Introduction	Sequence Classification	Support Vector

Applications Explanation and Visualization Summar

Genomic Signal Detection ... using Support Vector Machines

Machines

Sören Sonnenburg TU Berlin

joint work with Alexander Zien, Jonas Behr, Gabriele Schweikert, Konrad Rieck, Petra Philips, Gunnar Rätsch, Vojtech Franc



Introduction Sequence Classification Support Vector Machines Applications Explanation and Visualization

Outline



- 2 Sequence Classification
- Support Vector Machines

4 Applications

5 Explanation and Visualization





Applications Explanation and Visualization Summ

Outline



Genomic Signals

- Sequence Classification
 Support Vector Machines
- Support Vector Machines
- 4 Applications
 - Splice Site Recognition
 - TSS Recognition
 - Gene Finding mGene
 - Aligning Short Reads QPALMA
- 5 Explanation and Visualization
 - Introduction
 - Definition
 - Applications





equence Classification

Support Vector Machines

Applications Explanation and Visualization Su

Genome



Genomic Signals



Genomic Signal Detection

- Start/Stop of Genes
- Donor Splice Site (Exon-Intron-Boundary)
- Acceptor Splice Site (Intron-Exon-Boundary)





Discriminate true signal positions against all other positions



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

AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG AAGATTAAAAAAAACAAATTTTTAGCATTACAGATATAATAATCTAATT CACTCCCCAAATCAACGATATTTTAGTTCACTAACACATCCGTCTGTGCC TTAATTTCCACTTCCACATACTTCCAGATCATCCAAAACCAACAC TTGTTTTAATATTCAATTTTTACAGTAAGTTGCCAATTCAATGTTCCAC TACTAATTATGAAATTAAAATTCAGTGTGCCGATTGGAAACGGAGAAGTC

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



Sequence Classification Introduction 00000

Support Vector Machines

Applications Explanation and Visualization

Genomic Signals

Types of Signal Detection Problems I

Vague categorization

(based on positional variability of motifs)

Position Independent

 \rightarrow Motifs may occur anywhere,

AAACAAAAACGTAACTAATCTTTTAGAGAGAACGTTTCAACCATTTTGAG AAGATTAACTCATCACAGATTTCATTACATACAGATATAATTCAAAAATT CACTCCCCAAATCAACGATATTTAAAAATCACTAACACATCCGTCTGTGC

e.g. tissue classification using promotor region



Sequence Classification Introduction 00000

Support Vector Machines

Applications Explanation and Visualization

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,

AAACAAATAAGTAACTAATCTTTTAAGAAGAACGTTTCAACCATTTTGAG AAGATTAAAAAAAAAACAAATTTTT<mark>AA</mark>CATTACAGATATAATAATCTAATT CACTCCCCAAATCAACGATATTTTAATTCACTAACACATCCGTCTGTGCC

e.g. Splice Site Classification



sification Support 00000

Support Vector Machines

Applications Explanation and Visualization Sur

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



Sequence Classification

Support Vector Machines

Applications Explanation and Visualization Summ

Outline

- Introduction
 Genomic Signals
- Sequence Classification
 Support Vector Machines
 - 3 Support Vector Machines
- 4 Applications
 - Splice Site Recognition
 - TSS Recognition
 - Gene Finding mGene
 - Aligning Short Reads QPALMA
- 5 Explanation and Visualization
 - Introduction
 - Definition
 - Applications
- 6 Summary



Classification - Learning based on examples I

Given:

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

(pprox 1 billion neg. sequences; < 200.000 positive sequences)

Wanted:

(

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

 $\approx\!150$ nucleotides window around dimer

Aim: Accurate signal prediction for the whole genome

Introduction Sequence Classification Support Vector Machines Applications Explanation and Visualization Summa

Classification - Learning based on examples I

Given:

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

(≈ 1 billion neg. sequences; < 200.000 positive sequences)

Wanted:

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

 \approx 150 nucleotides window around dimer CT...GTCGTA...GAAGCTAGGAGCGC...ACGCGT...GA

