

Learning interpretable SVMs for biological sequence classification

Sören Sonnenburg[†], Gunnar Rätsch^{†,‡}, Christin Schäfer[†]



Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik



MAX-PLANCK-GESELLSCHAFT

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

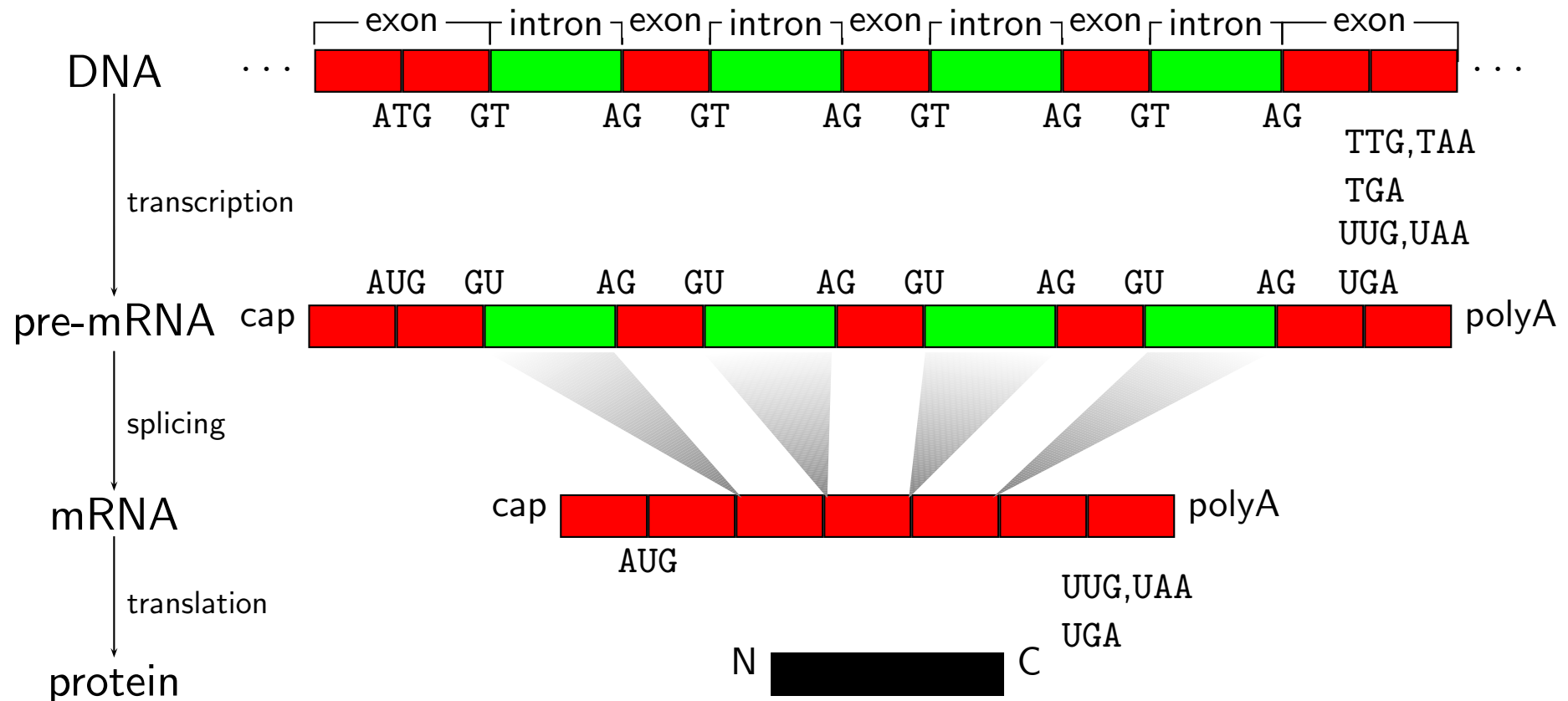
ROADMAP:

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

THE MOTIVATING APPLICATION

Splice sites are locations on DNA at boundaries of

- **exons** (which code for proteins)
- **introns** (which do not)



BIOLOGY: DETECTION OF SPLICE SITES

Intron Exon

```
AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG
AAGATTAAAAAAAACAAATTTTTCAGCATTACAGATATAATAATCTAATT
CACTCCCAAATCAACGATATTTTAGTTCACTAACACATCCGTCTGTGCC
TTAATTTCACTTCCACATACTTCCAGATCATCAATCTCCAAAACCAACAC
TTGTTTTAATATTCAATTTTTTACAGTAAGTTGCCAATTCAATGTTCCAC
TACCTAATTATGAAATTAATAATTCAGTGTGCTGATGGAAACGGAGAAGTC
```

