

Topics I am working on V0.01-2004-04-14

Soeren.Sonnenburg@first.fraunhofer.de



Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

ROADMAP

- how to spend time in the addnet project
 - getting any info out of the $9 + 2 + 2$ microarrays
 - proceeding on the sequence level of ≈ 17000 genes
 - combine expression with sequence level information
- Kernels for splice site recognition
 - TOP/Fisher Kernel
 - Histogram Kernel
 - Salzberg Kernel
 - Polynomial Kernel
 - Locality Improved (2 versions)
 - Weighted Degree
 - Weighted Degree Position
 - Spectrum Kernel

ADDNET

- getting any info out of the $9 + 2 + 2$ microarrays
 - Combining Mouse and Human Microarrays
 - * find Mouse sequences in Human Microarrays
 - * then use matching dimensions
 - Data normalization (what about BLANKs)

$$\min_{\alpha} \sum_{i,j} (\log(\cosh(\alpha_i \text{array}(i) - \alpha_j \text{array}(j))))$$

- proceeding on the sequence level of ≈ 17000 genes
 - find promoters and then TFBS in promoters ?
 - compute alignment scores ?
- combine expression level information with dna

KERNEL MACHINES

Kernel Machine, e.g. SVM

- sequences $\mathbf{x}_i \in \mathbb{X}$ ($i = 1, \dots, \ell$) with respective labels y_i
- SVM classifier:

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

Kernel is **key** ingredient

- tune to incorporate prior knowledge
- determines feature space

SPLICE DATA

```
AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG
AAGATTAAAAAAAACAAATTTTTAGCATTACAGATATAATAATCTAATT
CACTCCCAAATCAACGATATTTTAGTTCACTAACACATCCGTCTGTGCC
TTAATTTCACTTCCACATACTTCCAGATCATCAATCTCCAAAACCAACAC
TTGTTTTAATATTCAATTTTTTACAGTAAGTTGCCAATTCAATGTTCCAC
TACCTAATTATGAAATTAAAATTCAGTGTGCTGATGGAAACGGAGAAGTC
```

- conserved sequence, position dependent
- mutation
- insertion
- deletion
- moved

TOP/FISHER/HISTOGRAM KERNEL

- Kernels from generative models
 - compare objects using a generative model $\Pr(\mathbf{x}|\Theta)$
 - exploit probabilistic model for **discriminative** training

Fisher Kernel (Jaakkola and Haussler, 1998)

$$\mathbf{k}_{\text{FK}}(\mathbf{x}, \mathbf{x}') = \mathbf{s}(\mathbf{x}, \hat{\boldsymbol{\theta}})^\top \mathbf{Z}^{-1}(\hat{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}', \hat{\boldsymbol{\theta}})$$

$$\mathbf{s}(\mathbf{x}, \hat{\boldsymbol{\theta}}) = \nabla_{\boldsymbol{\theta}} \log \Pr(\mathbf{x}|\hat{\boldsymbol{\theta}}) \quad \text{Fisher score vector}$$

$$\mathbf{Z}(\hat{\boldsymbol{\theta}}) = \mathbb{E}_{\mathbf{x}} \left[\mathbf{s}(\mathbf{x}, \hat{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}, \hat{\boldsymbol{\theta}})^\top \mid \hat{\boldsymbol{\theta}} \right] \quad \text{Fisher information matrix}$$

Kernel (Tsuda et. al, 2002)

$$\mathbf{k}_{\text{TOP}}(\mathbf{x}, \mathbf{x}') = \mathbf{f}_{\hat{\boldsymbol{\theta}}}(\mathbf{x})^\top \mathbf{f}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}')$$

$$\mathbf{f}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) = (v(\mathbf{x}, \hat{\boldsymbol{\theta}}), \nabla_{\boldsymbol{\theta}} v(\mathbf{x}, \hat{\boldsymbol{\theta}}))^\top$$

$$v(\mathbf{x}, \hat{\boldsymbol{\theta}}) = \log(\Pr(y = +1|\mathbf{x}, \hat{\boldsymbol{\theta}})) - \log(\Pr(y = -1|\mathbf{x}, \hat{\boldsymbol{\theta}}))$$