Aim: Accurate signal prediction for the whole genome

Classification - Learning based on examples II

AAACAAATAAGTAACTAATCTTTT<mark>AG</mark>GAAGAACGTTTCAACCATTTTGAG





tion Support Vector Machines

Applications Explanation and Visualization S

Summary

Support Vector Machines

Classification - Learning based on examples III





Classification - Learning based on examples IV





Sequence Classification

Support Vector Machines

Applications Explanation and Visualization Summ

Outline

- Introduction
 - Genomic Signals
- 2 Sequence Classification
 - Support Vector Machines
- 3 Support Vector Machines
 - Applications
 - Splice Site Recognition
 - TSS Recognition
 - Gene Finding mGene
 - Aligning Short Reads QPALMA
- 5 Explanation and Visualization
 - Introduction
 - Definition
 - Applications
- 6 Summary



Sequence Classification

Support Vector Machines

Applications Explanation and Visualization Se

Support Vector Machines (SVMs) I





Sequence Classification

Support Vector Machines

Applications Explanation and Visualization S

Support Vector Machines (SVMs) II





Sequence Classification

Support Vector Machines

Applications Explanation and Visualization Se

Support Vector Machines (SVMs) III





Introduction Seque

ience Classification Si

Support Vector Machines

Applications Explanation and Visualization Sum

SVMs and Kernels





Sequence Classification Introduction

Support Vector Machines

Applications Explanation and Visualization

Support Vector Machines (SVMs)



• Support Vector Machines learn weights $\boldsymbol{\alpha} \in \mathbb{R}^N$ over training examples in kernel feature space $\Phi : \mathbf{x} \mapsto \mathbb{R}^D$,

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

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



Support Vector Machines

String Kernels

The Spectrum Kernel (Leslie et al. 2002)

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)

$$\mathcal{K}(\boldsymbol{x},\boldsymbol{x}') = \Phi_{\textit{sp}}(\boldsymbol{x}) \cdot \Phi_{\textit{sp}}(\boldsymbol{x}')$$

Example k = 3:

- x AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG
- x' TACCTAATTATGAAATTAAATTTCAGTGTGCTGATGGAAACGGAGAAGTC

3-mer	AAA	AAC	 CCA	CCC	 TTT
# in x	2	4	 1	0	 3
# in x ′	3	1	 0	0	 1



$$\mathbf{k}(\mathbf{x},\mathbf{x}') = 2 \cdot 3 + 4 \cdot 1 + \ldots 1 \cdot 0 + 0 \cdot 0 \ldots 3 \cdot 1$$

Support Vector Machines

Applications Explanation and Visualization

Summary

String Kernels

The Weighted Degree Kernel (Sonnenburg et al. 2005)

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



Introduction Sequence Classification Support Vector Machines Applications Explanation and Visualization Summary oceoco String Kernels The Weighted Degree Kernel with shifts (Raetsch, Sonnenburg et al. 2005)

Support Vector Machine

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



Support Vector Machines

Large Scale Learning

Accelerating String-Kernel-SVMs

Aim: Train and apply string-kernel SVM on all available data

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

Idea of the Linadd Algorithm (Sonnenburg et al., 2005): Store w and compute $w \cdot \Phi(x)$ efficiently

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

Possible for low-dimensional or sparse **Effort:** $\mathcal{O}(ML) \Rightarrow$ **speedup of factor** N_s (with $L := dim(\mathcal{X})$)



Support Vector Machines

Applications Explanation and Visualization S

Large Scale Learning

General Recipe — Howto

Collect Data

What features can be used to describe "my Signal"?

- Are these *position independent* k-mers? ⇒ Spectrum Kernel?
- Are these strongly position dependent k-mers? ⇒ WD Kernel?
- \bullet Are these partially position dependent k-mers? \Rightarrow WDK with shifts

Train SVM

- split data into training, validation, test
- model selection over hyperparameters C, k, s
- learn final model on train+validation, estimate performance on test



Support Vector Machines

Applications Explanation and Visualization Su

Large Scale Learning

General Recipe — Howto

Collect Data

What features can be used to describe "my Signal"?

- Are these *position independent* k-mers? \Rightarrow Spectrum Kernel?
- Are these *strongly position dependent* k-mers? ⇒ WD Kernel?
- Are these partially position dependent k-mers? \Rightarrow WDK with shifts

Train SVM

- split data into training, validation, test
- model selection over hyperparameters C, k, s
- learn final model on train+validation, estimate performance on test



Support Vector Machines

Applications Explanation and Visualization Su

Large Scale Learning

General Recipe — Howto

Collect Data

What features can be used to describe "my Signal"?