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

⇒ 2-class classification problem

APPROACH: STRING KERNEL + SVM

- use SVM classifier
$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^N y_i \alpha_i \mathbf{k}(\mathbf{x}, \mathbf{x}_i) + b \right)$$
- find parameters α by solving quadratic optimization problem:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j) \\ \text{subject to} \quad & \alpha_i \in [0, C], i = 1, \dots, N, \sum_{i=1}^N \alpha_i y_i = 0. \end{aligned}$$

- **Solution has no local minima**

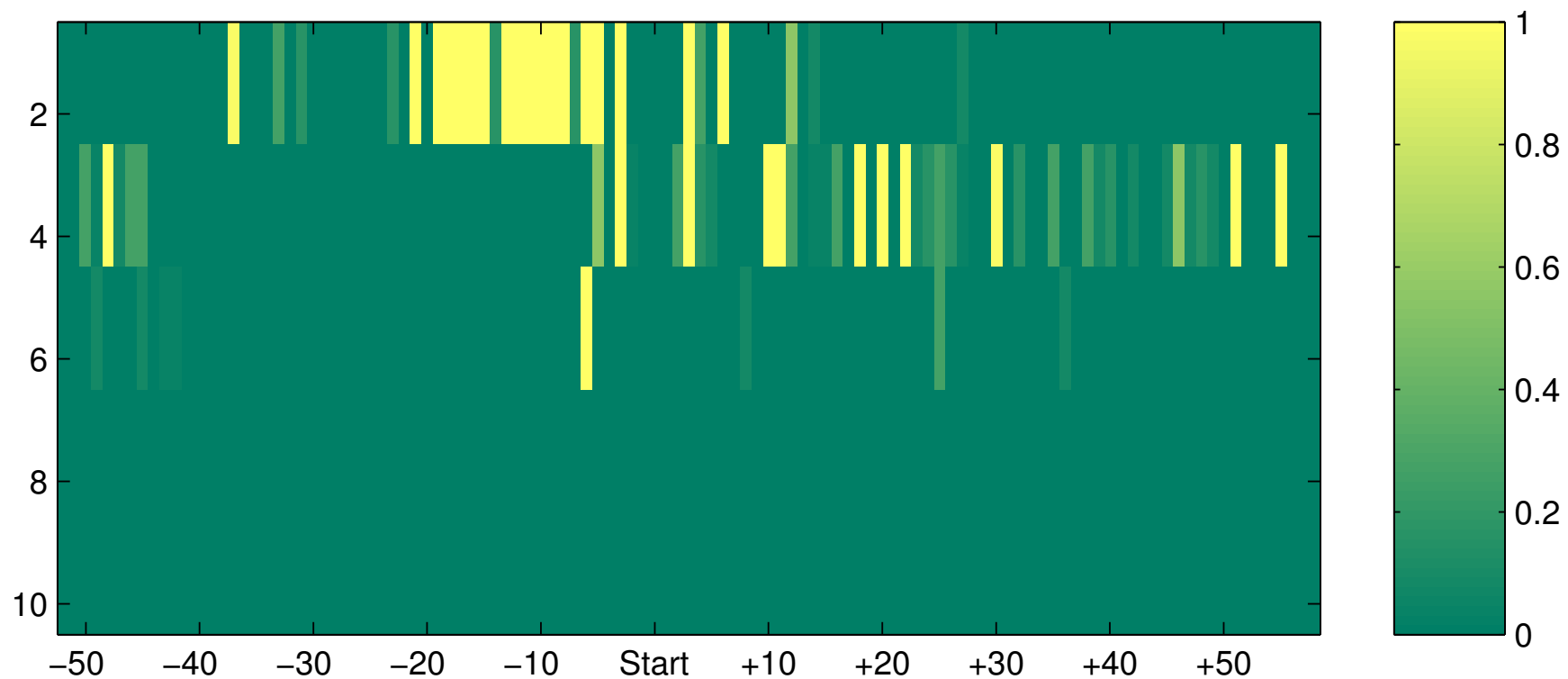
Key ingredient **Kernel** here “**Weighted Degree Kernel**”

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

POSITION DEPENDANT WEIGHTING β

Even more weights:

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



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 Multiple Kernel Learning ?

