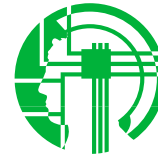


# Learning interpretable SVMs for biological sequence classification

Sören Sonnenburg<sup>†</sup>, Gunnar Rätsch<sup>†‡</sup>, Christin Schäfer<sup>†</sup>



† Fraunhofer FIRST.IDA, Kekuléstr. 7, 12489 Berlin, Germany  
‡ Friedrich Miescher Lab. of the Max Planck Society,  
Spemannstr. 39, 72076 Tübingen, Germany  
{Soeren.Sonnenburg, Christin.Schaefer}@first.fraunhofer.de,  
Gunnar.Raetsch@tuebingen.mpg.de

## ROADMAP:

- The Motivating Application
- Multiple Kernel Learning (MKL)
- Derivation of the MKL Optimization Problem
- Algorithms
- Significance Analysis
- Results
- Outlook and Conclusion

# BIOLOGY: DETECTION OF SPLICE SITES

Intron Exon

```
AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG
AAGATTAATAAAAAAAAAACAAATTTTTCATTACAGATATAATAATCTAATT
CACTCCCAAATCAACGATATTTTAGTTCACTAACACATCCGTCTGTGCC
TTAATTTCACTTCCACATACTTCCAGATCATCAATCTCCAAAACCAACAC
TTGTTTTAATATTCAATTTTTTACAGTAAGTTGCCAATTCAATGTTCCAC
TACCTAATTATGAAATTAATAATTCAGTGTGCTGATGGAAACGGAGAAGTC
```

- aligned sequences of fixed length (AG always at position 61 & 62)
- Task: distinguish splice sites from fake - splice sites

⇒ 2-class classification problem

## APPROACH: STRING KERNEL + SVM

- use SVM classifier 
$$f(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^N y_i \alpha_i \mathbf{k}(\mathbf{x}, \mathbf{x}_i) + b \right)$$
- key ingredient is kernel  $\mathbf{k}(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}') \Rightarrow$  gives means to compare 2 sequences
- Kernel here “Weighted Degree Kernel”  $\beta_k$  ( $k = 1, \dots, d$ ):

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \sum_{k=1}^d \beta_k \sum_{l=1}^{L-k} \mathbf{I}(\mathbf{u}_{k,l}(\mathbf{x}) = \mathbf{u}_{k,l}(\mathbf{x}')),$$

```

AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG
.|.|.||||.|..||.|.|..||||.||...|...|...|||.....|..
TACCTAATTATGAAATTAAAATTCAGTGTGCTGATGGAAACGGAGAAGTC

```

- $L$  length of the  $\mathbf{x}$ 's,  $d$  maximal “match length” and  $\mathbf{u}_{k,l}(\mathbf{x})$  subsequence of length  $k$  at position  $l$  of sequence  $\mathbf{x}$

# SUCCESS

## Choosing a particular weighting

$$\beta_k = \frac{2(d - k + 1)}{(d(d + 1))}$$

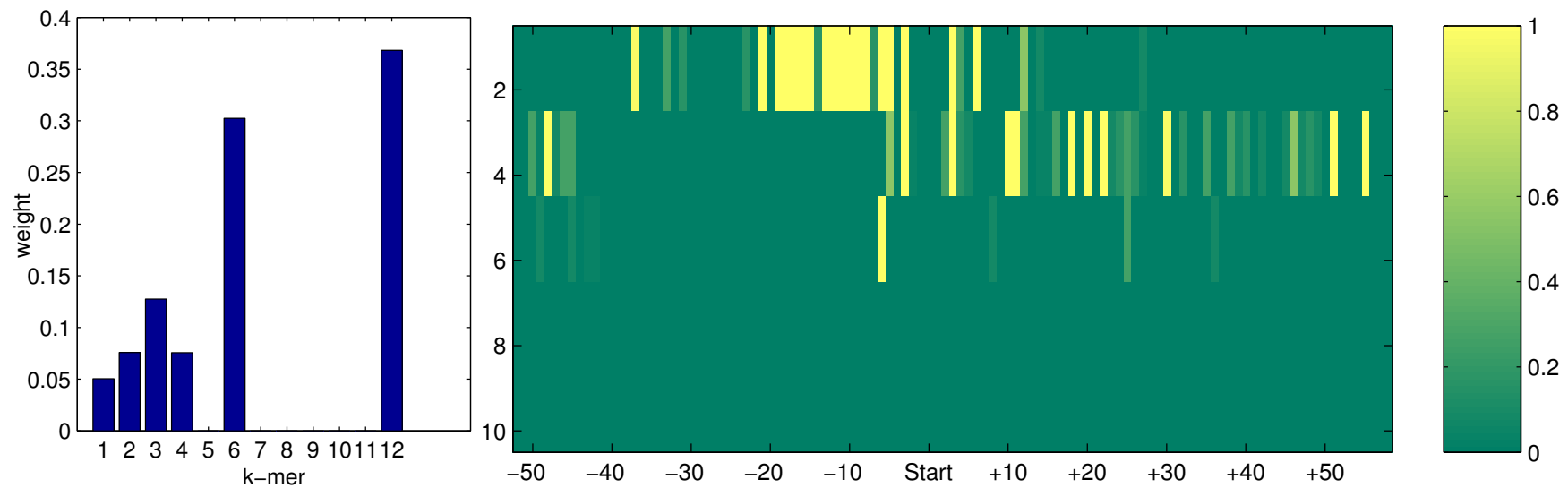
seems to solve the task: on 500,000 training examples test AUC 99,80%  
(test error 0.78%)