- Are these *position independent* k-mers? \Rightarrow Spectrum Kernel?
- Are these strongly position dependent k-mers? \Rightarrow WD Kernel?
- Are these partially position dependent k-mers? \Rightarrow WDK with shifts

Train SVM

- split data into training, validation, test
- model selection over hyperparameters C, k, s
- learn final model on train+validation, estimate performance on test



Sequence Classification

Support Vector Machines

Applications Explanation and Visualization Sum

Outline

- Introduction
 - Genomic Signals
- 2 Sequence Classification
 - Support Vector Machines
 - Support Vector Machines
- 4 Applications
 - Splice Site Recognition
 - TSS Recognition
 - Gene Finding mGene
 - Aligning Short Reads QPALMA
- 5 Explanation and Visualization
 - Introduction
 - Definition
 - Applications
- 6 Summary





Discriminate true signal positions against all other positions



• True sites: fixed window around a true splice site

• Decoy sites: all other consensus sites

AAACAAATAAGTAACTAATCTTTTA<mark>GGAAGAACGTTT</mark>CAACCATTTTGAG AAGATTAAAAAAAAACAAATTTTTA<mark>GCATTACAGATATAATAATCTAATT</mark> CACTCCCCAAATCAACGATATTTTA<mark>GTTCACTAACACATCCGTCTG</mark>TGCC TTAATTTCACTTCCACATACTTCCAG<mark>ATCATCAATCTCCAAAACCAACAC</mark>

- Create training sample from cDNA/EST alignments to genome (for human, e.g., 50 million examples)
- Sequences are compared via Weighted Degree Kernel
- Train SVM on up to 8 million examples.





Discriminate true signal positions against all other positions



• True sites: fixed window around a true splice site

• Decoy sites: all other consensus sites

AAACAAATAAGTAACTAATCTTTTA<mark>GGAAGAACGTTTCAACCA</mark>TTTTGAG AAGATTAAAAAAAAACAAATTTTTA<mark>GCATTACAGATATAATAATCTAATT</mark> CACTCCCCAAATCAACGATATTTTA<mark>GTTCACTAACACATCCGTCTGTGCC</mark> TTAATTTCACTTCCACATACTTCCA<mark>GATCATCAACACCAAAACCCAACAC</mark>

- Create training sample from cDNA/EST alignments to genome (for human, e.g., 50 million examples)
- Sequences are compared via Weighted Degree Kernel
- Train SVM on up to 8 million examples.



Introduction Sequence Classification Supposed 0000 000

Support Vector Machines

Applications Explanation and Visualization S

Splice Site Recognition

Splice Site Recognition - Results

- $\bullet\,$ Human splice sites: $5\cdot 10^7$ strings of length ≈ 141
- Note: Raw data is already 7GB in size



SVM \approx 3 times more accurate than IMCs (54.4% vs. 16.2% auPRC; Sonnenburg et al. 2007)



Introduction Sequence Classification Sequence Classification

Support Vector Machines

Applications Explanation and Visualization Su

Splice Site Recognition

Splice Site Recognition - Results

- $\bullet\,$ Human splice sites: $5\cdot 10^7$ strings of length ≈ 141
- Note: Raw data is already 7GB in size



Support Vector Machines

Applications Explanation and Visualization S

Splice Site Recognition

Detecting Transcription Start Sites



- $\bullet\,$ POL II binds to a rather vague region of $\approx [-20,+20]$ 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
- \Rightarrow Promoter Prediction is non-trivial



Support Vector Machines

Applications Explanation and Visualization Sur

TSS Recognition

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)
- TFBS in Promoter region, Exon content in UTR 5" region
- Distance to first donor splice site

Idea:

Combine weak features to build strong promoter predictor

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



- TSS 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


Support Vector Machines

Applications Explanation and Visualization Su

TSS Recognition

Training – Data Generation

Training and Validation Data (50% : 50%)

True TSS (8508 positive)

• From dbTSSv4 (based on hg16) extract putative TSS windows of size [-1000, +1000]

