 Large Scale Learning
 TSS recognition
 Discussion

 •ocococ
 cocococ
 cocococ
 cocococ

 Genomic Signals
 Recognizing Genomic Signals
 cococc

Discriminate true signal positions against all other positions

 $\approx\!150$ nucleotides window around dimer

CT...GTCGTA...GAAGCTAGGAGCGC...ACGCGT...GA

- True sites: fixed window around a true site
- Decoy sites: all other consensus sites

AAACAAATAAGTAACTAATCTTTTA<mark>GGAAGAACGTTTCAACCA</mark>TTTTGAG AAGATTAAAAAAAAACAAATTTTTA<mark>GCATTACAGATATAATAATCTAATT</mark> CACTCCCCCAAATCAACGATATTTTA<mark>GTTCACTAACACACTCCGTCTG</mark>TGCC TTAATTTCACTTCCACATACTTCCAGATCATCAACCAAAACCAACAC

Examples: Transcription start site finding, splice site prediction, alternative splicing prediction, trans-splicing, polyA signal detection, translation initiation site detection

Large Scale Learning

TSS recognition

Discussion

Genomic Signals

Types of Signal Detection Problems I

Vague categorization

(based on positional variability of motifs)

Position Independent

 \rightarrow Motifs may occur anywhere,

AAACAAAAACCGTAACTAATCTTTTAGAGAGAACGTTTCAACCATTTTGAG AAGATTAACTCATCACAGATTTCATTACATACAGATATAATTCAAAAATT CACTCCCCCAAATCAACGATATTTAAAAATCACCACACCCGTCTGTGC

e.g. tissue classification using promotor region



Large Scale Learning

TSS recognition

Genomic Signals

Types of Signal Detection Problems II

Vague categorization

(based on positional variability of motifs)

Position Dependent

 \rightarrow Motifs very stiff, almost always at same position,

e.g. Splice Site Classification



Large Scale Learning

TSS recognition

Discussion

Genomic Signals

Types of Signal Detection Problems III

Vague categorization

(based on positional variability of motifs)

Mixture Position Dependent/Independent

 \rightarrow variable but still positional information

e.g. Promoter Classification



Large Scale Learning

TSS recognition

Discussion

Support Vector Machines

Classification - Learning based on examples

Given:

Training examples $(\mathbf{x}_i, y_i)_{i=1}^N \in (\{A, C, G, T\}^L, \{-1, +1\})^N$

AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG AAGATTAAAAAAAAACAAATTTTTAGCATTACAGATATAATAATCTAATT CACTCCCCAAATCAACGATATTTTAGTTCACTAACACATCCGTCTGTGCC TTAATTTCACTTCCACATACTTCCAGATCATCAATCTCCAAAACCAACAC TTGTTTTAATATTCAATTTTTTACAGTAAGTTGCCAATTCAATGTTCCAC TACCTAATTATGAAATTAAAATTCAGTGTGCTGATGGAAACGGAGAAGTC

Wanted:

Function (Classifier) $f(\mathbf{x}) : \{A, C, G, T\}^L \mapsto \{-1, +1\}$



 Introduction
 Large Scale Learning
 TSS recognition
 Discussion

 00000●
 00000
 0000000
 0000000

 Support Vector Machines
 Vector Machines
 Vector

Support Vector Machines (SVMs)



 Support Vector Machines learn weights α ∈ ℝ^N over training examples in kernel feature space Φ : x → ℝ^D,

$$f(\mathbf{x}) = \operatorname{sign}\left(\sum_{i=1}^{N} y_i \alpha_i \mathbf{k}(\mathbf{x}, \mathbf{x}_i) + b\right),$$

with kernel $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$



Introduction 000000	Large Scale Learning ●0000	TSS recognition	Discussion
String Kernels			
The Spectr	rum Kernel		

Support Vector Machine

$$f(\mathbf{x}) = \operatorname{sign}\left(\sum_{i=1}^{N} y_i \alpha_i \mathbf{k}(\mathbf{x}, \mathbf{x}_i) + b\right),$$

Spectrum Kernel (with mismatches, gaps)

$$K(\mathbf{x}, \mathbf{x}') = \Phi_{sp}(\mathbf{x}) \cdot \Phi_{sp}(\mathbf{x}')$$

AAACAAAAACGTAACTAATCTTTTAGAGAGAACGTTTCAACCATTTTGAG AAGATTAACTCATCACAGATTTCATTACATACAGATATAATTCAAAAAATT CACTCCCCCAAATCAACGATATTTAAAAATCACTAACACATCCGTCTGTGC



Large Scale Learning

TSS recognition

String Kernels

The Weighted Degree Kernel

Support Vector Machine

$$f(\mathbf{x}) = \operatorname{sign}\left(\sum_{i=1}^{N} y_i \alpha_i \mathbf{k}(\mathbf{x}, \mathbf{x}_i) + b\right),$$

$$\mathbf{k}(\mathbf{x},\mathbf{x}') = \sum_{k=1}^{K} \beta_k \sum_{i=1}^{L-k+1} \mathbb{I}\left\{\mathbf{x}[i]^k = \mathbf{x}'[i]^k\right\}.$$

Example: K = 3: $\mathbf{k}(\mathbf{x}, \mathbf{x}') = \beta_1 \cdot 21 + \beta_2 \cdot 8 + \beta_3 \cdot 3$



Large Scale Learning

TSS recognition

Discussion

String Kernels

The Weighted Degree Kernel with *shifts*

Support Vector Machine

$$f(\mathbf{x}) = \operatorname{sign}\left(\sum_{i=1}^{N} y_i \alpha_i \mathbf{k}(\mathbf{x}, \mathbf{x}_i) + b\right),$$