## Open questions

- Why does that weighting make sense ?
- Is there a better weighting ?
  - in terms of sense
  - classification performance
- Can we learn that weighting ?
- What does that have to do with optimization ?

# GAIN

- automated model selection
- interpretability
  - which  $k$ -mers are of higher importance
  - even which  $k$ -mers where in the sequence are of higher importance



## REFORMULATION: MULTIPLE KERNEL LEARNING

The Weighted Degree kernel is a linear combination of kernels !

$$\begin{aligned} \mathbf{k}(\mathbf{x}, \mathbf{x}') &= \sum_{k=1}^d \beta_k \sum_{l=1}^{L-k} \mathbb{I}(\mathbf{u}_{k,l}(\mathbf{x}) = \mathbf{u}_{k,l}(\mathbf{x}')) \\ &= \sum_{k=1}^d \beta_k \mathbf{k}_k(\mathbf{x}, \mathbf{x}') \end{aligned}$$

with

$$\mathbf{k}_k(\mathbf{x}, \mathbf{x}') = \sum_{l=1}^{L-k} \mathbb{I}(\mathbf{u}_{k,l}(\mathbf{x}) = \mathbf{u}_{k,l}(\mathbf{x}')).$$

(This also holds for other String Kernels like the Spectrum Kernel)

⇒ need to solve the so called **Multiple Kernel Learning Problem**,  
i.e. determine  $(\beta, \alpha, b)$  simultaneously.

# MULTIPLE KERNEL LEARNING

## Motivation

- Kernel  $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$  used in standard SVM Classifier

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

- Now: linear combination of kernels (again a kernel)

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^M \beta_j k_j(\mathbf{x}, \mathbf{x}'), \beta_j \geq 0$$

- **useful:** Polynomial kernels of different degree, kernels on different domain
- **but:** How to **learn** and **constrain** weights  $\beta_j$  ?



# CONSTRAINING THE WEIGHTS

## $L_2$ -vs. $L_1$ -Norm

- in max problem weights certain  $\beta_j$  would grow infinitely (min shrink to zero)  $\Rightarrow$  **constraining  $\beta_j$  necessary**
- Dense  $\|\beta\|_2 = 1$  (Lin and Zhang 2004)
- vs. Sparse  $\|\beta\|_1 = 1$  (Bach, Lanckriet and Jordan 2004)
  - convex combination of kernels
  - sparse solution in terms of kernels
  - allows for interpretation of result

**constraints on  $\beta_j$ :**

$$\sum_{j=1}^N \beta_j = 1, \beta_j \geq 0$$

# STANDARD SVM OPTIMIZATION PROBLEM

## This we all know

SVM Primal formulation:

$$\min \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i$$

$$\text{w.r.t.} \quad \mathbf{w} \in \mathbb{R}^k, \boldsymbol{\xi} \in \mathbb{R}_+^N, b \in \mathbb{R}$$

$$\text{s.t.} \quad y_i (\mathbf{w}^\top \Phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \forall i = 1, \dots, N$$