Decoy TSS (85042)

- Annotate dbTSSv4 with transcription-stop (via *BLAT* alignment of mRNAs)
- From the interior of the gene (+100*bp* to gene end) sample negatives for training (10 per positive)

Fair genome-wide evaluation

- Compare against FirstEF, McPromoter, Eponine
- Only consider "new" TSS from dbTSSv5-dbTSSv4, with max 30% overlap



Support Vector Machines

Applications Explanation and Visualization Sum

TSS Recognition

State-of-the-art Performance

Receiver Operator Characteristic and Precision Recall Curve



ARTS (Sonnenburg et al. 2006) twice as accurate!

Independent evaluation of 17 methods (Abeel et al. ISMB, 2009) **TSS detector (ARTS) winner in evaluation of 17 methods.**



Support Vector Machines

Applications Explanation and Visualization Sum

TSS Recognition

State-of-the-art Performance

Receiver Operator Characteristic and Precision Recall Curve



ARTS (Sonnenburg et al. 2006) twice as accurate!

Independent evaluation of 17 methods (Abeel et al. ISMB, 2009) **TSS detector (ARTS) winner in evaluation of 17 methods.**

Support Vector Machines

Applications Explanation and Visualization S

TSS Recognition

Beauty in Generality



• Transcription Start (Sonnenburg et al., 2006/Down et al. 2002)

- Acceptor Splice Site (Sonnenburg et al., 2007/Baten et al. 2006)
- Donor Splice Site (Sonnenburg et al., 2007/Baten et al. 2006)
- Alternative Splicing (Rätsch, Sonnenburg et al., 2005/-)
- Transsplicing (Schweikert, Sonnenburg et al., 2009/-)
- Translation Initiation (Sonnenburg et al., 2008/Saeys et al., 2007)

SVM with string kernel often most accurate method.



Support Vector Machines

Applications Explanation and Visualization Sun

Gene Finding - mGene

Individual Signal Predictions in UCSC Browser



(Raetsch, Sonnenburg et al. 2007, Schweikert et al. 2009)

Support Vector Machines

Applications Explanation and Visualization Sun

Gene Finding - mGene

Individual Signal Predictions in UCSC Browser



(Raetsch, Sonnenburg et al. 2007, Schweikert et al. 2009)

Support Vector Machines

Applications Explanation and Visualization Summary

Gene Finding - mGene

Individual Signal Predictions in UCSC Browser



(Raetsch, Sonnenburg et al. 2007, Schweikert et al. 2009

Introduction Sequence Classification Support Vector Machines Ococo Company Control Con

Results using mGene (Schweikert et al. 2009)

- Most accurate *ab initio* method in the nGASP genome annotation challenge for *C. elegans*
- Validation of gene predictions for *C. elegans*:

	No. of genes	No. of genes	Frac. of genes
New genes	2,197		pprox 42%
Missing unconf. genes		24	

• Annotation of other nematode genomes:

	No. of	No. exons/gene		
size [Mbp]				
235.94			96.6%	
	20121		93.3%	
	41129	5.4	93.1%	
			87.0%	



 Introduction
 Sequence Classification
 Support Vector Machines
 Applications
 Explanation and Visualization
 Summary

 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000</td

Results using mGene (Schweikert et al. 2009)

- Most accurate *ab initio* method in the nGASP genome annotation challenge for *C. elegans*
- Validation of gene predictions for *C. elegans*:

	No. of genes	No. of genes	Frac. of genes
		analyzed	w/ expression
New genes	2,197	57	pprox 42%
Missing unconf. genes	205	24	pprox 8%

• Annotation of other nematode genomes:

	No. of	No. exons/gene		
size [Mbp]				
235.94			96.6%	
	20121		93.3%	
	41129	5.4	93.1%	
			87.0%	



 Introduction
 Sequence Classification
 Support Vector Machines
 Applications
 Explanation and Visualization
 Summary

 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000</td

Results using mGene (Schweikert et al. 2009)

- Most accurate *ab initio* method in the nGASP genome annotation challenge for *C. elegans*
- Validation of gene predictions for *C. elegans*:

	No. of genes	No. of genes	Frac. of genes
		analyzed	w/ expression
New genes	2,197	57	pprox 42%
Missing unconf. genes	205	24	pprox 8%

• Annotation of other nematode genomes:

Genome	Genome	No. of	No. exons/gene	mGene	best other
	size [Mbp]	genes	(mean)	accuracy	accuracy
C. remanei	235.94	31503	5.7	96.6%	93.8%
C. japonica	266.90	20121	5.3	93.3%	88.7%
C. brenneri	453.09	41129	5.4	93.1%	87.8%
C. briggsae	108.48	22542	6.0	87.0%	82.0%



 Introduction
 Sequence Classification
 Support Vector Machines
 Applications
 Explanation and Visualization
 Summary

 0000
 000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 00000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000</td

Results using mGene (Schweikert et al. 2009)

- Most accurate *ab initio* method in the nGASP genome annotation challenge for *C. elegans*
- Validation of gene predictions for *C. elegans*:

	No. of genes	No. of genes	Frac. of genes
		analyzed	w/ expression
New genes	2,197	57	pprox 42%
Missing unconf. genes	205	24	pprox 8%

• Annotation of other nematode genomes:

Genome	Genome	No. of	No. exons/gene	mGene	best other
	size [Mbp]	genes	(mean)	accuracy	accuracy
C. remanei	235.94	31503	5.7	96.6%	93.8%
C. japonica	266.90	20121	5.3	93.3%	88.7%
C. brenneri	453.09	41129	5.4	93.1%	87.8%
C. briggsae	108.48	22542	6.0	87.0%	82.0%



Support Vector Machines

Applications Explanation and Visualization Su

Aligning Short Reads - QPALMA

mRNA Deep Sequencing





Support Vector Machines

Applications Explanation and Visualization Su

Aligning Short Reads - QPALMA

RNA-Seq Read Alignment



Support Vector Machines

Applications Explanation and Visualization Sum

Aligning Short Reads - QPALMA

RNA-Seq Read Alignment

... GCAAACCAGTGACCTGACTACGTCGTCACGTACGTACACGGTAGCT....CCGTAGAATTGACTGTGTTG... GCAAACCAGTGACCTGACTACGTCGTCGTAACGTAC



Support Vector Machines

Applications Explanation and Visualization Sun

Aligning Short Reads - QPALMA

RNA-Seq Read Alignment

... GCAAACCAGTGACCTGACTACCTACGTCGTAACGTACACGGTAGCT....CCGTAGAATTGACTGTGTTG... GCAAACCAGTGACCTGACTACGTCGTAACGTAC CAAACCAGTGACCTGACTACGTCGTAACGTACA AAACCAGTGACCTGACTACTACGTCGTAACGTACAC AACCAGTGACCTGACTACTACGTCGTAACGTACACG





Support Vector Machines

Applications Explanation and Visualization Su

Summar

Aligning Short Reads - QPALMA

RNA-Seq Read Alignment



 Introduction
 Sequence Classification
 Support Vector Machines
 Applications
 Explanation and Visualization
 Summary

 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000</td

Extended Smith-Waterman Algorithm (De Bona et al. 2008)

 \Rightarrow Combine heterogeneous evidence for accurate alignment



Classical scoring $f: \Sigma \times \Sigma \rightarrow \mathbb{R}$



				•		
	gap	А	С	G	Т	Ν
gap	0.33	0.3	0.12	0.3	0.3	0.55
А	0,31	0,12	0,12	0,3	0,55	0.33
С	0.44	0.12	0.44	0.3	0.59	0.12
G	0,13	0,85	0,31	0,33	0.51	0.3
Т	0,55	0,12	013	0,12	0,11	0.1
Ν	0.12	0.01	0.3	0.12	0.3	0.01
	^{gap} A C G T N	gap gap gap A O.33 C GA G.44 G.35 T N O.12	gap A gap 0.33 0.3 A 0.34 0.12 C 0.44 0.12 G 0.13 0.35 T 0.55 0.12 N 0.12 0.13	gap A C gap 0.33 0.33 0.12 A 0.31 0.12 0.12 C 0.44 0.12 0.44 G 0.13 0.85 0.41 T 0.55 0.12 0.13 N 0.15 0.12 0.33	gap A C G gap 0.33 0.31 0.12 0.33 A 0.31 0.12 0.12 0.33 C 0.44 0.12 0.44 0.33 G 0.43 0.45 0.44 0.33 G 0.43 0.45 0.43 0.43 T 0.55 0.42 0.43 0.42 N 0.42 0.42 0.43 0.43	gap A C G T gap 0.33 0.3 0.12 0.3 0.3 A 0.33 0.12 0.12 0.3 0.35 A 0.34 0.12 0.44 0.3 0.59 G 0.44 0.12 0.44 0.3 0.59 G 0.13 0.12 0.44 0.3 0.59 G 0.13 0.15 0.14 0.13 0.11 0.11 N 0.12 0.12 0.13 0.13 0.13 0.13