HISTOGRAM KERNEL

Histogram Kernel

- it is just a TOP Kernel from a Markov chain
- very fast compared with TOP from HMM

Computing partial derivatives becomes counting matches weighted with a constant

$$\partial_{\theta_{i,j}} v(\mathbf{x}, \boldsymbol{\theta}) = \mathbb{I}(x_i = j) / \theta_{i,j}$$

since the probabilistic model is just

$$\begin{aligned} P(\mathbf{x} | \boldsymbol{\theta}^{\pm}) &= P(x_1, \dots, x_N | \boldsymbol{\theta}^{\pm}) \\ &= P(x_1, \dots, x_{\omega} | \boldsymbol{\theta}^{\pm}) \prod_{i=\omega+1}^N P(x_i | x_{i-1}, \dots, x_{i-\omega}, \boldsymbol{\theta}^{\pm}). \end{aligned}$$

KERNELS - POLYNOMIAL KERNEL

Polynomial Kernel of degree d :

$$k_{\text{POLY}}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}, \mathbf{x}')^d$$

- not directly applicable
- common approach convert input into binary space, or equivalently count matches

$$k_{\text{POLY}}(\mathbf{x}, \mathbf{x}') = \left(\sum_{p=1}^N l_p(\mathbf{x}, \mathbf{x}') \right)^d$$

⇒ Feature space all d -th order monomials, i.e.

$$x_{i_1} \cdot x_{i_2} \cdot \dots \cdot x_{i_d}, \quad i_{1\dots d} \in 1 \dots N$$

⇒ Ultra high dimensional, *not robust*

KERNELS - LOCALITY IMPROVED KERNEL

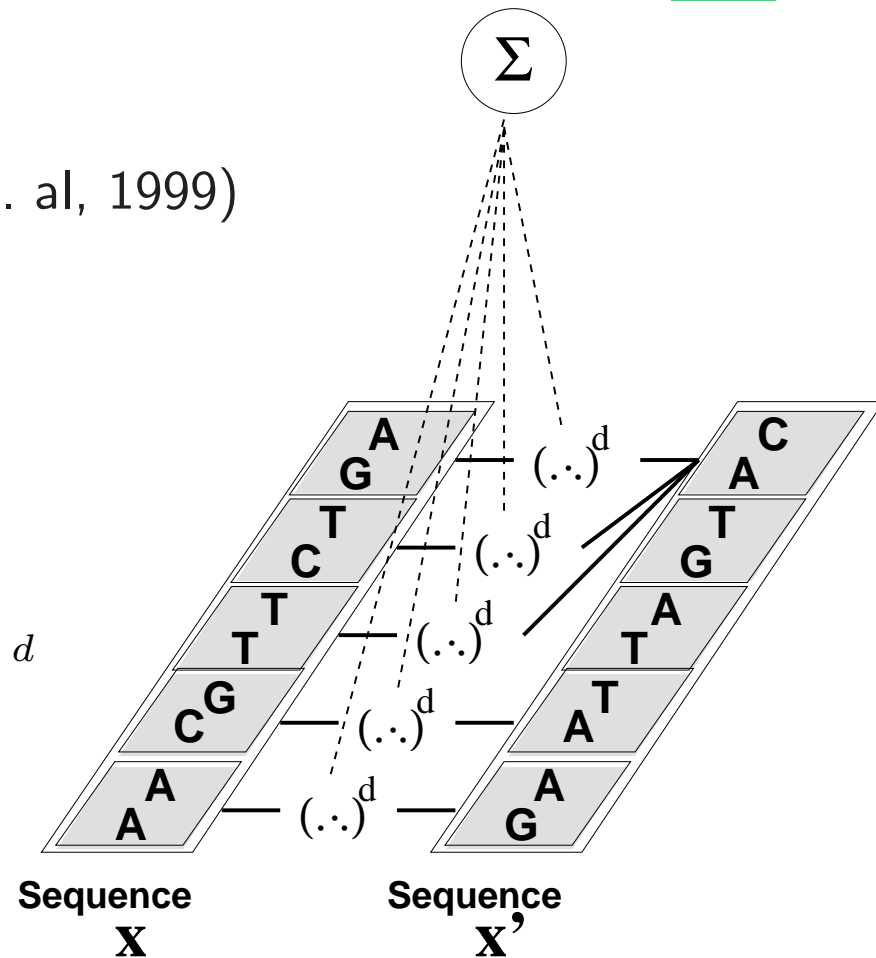
Locality Improved Kernel (Zien et. al, 1999)