# MKL OPTIMIZATION PROBLEM I

MKL Primal formulation:

$$\min \quad \frac{1}{2} \left( \sum_{j=1}^M \beta_j \|\mathbf{w}_j\|_2 \right)^2 + C \sum_{i=1}^N \xi_i$$

$$\text{w.r.t.} \quad \mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_M), \mathbf{w}_j \in \mathbb{R}^{k_j}, \boldsymbol{\xi} \in \mathbb{R}_+^N, \boldsymbol{\beta} \in \mathbb{R}_+^M, b \in \mathbb{R}$$

$$\text{s.t.} \quad y_i \left( \sum_{j=1}^M \beta_j \mathbf{w}_j^\top \Phi_j(\mathbf{x}_i) + b \right) \geq 1 - \xi_i, \forall i = 1, \dots, N$$

$$\sum_{j=1}^M \beta_j = 1$$

**Properties:** equivalent to SVM for  $M = 1$ ; solution sparse in “blocks”; each block  $j$  corresponds to one kernel

# MKL OPTIMIZATION PROBLEM II

**Dual Formulation (Bach, Lanckriet, Jordan 2004):**

$$\begin{aligned}
 \min \quad & \frac{1}{2}\gamma^2 - \sum_{i=1}^N \alpha_i \\
 \text{w.r.t.} \quad & \gamma \in \mathbb{R}, \boldsymbol{\alpha} \in \mathbb{R}^N \\
 \text{s.t.} \quad & 0 \leq \boldsymbol{\alpha} \leq C, \sum_{i=1}^N \alpha_i y_i = 0 \\
 & \underbrace{\sum_{r=1}^N \sum_{s=1}^N \alpha_r \alpha_s y_r y_s K_j(\mathbf{x}_r, \mathbf{x}_s)}_{=: S_j(\boldsymbol{\alpha})} - \gamma^2 \leq 0, \quad \forall j = 1, \dots, M
 \end{aligned}$$

“partial Lagrangian:”

$$L := \frac{1}{2}\gamma^2 - \sum_{i=1}^N \alpha_i + \sum_{j=1}^M \beta_j (S_j(\boldsymbol{\alpha}) - \gamma^2)$$

# MKL OPTIMIZATION PROBLEM II

## Reformulation as Semi-Infinite Linear Program:

$$\begin{aligned} & \max_{\beta} \min_{\alpha} \sum_{j=1}^M \beta_j \left( \frac{1}{2} S_j(\alpha) - \sum_{i=1}^N \alpha_i \right) \\ & \text{s.t. } 0 \leq \alpha \leq C, \sum_{i=1}^N \alpha_i y_i = 0, \sum_{j=1}^M \beta_j = 1 \end{aligned}$$

$$\begin{aligned} & \max \quad \theta \\ & \text{w.r.t.} \quad \theta \in \mathbb{R}, \beta \in \mathbb{R}_+^M \text{ with } \sum_{j=1}^M \beta_j = 1 \\ & \text{s.t.} \quad \sum_{j=1}^M \beta_j \left( \frac{1}{2} S_j(\alpha) - \sum_{i=1}^N \alpha_i \right) \geq \theta \\ & \quad \text{for all } \alpha \text{ with } 0 \leq \alpha \leq C \text{ and } \sum_{i=1}^N y_i \alpha_i = 0 \end{aligned}$$

⇒ **Linear, but infinitely many constraints**

# THE SEMI-INFINITE LINEAR PROGRAM I

$$\begin{aligned}
 & \max && \theta \\
 & \text{w.r.t.} && \theta \in \mathbb{R}, \boldsymbol{\beta} \in \mathbb{R}_+^M \text{ with } \sum_{j=1}^M \beta_j = 1 \\
 & \text{s.t.} && \sum_{j=1}^M \beta_j \left( \frac{1}{2} S_j(\boldsymbol{\alpha}) - \sum_{i=1}^N \alpha_i \right) \geq \theta \\
 & && \text{for all } \boldsymbol{\alpha} \text{ with } 0 \leq \boldsymbol{\alpha} \leq C \text{ and } \sum_{i=1}^N y_i \alpha_i = 0
 \end{aligned}$$

Properties:

- optimize a convex combination
- infinitely many constraints
- quite easy to identify violated constraints