Source of Information

- Sequence matches
- Computational splice site predictions
- Intron length model
- Read quality information

Classical scoring $f: \Sigma \times \Sigma \to \mathbb{R}$

 Introduction
 Sequence Classification
 Support Vector Machines
 Applications
 Explanation and Visualization
 Summary

 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000</td

Extended Smith-Waterman Algorithm (De Bona et al. 2008)

 \Rightarrow Combine heterogeneous evidence for accurate alignment



Classical scoring $f: \Sigma \times \Sigma \to \mathbb{R}$

Extended Smith-Waterman Algorithm (De Bona et al. 2008)

\Rightarrow Combine heterogeneous evidence for accurate alignment



(De Bona et al. 2008)

Quality scoring $f : (\Sigma \times \mathbb{R}) \times \Sigma \to \mathbb{R}$

Introduction Sequence Classification Support Vector Machines

Applications Explanation and Visualization

Summary

Aligning Short Reads - QPALMA

QPalma's Accurate Alignments (De Bona et al. 2008)

Generate set of artificially spliced reads

- Genomic reads with quality information
- Genome annotation for artificially splicing the reads
- Use 10,000 reads for training and 30,000 for testing





Introduction 00000 Sequence Classification 0000

Support Vector Machines

Applications Explanation and Visualization Su

Outline

- Introduction
 - Genomic Signals
- Sequence Classification
 Support Vector Machines
- **3** Support Vector Machines
- 4 Applications
 - Splice Site Recognition
 - TSS Recognition
 - Gene Finding mGene
 - Aligning Short Reads QPALMA
- 5 Explanation and Visualization
 - Introduction
 - Definition
 - Applications





Support Vector Machines

Applications Explanation and Visualization

Summary

Introduction

Understanding Support Vector Machines

Goal

For PWMs we have sequence logos:



We would like to have similar means to understand Support Vector Machines.



SVM decision function is α weighting of training points

$$s(\mathbf{x}) = \sum_{i=1}^{N} lpha_i y_i \, \mathsf{k}(\mathbf{x}_i, \mathbf{x}) + b$$

 $\alpha_1 \cdot \textbf{AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG}$

 α_2 · AAGATTAAAAAAAAAAAAAAAAATTTTTAGCATTACAGATATAATAATAATCTAATT

 $\alpha_{3} \cdot \textbf{CACTCCCCAAATCAACGATATTTTAGTTCACTAACACATCCGTCTGTGCC}$

 α_N · TTAATTTCACTTCCACATACTTCCAGATCATCAATCTCCAAAACCAACAC

But we are interested in weights over features.



Introduction

equence Classification

Support Vector Machines

Applications Explanation and Visualization

Summary

Introduction

SVM Scoring Function

$$\mathbf{w} = \sum_{i=1}^{N} \boldsymbol{\alpha}_{i} y_{i} \Phi(\mathbf{x}_{i}) \qquad s(\mathbf{x}) := \sum_{k=1}^{K} \sum_{i=1}^{L-k+1} w(\mathbf{x}[i]^{k}, i) + b$$

k-mer	pos. 1	pos. 2	pos. 3	pos. 4	
Α	+0.1	-0.3	-0.2	+0.2	
С	0.0	-0.1	+2.4	-0.2	
G	+0.1	-0.7	0.0	-0.5	
т	-0.2	-0.2	0.1	+0.5	
AA	+0.1	-0.3	+0.1	0.0	
AC	+0.2	0.0	-0.2	+0.2	
÷	:	÷	÷	÷	·
тт	0.0	-0.1	+1.7	-0.2	
AAA	+0.1	0.0	0.0	+0.1	
AAC	0.0	-0.1	+1.2	-0.2	
÷	÷	÷	:	:	·
TTT	+0.2	-0.7	0.0	0.0	