Large Scale Learning

TSS recognition

Fast SVM Training and Evaluation

Accelerating String-Kernel-SVMs

- Linear run-time of the kernel
- Accelerating linear combinations of kernels

Idea of the Linadd Algorithm:

Store w and compute $\mathbf{w} \cdot \Phi(\mathbf{x})$ efficiently

$$f(\mathbf{x}_j) = \sum_{i=1}^{N} \alpha_i y_i \, \mathsf{k}(\mathbf{x}_i, \mathbf{x}_j) = \underbrace{\sum_{i=1}^{N} \alpha_i y_i \Phi(\mathbf{x}_i)}_{\mathbf{w}} \cdot \Phi(\mathbf{x}_j) = \mathbf{w} \cdot \Phi(\mathbf{x}_j)$$

Possible for low-dimensional or sparse ${f w}$

Effort: $\mathcal{O}(NL) \Rightarrow$ **speedup of factor** N \Rightarrow Training on millions of examples, evaluation on billions.

Laboratory

Large Scale Learning

TSS recognition

Fast SVM Training and Evaluation

Accelerating String-Kernel-SVMs II

Recent work:

Further drastic speedup using advances of primal SVMs solvers

Acceleration using fast primal SVMs

- Idea: Train SVM in primal using kernel feature space
- Problem: > 12 million dims; 50 million examples
- Only $\mathbf{w} \leftarrow \mathbf{w} + \alpha \Phi(\mathbf{x})$ and $\mathbf{w} \cdot \Phi(\mathbf{x})$ required.
- Compute $\Phi(\mathbf{x})$ on-the-fly and parallelize!

Results

- \Rightarrow Computations are simple "table lookups" of *k*-mers weights
- \Rightarrow Allows training on 50 million examples



Detecting Transcription Start Sites



- POL II indirectly binds to a rather vague region of $\approx [-20,+20]~{\rm bp}$
- Upstream of TSS: promoter containing transcription factor binding sites
- Downstream of TSS: 5' UTR, and further downstream coding regions and introns (different statistics)
- 3D structure of the promoter must allow the transcription factors to bind

Several weak features \Rightarrow Promoter prediction is non-trivial

Large Scale Learning

TSS recognition

Incorporating Prior Knowledge

Features to describe the TSS

- TFBS in Promoter region
- condition: DNA should not be too twisted
- CpG islands (often over TSS/first exon; in most, but not all promoters)
- TSS with TATA box (pprox -30 bp upstream)
- Exon content in UTR 5" region
- Distance to first donor splice site

Idea:

Combine weak features to build strong promoter predictor

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') \!=\! k_{TSS}(\mathbf{x}, \mathbf{x}') \!+\! k_{CpG}(\mathbf{x}, \mathbf{x}') \!+\! k_{coding}(\mathbf{x}, \mathbf{x}') \!+\! k_{energy}(\mathbf{x}, \mathbf{x}') \!+\! k_{twist}(\mathbf{x}, \mathbf{x}')$$

000000	Large Scale Learning	OOOOOO	Discussion
Incorporating Prior Knowl	edge		
The 5 sub-k	ernels		

- Signal (including parts of core promoter with TATA box)
 - use Weighted Degree Shift kernel
- Opg Islands, distant enhancers and TFBS upstream of TSS
 - use **Spectrum kernel** (large window upstream of TSS)
- Model coding sequence TFBS downstream of TSS
 - use another **Spectrum kernel** (small window downstream of TSS)
- Stacking energy of DNA
 - use *btwist* energy of dinucleotides with Linear kernel
- Twistedness of DNA
 - use btwist angle of dinucleotides with Linear kernel



Large Scale Learning

TSS recognition

Laboratory

Results

State-of-the-art Performance

Receiver Operator Characteristic Curve and Precision Recall Curve



 \Rightarrow 35% true positives at a false positive rate of 1/1000 (best other method find about a half (18%))

Introduction 000000	Large Scale Learning 00000	TSS recognition ○○○○○○○	
General			
Beauty in Ge	nerality		



- Transcription Start (Sonnenburg et al., Eponine Down et al.)
- Acceptor Splice Site (Schweikert et al.)
- Donor Splice Site (Schweikert et al.)
- Alternative Splicing (Rätsch et al., -)
- Transsplicing (Schweikert et al., -)
- Translation Initiation (Sonnenburg et al., Saeys et al.)

Large Scale Learning

TSS recognition ○○○○○●○ Discussion

Positional Oligomer Importance Matrices

Positional Oligomer Importance Matrices (POIMs)

Determine importance of k-mers at one glance:

• Given k-mer **z** at position j in the sequence, compute expected score $\mathbb{E}[s(\mathbf{x}) | \mathbf{x}[j] = \mathbf{z}]$ (for small \mathbf{k})

• Normalize with expected score over all sequences

POIMs

$$Q(\mathbf{z}, j) := \mathbb{E}[s(\mathbf{x}) | \mathbf{x}[j] = \mathbf{z}] - \mathbb{E}[s(\mathbf{x})]$$



Large Scale Learning

TSS recognition

Interpretable

Interpretable via Positional Oligomer Importance Matrices

Example: Drosophila Transcription Starts



Differential POIM Overview - Drosophila TSS

Introduction	

TSS recognition

Conclusions

Support Vector Machines with string kernels

- General
- Fast: Applicable to genome-sized datasets
- Often are state-of-the art signal detectors
 - TSS
 - Acceptor and Acceptor Splice Site
 - ...
- Used in mGene gene finder http://www.mgene.org
- Positional Oligomer Importance Matrices help making SVMs interpretable

Galaxy web-interface http://galaxy.fml.tuebingen.mpg.de Efficient implementation http://www.shogun-toolbox.org More machine learning software http://mloss.org