# THE SEMI-INFINITE LINEAR PROGRAM II

$$\begin{aligned}
 & \max && \theta \\
 & \text{w.r.t.} && \theta \in \mathbb{R}, \boldsymbol{\beta} \in \mathbb{R}_+^M \text{ with } \sum_{j=1}^M \beta_j = 1 \\
 & \text{s.t.} && \sum_{j=1}^M \beta_j \left( \frac{1}{2} S_j(\boldsymbol{\alpha}) - \sum_{i=1}^N \alpha_i \right) \geq \theta \\
 & && \text{for all } \boldsymbol{\alpha} \text{ with } 0 \leq \boldsymbol{\alpha} \leq C \text{ and } \sum_{i=1}^N y_i \alpha_i = 0
 \end{aligned}$$

Solving the SILP:

- Column Generation
  - fast, but no convergence rate
- Use Boosting like techniques: Arc-GV or AdaBoost\*
  - known convergence rate  $\mathcal{O}(\log(M)/\varepsilon^2)$
- SMO like algorithm
  - consider suboptimal SVM solutions: empirically 3-5 times faster

# SOLVING THE SILP: COLUMN GENERATION I

$$\begin{aligned}
 & \max && \theta \\
 & \text{w.r.t.} && \theta \in \mathbb{R}, \boldsymbol{\beta} \in \mathbb{R}_+^M \text{ with } \sum_{j=1}^M \beta_j = 1 \\
 & \text{s.t.} && \sum_{j=1}^M \beta_j \left( \frac{1}{2} S_j(\boldsymbol{\alpha}) - \sum_{i=1}^N \alpha_i \right) \geq \theta \\
 & && \text{for all } \boldsymbol{\alpha} \text{ with } 0 \leq \boldsymbol{\alpha} \leq C \text{ and } \sum_{i=1}^N y_i \alpha_i = 0
 \end{aligned}$$

- solved by taking set of most violated constraints into account
- most violated constraints given by SVM solution for fixed  $\boldsymbol{\beta}$

$$\sum_{j=1}^M \beta_j \left( \frac{1}{2} S_j(\boldsymbol{\alpha}) - \sum_{i=1}^N \alpha_i \right) = \frac{1}{2} \sum_{r=1}^N \sum_{s=1}^N \alpha_r \alpha_s y_r y_s \sum_{j=1}^M \beta_j k_j(\mathbf{x}_r, \mathbf{x}_s) - \sum_{i=1}^N \alpha_i,$$

- iteratively find most violated constraints, solve linear program with current constraints, . . . , till convergence to the global optimum



# SOLVING THE SILP: COLUMN GENERATION II

$$D^0 = 1, \theta^1 = 0, \beta_k^1 = \frac{1}{M} \text{ for } k = 1, \dots, M$$

**for**  $t = 1, 2, \dots$  **do**

  obtain SVM's  $\alpha^k$  with kernel

$$k^t(\mathbf{x}_i, \mathbf{x}_j) := \sum_{k=1}^M \beta_k^t k_k(\mathbf{x}_i, \mathbf{x}_j)$$

**for**  $k = 1, \dots, M$  **do**

$$D_k^t = \frac{1}{2} \sum_{r,s} \alpha_r^t \alpha_s^t y_r y_s k_k(\mathbf{x}_r, \mathbf{x}_s) - \sum_r \alpha_r^t$$

**end for**

$$D^t = \sum_{k=1}^M \beta_k^t D_k^t$$

$$(\beta^{t+1}, \theta^{t+1}) = \operatorname{argmax} \theta$$

$$\text{w.r.t. } \beta \in \mathbb{R}_+^M, \theta \in \mathbb{R} \text{ with } \sum_k \beta_k = 1$$

$$\text{s.t. } \sum_{k=1}^M \beta_k D_k^r \geq \theta \text{ for } r = 1, \dots, t$$

**if**  $|1 - \frac{\theta^{t+1}}{D^t}| \leq \epsilon$  **then break**

**end for**

# SOLVING THE SILP: BOOSTING I

## Boosting

### Primal

$$\begin{aligned} & \max_{\alpha, \rho} \quad \rho \\ \text{subject to} \quad & y_n \sum_{j=1}^J \alpha_j h_j(\mathbf{x}_n) \geq \rho \\ & \alpha_j \geq 0 \\ & \sum_j \alpha_j = 1 \end{aligned}$$

### Dual

$$\begin{aligned} & \min_{d, \delta} \quad \delta \\ \text{subject to} \quad & \sum_{n=1}^N d_n y_n h_j(\mathbf{x}_n) \leq \delta \\ & d_n \geq 0 \\ & \sum_n d_n = 1 \end{aligned}$$

- max w.r.t.  $-\delta$ ;  $d_n$  corresponds to  $\beta_n$ ;  $h_j$  corresponds to expression using a certain  $\alpha$
- Arc-GV and AdaBoost\*:  $\rho$  and  $\delta$  are optimized

The number of hypotheses can be very large or infinite!