Introduction Sequence Classification Support Vector Machines Applications Explanation and Visualization Summary

The Scoring System - Examples

$$s(\mathbf{x}) := \sum_{k=1}^{K} \sum_{i=1}^{L-k+1} w\left(\mathbf{x}[i]^k, i\right) + b$$

Examples:

- WD-kernel (Rätsch, Sonnenburg, 2005)
- WD-kernel with shifts (Rätsch, Sonnenburg, 2005)
- Spectrum kernel (Leslie, Eskin, Noble, 2002)
- Oligo Kernel (Meinicke et al., 2004)

Not limited to SVM's:

• Markov Chains (higher order/inhomogeneous/mixed order)



Introduction 00000 Sequence Classification 0000

Support Vector Machines

Applications Explanation and Visualization

Summary

Introduction

The SVM Weight Vector w



- Explicit representation of **w** allows for (some) interpretation!
- String kernel SVMs capable of efficiently dealing with large k-mers k > 10

But: Weights for substrings not independent



Introduction Se 00000 0

Sequence Classification

Support Vector Machines

Applications Explanation and Visualization

Introduction

Interdependence of k-mer Weights



What is the score for TAC?

- Take *w_{TAC}*?
- But substrings and overlapping strings contribute too!

Problem

The SVM-w does NOT reflect the score for a motif



Positional Oligomer Importance Matrices (POIMs)

Idea:

• 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 (Sonnenburg et al. (2008) $Q(\mathbf{z}, j) := \mathbb{E}[s(\mathbf{x}) | \mathbf{x}[j] = \mathbf{z}] - \mathbb{E}[s(\mathbf{x})]$

\Rightarrow Needs efficient algorithm for computation

Efficient Computation

Effort of naive approach exponential $\mathcal{O}(|\Sigma|^{L} + L|\Sigma|^{k})$ (e.g. Splice Sites 10¹²⁰)

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

- Number of k-mers grows linearly with size of input
- Only features which are dependent on (z, j) matter
- Computation can be split in contributions from 4 cases

Efficient Recursive Algorithm: Effort linear in length of input: $O(LN + L|\Sigma|^k)$

Introduction	Sequence Classification	Support Vector Machines	Applications	Explanation and Visualization	Sum
				000000000000000000000000000000000000000	
Definition					

Ranking Features and Condensing Information

- Obtain highest scoring z from Q(z, i) (Enhancer or Silencer)
- Visualize POIM as heat map; x-axis: position y-axis: k-mer color: importance
- For large k: Differential POIMs; x-axis: position y-axis: k-mer length color: importance

z	i	$Q(\mathbf{z}, i)$
GATTACA	10	+30
AGTAGTG	30	+20
AAAAAAA	10	-10







Introduction 00000 Sequence Classification

Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

GATTACA and AGTAGTG at Fixed Positions 10 and 30

TGAGCGCGT<mark>GATTACA</mark>GTCCGTCTGGGCC<mark>AGTAGTG</mark>CGTAGTCGCCGGGA GGCATGGTC<mark>GATTACA</mark>AACGAGCCCTCTC<mark>AGTAGTG</mark>GGGGAGCCACGAAA CCCGTCGAA<mark>GATTACA</mark>CACGGGGGCGTGGG<mark>AGTAGTG</mark>GCGATTACGGGCTC GGTCGGCAG<mark>GATTACA</mark>CGACGCGTTTACG<mark>AGTAGTG</mark>AACACTGACTCCTC



Introduction Sequence

Sequence Classification

Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

w

GATTACA and AGTAGTG at fixed positions 10 and 30

TGAGCGCGT<mark>GATTACA</mark>GTCCGTCTGGGCC<mark>AGTAGTG</mark>CGTAGTCGCCGGGA GGCATGGTC<mark>GATTACA</mark>AACGAGCCCTCTC<mark>AGTAGTG</mark>GGGGAGCCACGAAA CCCGTCGAA<mark>GATTACA</mark>CACGGGGCGTGGG<mark>AGTAGTG</mark>GCGATTACGGGCTC GGTCGGCAG<mark>GATTACA</mark>CGACGCGTTTACG<mark>AGTAGTG</mark>AACACTGACTCCTC