REFORMULATION: MULTIPLE KERNEL LEARNING

The Weighted Degree kernel is a linear combination of kernels !

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

with

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

(Can also be applied to other String Kernels)

⇒ need to solve the so called **Multiple Kernel Learning Problem**, i.e. determine (β, α, 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 β_j ?

CONSTRAINING THE WEIGHTS

L_2 -vs. L_1 -Norm

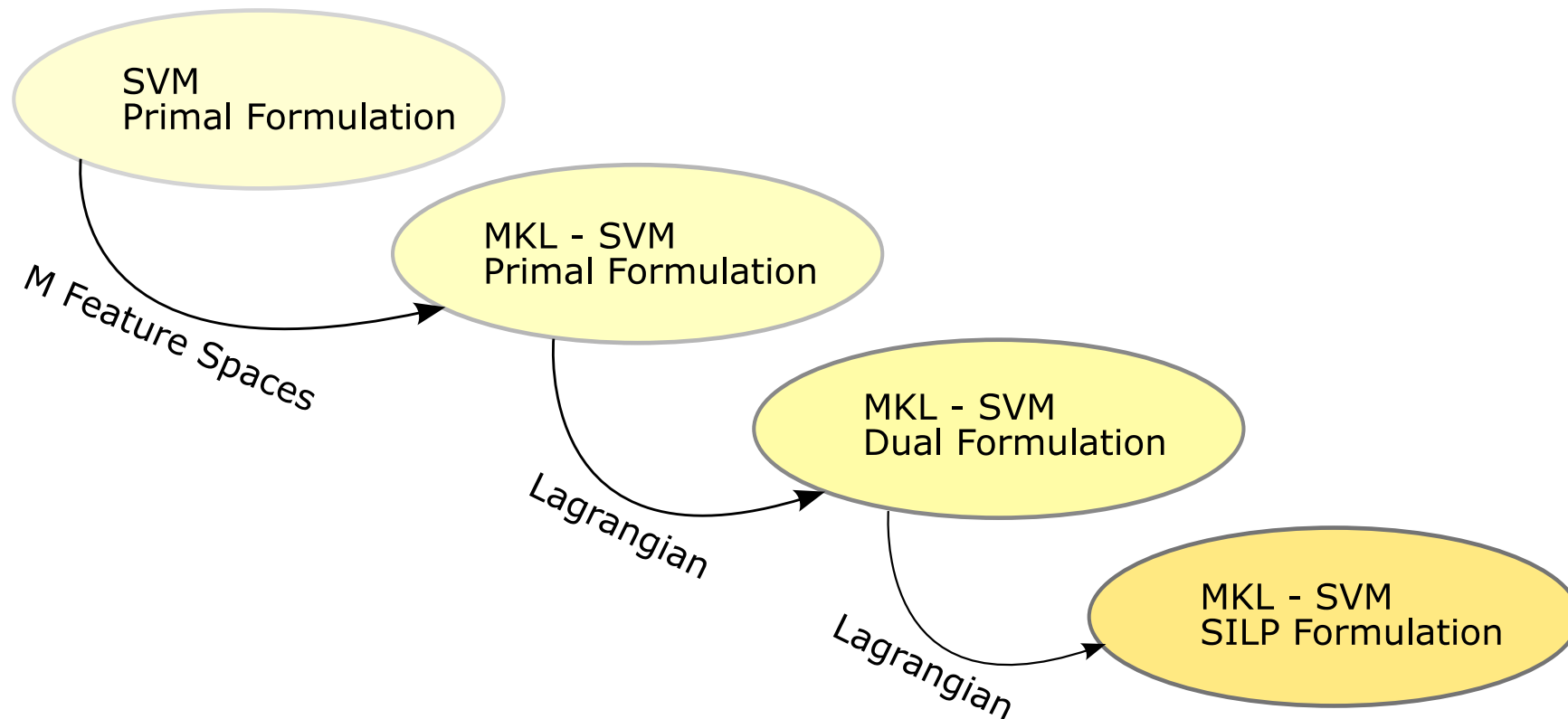
- in max problem weights certain β_j would grow infinitely (min shrink to zero) \Rightarrow **constraining β_j necessary**
- Dense $\|\beta\|_2 = 1$ (Lin and Zhang 2004)
- vs. Sparse $\|\beta\|_1 = 1$ (Bach, Lanckriet and Jordan 2004)
 - convex combination of kernels
 - sparse solution in terms of kernels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

constraints on β_j :

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

REWRITING THE MKL FORMULATION

From the Standard SVM Primal to the Semi-Infinite Linear Program:



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:

$$\max_{\beta} \min_{\alpha} \sum_{j=1}^M \beta_j \left(\frac{1}{2} S_j(\alpha) - \sum_{i=1}^N \alpha_i \right)$$

$$\text{s.t. } 0 \leq \alpha \leq C, \sum_{i=1}^N \alpha_i y_i = 0, \sum_{j=1}^M \beta_j = 1$$

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

⇒ **Linear, but infinitely many constraints**

THE SEMI-INFINITE LINEAR PROGRAM

Properties:

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

Solving the SILP:

- Use Boosting like techniques: Arc-GV or AdaBoost*
- Column Generation
- SMO like algorithm

COLUMN GENERATION I

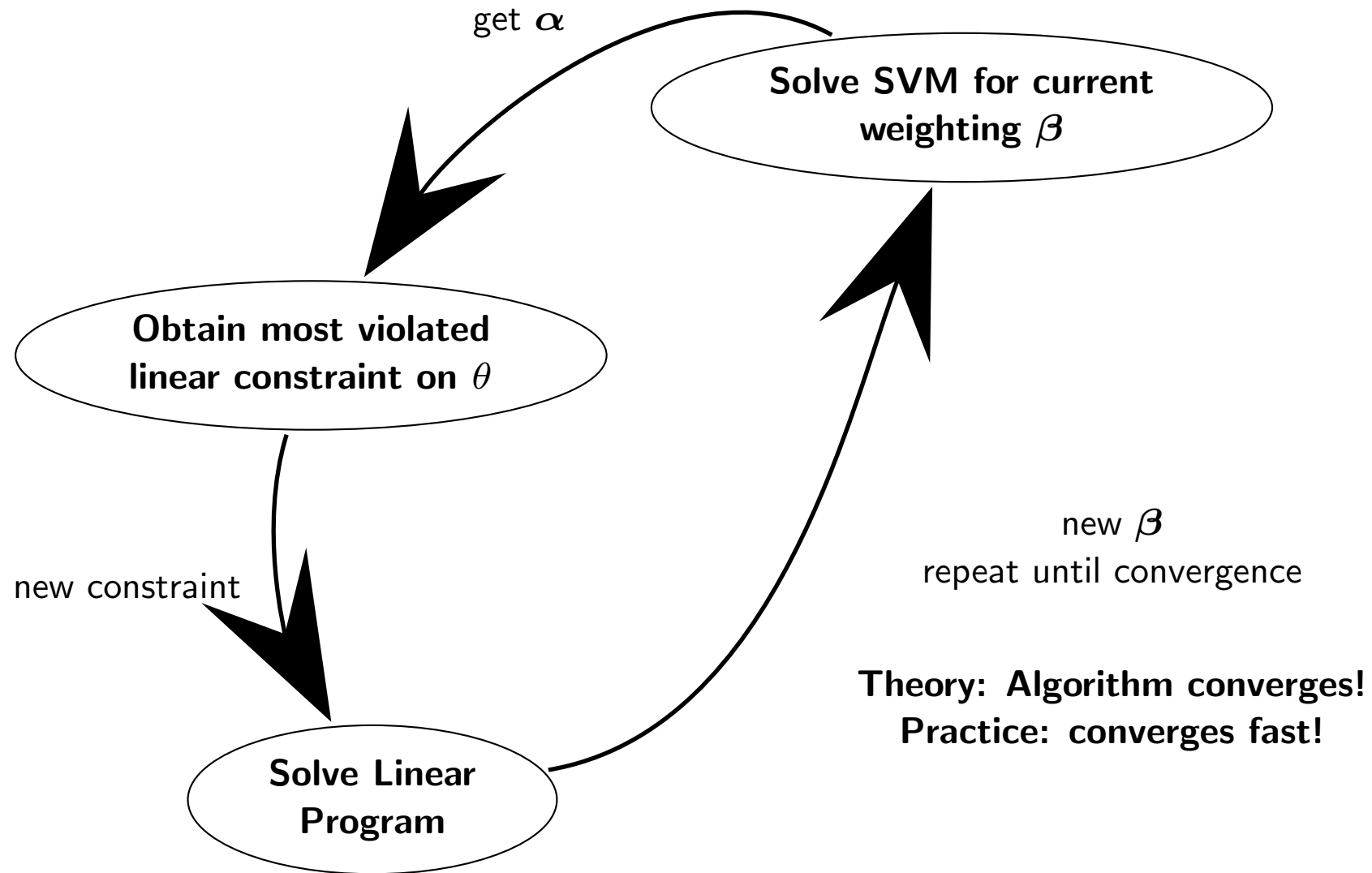
$$\begin{aligned}
 & \max && \theta \\
 & \text{w.r.t.} && \theta \in \mathbb{R}, \boldsymbol{\beta} \in \mathbb{R}_+^M \text{ with } \sum_{j=1}^M \beta_j = 1 \\
 & \text{s.t.} && \sum_{j=1}^M \beta_j \left(\frac{1}{2} S_j(\boldsymbol{\alpha}) - \sum_{i=1}^N \alpha_i \right) \geq \theta \\
 & && \text{for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