# SOLVING THE SILP: BOOSTING I

## An Arc-GV like Algorithm

$D^0 = 1, \rho^1 = \tau_k^1 = 0, \beta_k^1 = \frac{1}{M}$  for  $k = 1, \dots, M$

**for**  $t = 1, 2, \dots$  **do**

obtain SVM's  $\alpha^k$  with kernel

$$k^t(\mathbf{x}_i, \mathbf{x}_j) := \sum_{k=1}^M \beta_k^t k_k(\mathbf{x}_i, \mathbf{x}_j)$$

**for**  $k = 1, \dots, M$  **do**

$$D_k^t = \frac{1}{2} \sum_{r,s} \alpha_r^t \alpha_s^t y_r y_s k_k(\mathbf{x}_r, \mathbf{x}_s) - \sum_r \alpha_r^t$$

**end for**

$$D^t = \sum_{k=1}^M \beta_k^t D_k^t$$

$$\gamma_t = \operatorname{argmin}_{\gamma \in [0,1]} \sum_{k=1}^M \beta_k^t \exp \left\{ \gamma (D_k - \rho^t) \right\}$$

**for**  $k = 1, \dots, M$  **do**

$$\tau_k^{t+1} = \tau_k^t + \gamma_t D_k^t$$

$$\beta_k^{t+1} = \beta_k^t \exp(\gamma_t D_k^t) / \left( \sum_{k'} \beta_{k'}^t \exp(\gamma_t D_{k'}^t) \right)$$

**end for**

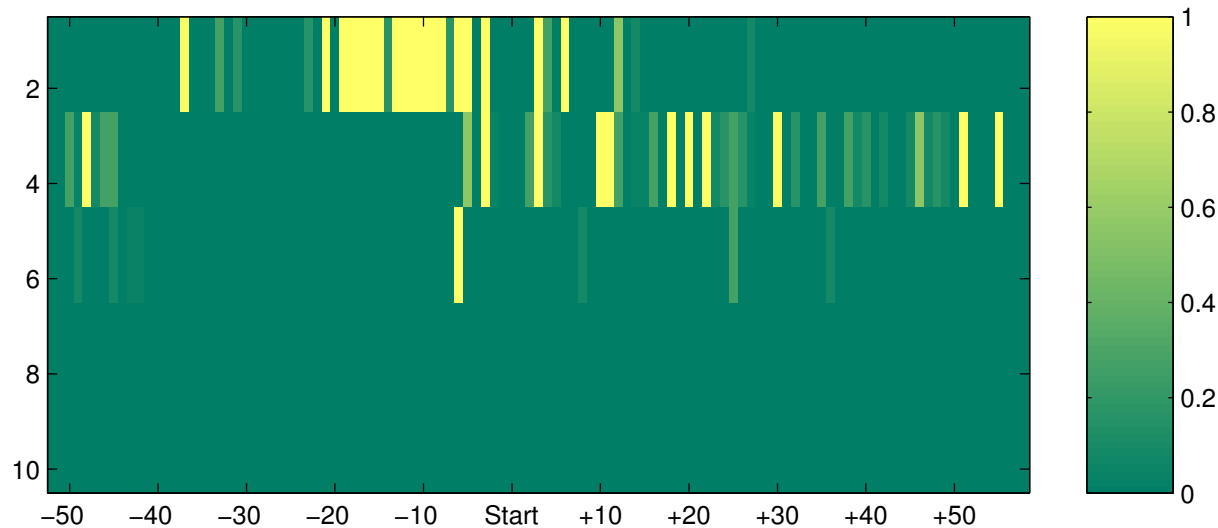
$$\rho^{t+1} = \max_k \tau_k^{t+1} / \sum_{r=1}^t \gamma_r$$