Introduction Sequence C 00000 0000

Sequence Classification

Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

GATTACA and AGTAGTG at fixed positions 10 and 30

TGAGCGCGT<mark>GATTACA</mark>GTCCGTCTGGGCC<mark>AGTAGTG</mark>CGTAGTCGCCGGGA GGCATGGTC<mark>GATTACA</mark>AACGAGCCCTCTC<mark>AGTAGTG</mark>GGGGAGCCACGAAA CCCGTCGAA<mark>GATTACA</mark>CACGGGGGCGTGGG<mark>AGTAGTG</mark>GCGATTACGGGCTC GGTCGGCAG<mark>GATTACA</mark>CGACGCGTTTACG<mark>AGTAGTG</mark>AACACTGACTCCTC





Introduction 9

Sequence Classification

Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

GATTACA and AGTAGTG at fixed positions 10 and 30

TGAGCGCGT<mark>GATTACA</mark>GTCCGTCTGGGCC<mark>AGTAGTG</mark>CGTAGTCGCCGGGA GGCATGGTC<mark>GATTACA</mark>AACGAGCCCTCTC<mark>AGTAGTG</mark>GGGGAGCCACGAAA CCCGTCGAA<mark>GATTACA</mark>CACGGGGGCGTGGG<mark>AGTAGTG</mark>GCGATTACGGGCTC GGTCGGCAG<mark>GATTACA</mark>CGACGCGTTTACG<mark>AGTAGTG</mark>AACACTGACTCCTC



Introduction 00000 Sequence Classification 0000

Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

GATTACA at variable positions

TGAGCGCGTGATTACAGTCCGTCT GGCTCGATCACAAACGAGCCCGAT CCCGTCGAACAGGATTACACACGG GGTCGGCAGCTTACACGACAGCGT


$\begin{array}{c} \text{Sequence Classification} \\ \text{0000} \end{array}$

Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

GATTACA at variable positions





weblogo.berkeley.edu



Sequence Classification 0000

Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

GATTACA at variable positions

TGAGCGCGTGATTACAGTCCGTCT GGCTCGATCACAAACGAGCCCGAT CCCGTCGAACAGGATTACACACGG GGTCGGCAGCTTACACGACAGCGT

Differential POIM Overview - GATTACA shift





Introduction Sequence Classifica

cation Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

C.elegans Acceptor Splice Site Recognition



•	Upstream	
	AG GT AAGT	-44/+
	GGGGGG	-16/-
	TAATAA	-16/+

- ++ Donor - Silencer? ++ Branch
- Central TTTTTC -06/+TTTCAG $\frac{A}{G}$ -03/++
 - + Acceptor
- Downstream TTTTTTT +07/- -TTTTT +26/- -



Sequence Classification

Support Vector Machines

Applications Explanation and Visualization

Summary

Applications

Drosophila Transcription Starts



Differential POIM Overview - Drosophila TSS

equence Classification

Support Vector Machines

Applications Explanation and Visualization Summary

Outline

- Introduction
 - Genomic Signals
- Sequence Classification
 Support Vector Machines
- 3 Support Vector Machines
- 4 Applications
 - Splice Site Recognition
 - TSS Recognition
 - Gene Finding mGene
 - Aligning Short Reads QPALMA
- 5 Explanation and Visualization
 - Introduction
 - Definition
 - Applications





Introduction Sequence Classification Support Vector Machines Applications Explanation and Visualization

Conclusions

Support Vector Machines with String Kernels

- General and often state-of-the art signal detectors
- Applicable to genome-sized datasets
- Using POIMs SVMs are interpretable

Software Available

- TSS Detector ARTS http://mloss.org/software/view/191/
- Accurate splice detector http://mloss.org/software/view/192/
- mGene http://mgene.org
- QPalma http://www.fml.tuebingen.mpg.de/raetsch/suppl/qpalma
- ML implemented in http://www.shogun-toolbox.org

More machine learning software http://mloss.org



Summarv