- essentially polynomial kernel which tries to emphasize local information

$$k_{LI}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^N \text{win}_p(\mathbf{x}, \mathbf{x}')$$

$$\text{win}_p(\mathbf{x}, \mathbf{x}') = \left(\sum_{j=-l}^{+l} p_j I_{p+j}(\mathbf{x}, \mathbf{x}') \right)^d$$

$$I_i(\mathbf{x}, \mathbf{x}') = \begin{cases} 1, & x_i = x'_i \\ 0, & \text{otherwise} \end{cases}$$



⇒ in the end, counts matches of $2l + 1$ -tuples

KERNELS - SIMPLE LOCALITY IMPROVED KERNEL

Simple Locality Improved Kernel

- different weighting scheme, weight per window

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \sum_{p=l+1}^{N-l} w_p \text{win}_p(\mathbf{x}, \mathbf{x}')$$

with

$$\text{win}_p(\mathbf{x}, \mathbf{x}') = \left(\frac{1}{2l+1} \sum_{j=-l}^{+l} \mathbb{I}(x_{p+j} = x'_{p+j}) \right)^d$$

WEIGHTED DEGREE KERNEL

- count matches between two sequences \mathbf{x} and \mathbf{x}' between the words $\mathbf{u}_{\omega,i}(\mathbf{x})$ and $\mathbf{u}_{\omega,i}(\mathbf{x}')$
- $\mathbf{u}_{\omega,i}(\mathbf{x}) = x_i x_{i+1} \dots x_{i+\omega-1}$ for all i and $1 \leq \omega \leq d$.
- ω denotes the order (length of the word) to be compared.

Weighted Degree Kernel

$$k_{WD}(\mathbf{x}, \mathbf{x}') = \sum_{\omega=1}^d w_{\omega} \sum_{i=1}^{N-d} \mathbb{I}(\mathbf{u}_{\omega,i}(\mathbf{x}) = \mathbf{u}_{\omega,i}(\mathbf{x}'))$$

- w_{ω} denotes the weight for a match of length ω
- mismatches are also incorporated, resulting in lower weighting of the imperfect match

Weighted Degree Position Kernel

$$k_{WDpos}(\mathbf{x}, \mathbf{x}') = \sum_{\omega=1}^d w_{\omega} \sum_{i=1}^{N-d-\text{shift}_i} \sum_{j=1}^{\text{shift}_i} \mathbb{I}(\mathbf{u}_{\omega, i+j}(\mathbf{x}) = \mathbf{u}_{\omega, i+j}(\mathbf{x}'))$$

Can cope with

- conserved sequence, position dependent
- mutation, insertion, deletion
- moved

SPECTRUM KERNEL

Spectrum Kernel (Leslie et. al, 2002):

$$k_{spec} = f(\mathbf{x})^\top f(\mathbf{x})$$

Inner product of

$$f(\mathbf{x}) = (\#A(\mathbf{x}))_{A \in \Sigma^k}$$

where $\#A(\mathbf{x}) =$ number of occurrences of A in \mathbf{x} .

- counts matches of k -mers
- no position dependend information