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

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

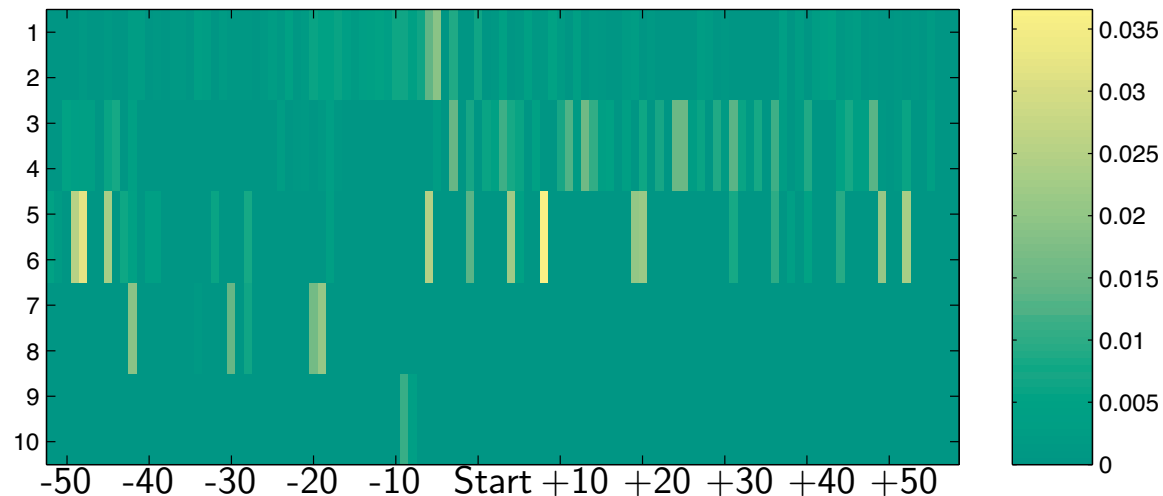
- iteratively find most violated constraints, solve linear program with current constraints, . . . , till convergence to the global optimum

COLUMN GENERATION II



RESULTS ON SPLICE SITES

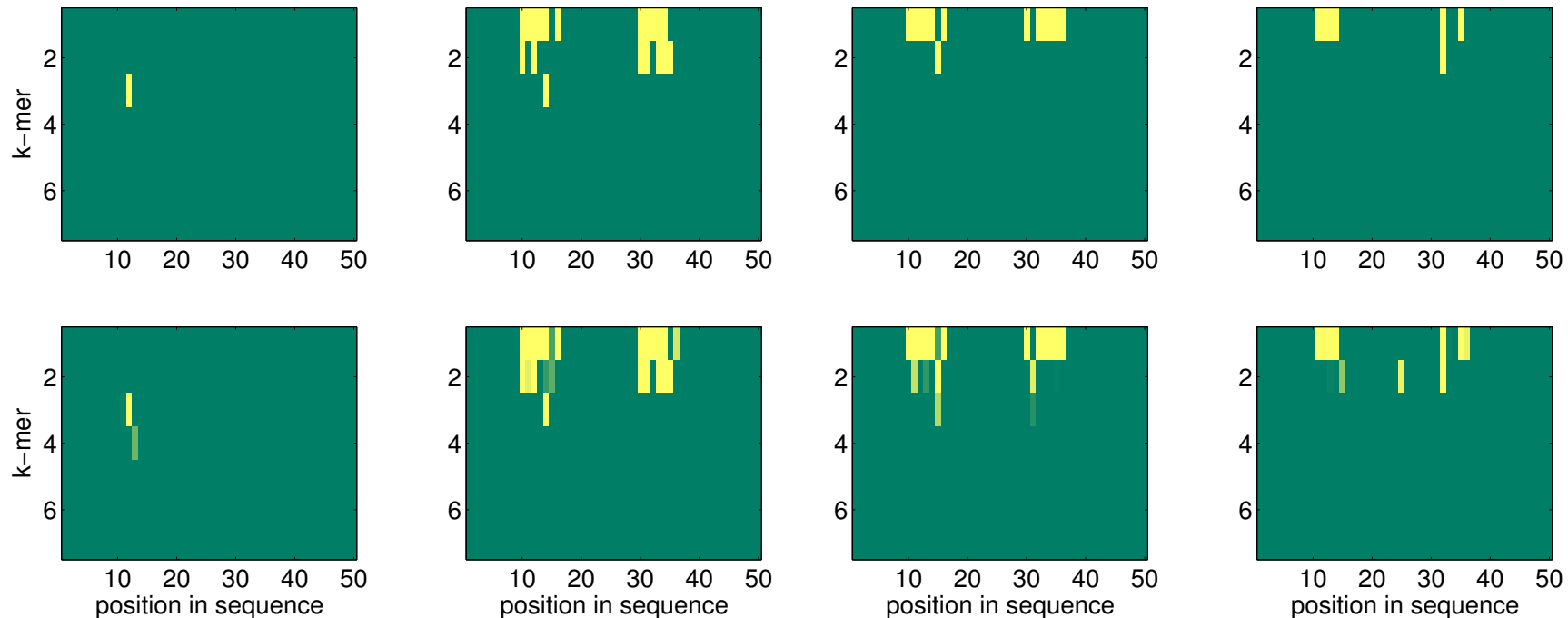
What characterizes the positions:



- Can we use that weighting to interpret the SVM solution? \Rightarrow **Not yet!**
 - Stability of the weighting β ?
 - Which weights are significant?

Use a statistical test to investigate significance of weights

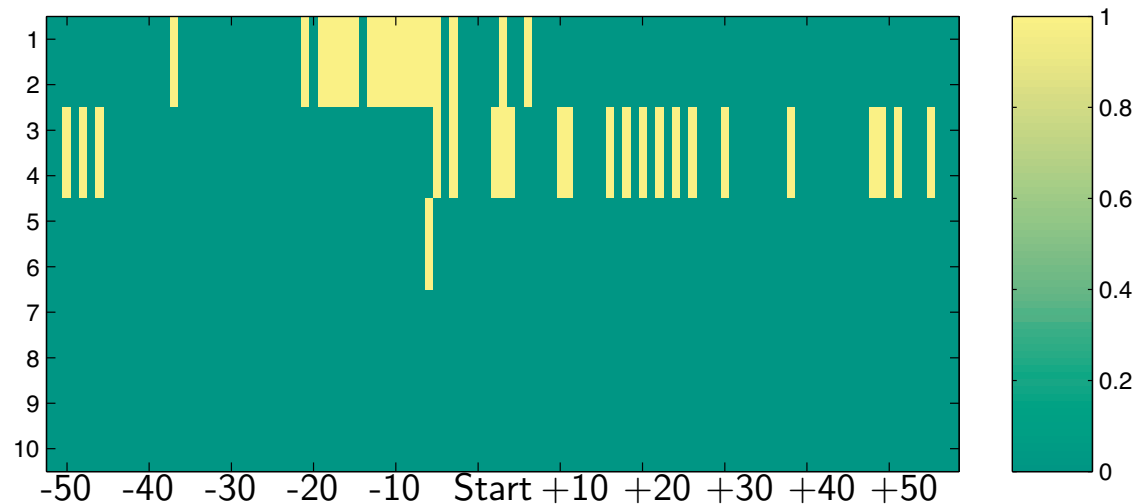
TOY DATASET



- DNA sequences of length 50 with hidden motifs at 10-16 and 30-36
- 8×50 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

RESULTS ON SPLICE SITES

What characterizes the positions:



- Can we use that weighting to interpret the SVM solution? \Rightarrow **NOW !**
 - Stability of the weighting β ?
 - Which weights are significant?

Use a statistical test to investigate significance of weights

CONCLUSION

Conclusion:

- MKL learns convex combination of kernels
 - ⇒ allows for interpreting SVM result
 - ⇒ matches prior knowledge about splice sites
- simple iterative algorithm
- suitable for large scale problems ($> 100,000$ examples)

Discussion:

- Can we improve classification using MKL ?
 - ⇒ $\|\cdot\|_1 = 1$ good choice ?
- Does MKL overfit and if so when ?
 - ⇒ How to regularize complexity ?
- Can we do model selection via MKL ?