**if**  $\left| 1 - \frac{\rho^{t+1}}{D^t} \right| \leq \epsilon$  **then break**

**end for**

## STABILITY OF THE SOLUTION ?

We obtain a weighting

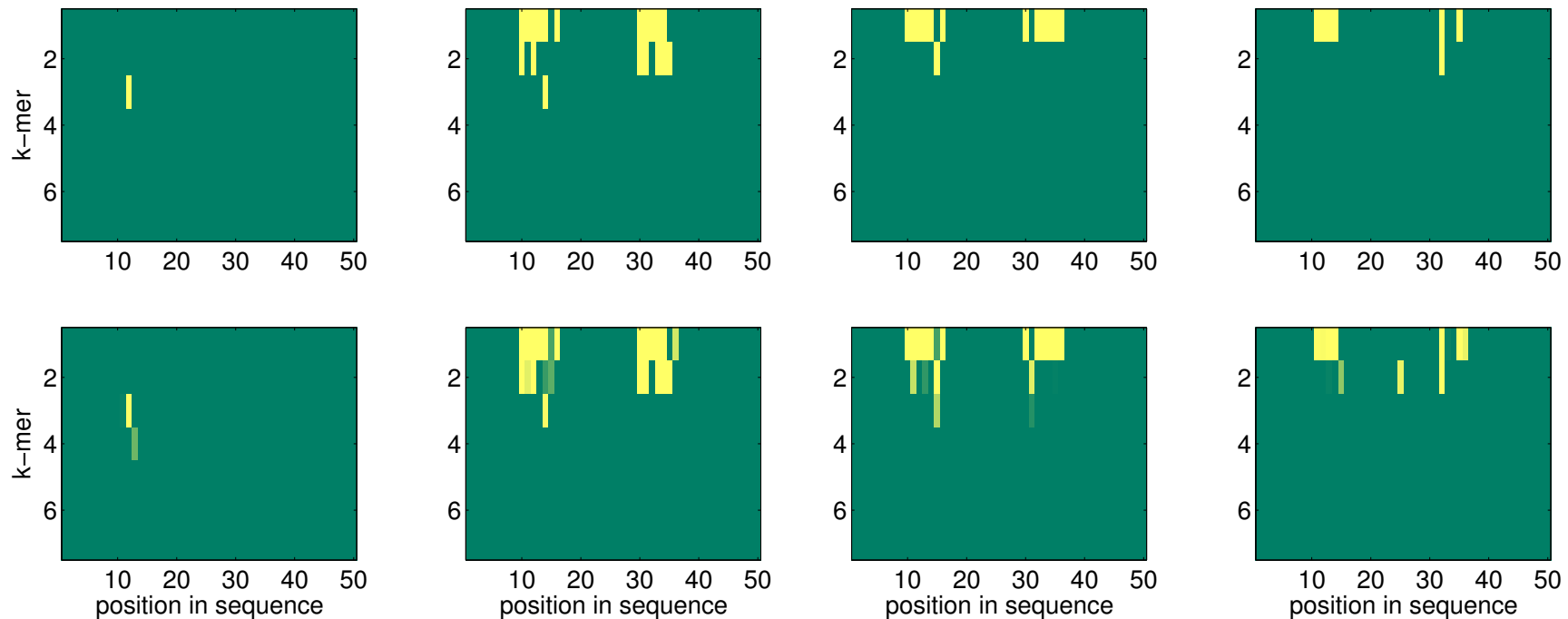


Questions:

- Can we use that weighting to interpret the SVM solution ?
  - Stability of the weighting  $\beta$  ?
  - Which weights are significant ?

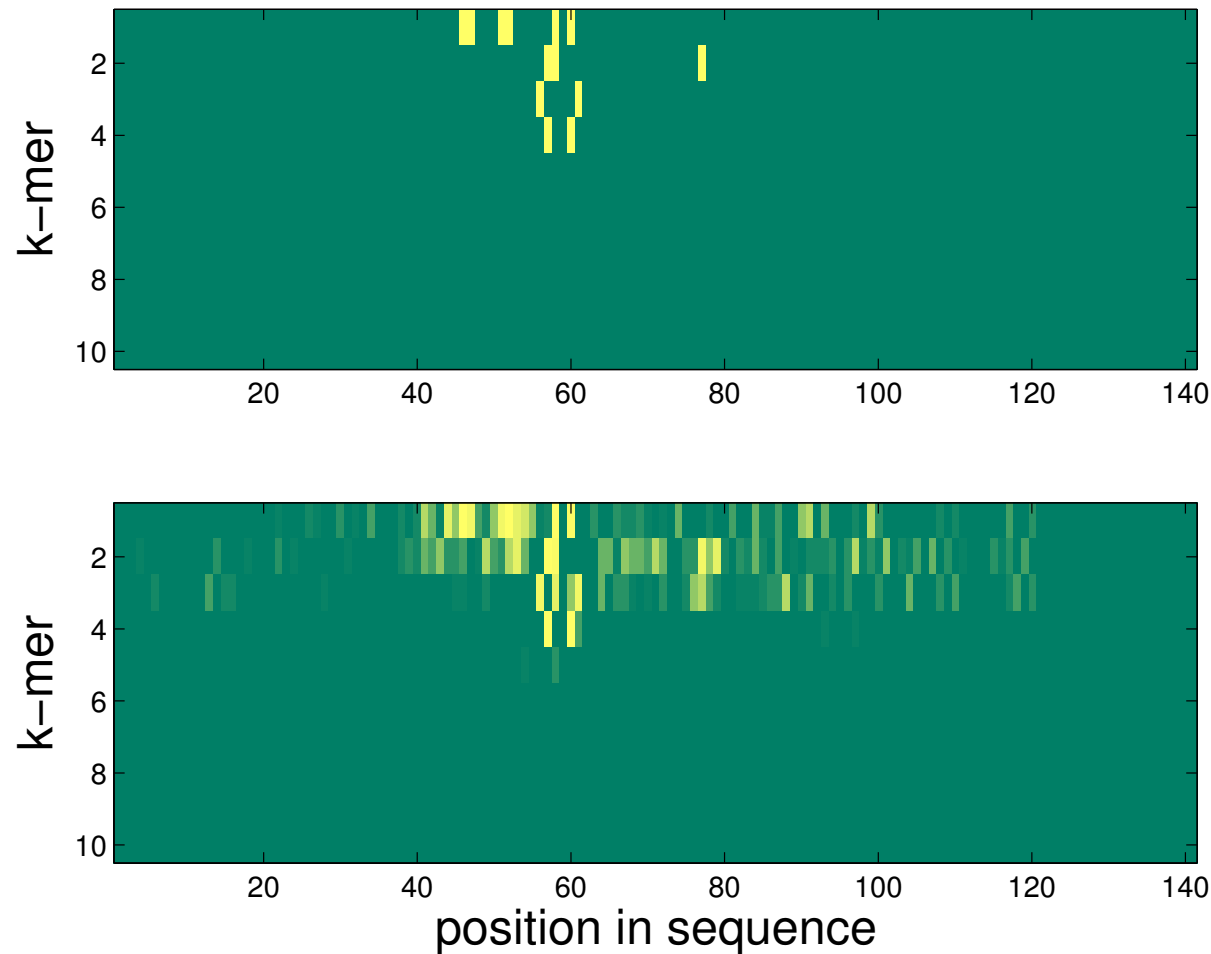
**We used a statistical test to investigate significance of weights**

# TOY DATASET



- DNA sequences of length 50 with hidden motifs at 10-16 and 30-36
- $8 \times 50$  string kernels with max. word-length 8
- compute significance level by bootstrapping
- columns  $\equiv$  noise level
- subplot columns weights used at certain position, rows oligomer length

# APPLICATION TO ACCEPTOR SPLICE SITES



- *C.elegans* splice dataset (AG at position 61 & 62)
- subplot columns weights used at certain position, rows oligomer length

## CONCLUSION

### Conclusion:

- suitable for large scale (WD kernel of degree 20 trained on 200000 examples)
- fast, because  $t(N, M, \epsilon) \approx 5 \cdot 10^{-11} M^{2.22} N^{1.68} \log_2(1/\epsilon)^{2.52}$
- for certain kernels: allows for interpreting SVM result

### Future Work:

- feature selection
- apply to regression

**Acknowledgements:** Partial support from the PASCAL Network of Excellence (EU #506778), DFG grants JA379/13-2 and MU987/2-1. Thanks to Alexander Zien, Gökhan H. Bakır and K.-R. Müller.