

Support Vector Machines and Multiple Kernel Learning

(lets do it large scale with shogun)

Sören Sonnenburg

Shogun Toolbox Data Days 2013, Juli 12, Berlin

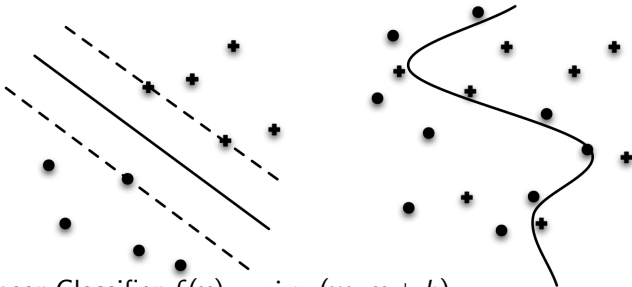
將軍

Outline

- 1 SVMs
- 2 Large Scale SVMs
- 3 Applications
- 4 MKL
- 5 MKL Extensions

Classification

Given training examples $(\mathbf{x}_i, y_i)_{i=1}^N \in (\mathcal{X}, \{-1, +1\})^N$



- Linear Classifier $f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$
- Kernel Machine (e.g. Support Vector Machine), learn weighting $\alpha \in \mathbb{R}^N$ on training examples in kernel feature space $\Phi(\mathbf{x})$

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

where **Kernel** $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$

Classification using Support Vector Machines

Via **kernel**: **non-linear** discrimination in input space

- learn weighting $\alpha \in \mathbb{R}^N$ on training examples $(\mathbf{x}_i, y_i)_{i=1}^N$ in kernel feature space $\Phi(\mathbf{x})$

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

- where **Kernel** $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$
- still linear classifier $f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$ in kernel feature space, with weighting $\mathbf{w} = \sum_{i=1}^N y_i \alpha_i \Phi(\mathbf{x}_i)$ and examples $\mathbf{x} \mapsto \Phi(\mathbf{x})$

SVMs rock!

- Kernels - flexible!
- In many applications SVMs define the state-of-the-art!

Classification using Support Vector Machines

Via **kernel**: **non-linear** discrimination in input space

- learn weighting $\alpha \in \mathbb{R}^N$ on training examples $(\mathbf{x}_i, y_i)_{i=1}^N$ in kernel feature space $\Phi(\mathbf{x})$

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

- where **Kernel** $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$
- still linear classifier $f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$ in kernel feature space, with weighting $\mathbf{w} = \sum_{i=1}^N y_i \alpha_i \Phi(\mathbf{x}_i)$ and examples $\mathbf{x} \mapsto \Phi(\mathbf{x})$

SVMs rock!

- Kernels - flexible!
- In many applications SVMs define the state-of-the-art!

Support Vector Machines in Shogun I

Implementations supporting Kernels

- LibSVM, LibSVMOneClass, MulticlassLibSVM, LibSVR (sequential minimal optimization)
- GNPPSVM (ℓ_2 -SVM)
- GPBTSVM (chunking algorithm based SVM)
- GMNPSVM (true multi-class SVM)
- LaRank (true multi-class SVM)
- MPDSVM (minimal primal-dual SVM)
- SVMLight, SVMLightOneClass, SVRLight (chunking based, thread-parallel non-GPL)
- ScatterSVM (multi-class approximation)

Support Vector Machines in Shogun I

Implementations supporting Kernels

- **LibSVM**, **LibSVMOneClass**, **MulticlassLibSVM**, **LibSVR** (sequential minimal optimization)
- **GNPPSVM** (ℓ_2 -SVM)
- **GPBTSVM** (chunking algorithm based SVM)
- **GMNPSVM** (true multi-class SVM)
- **LaRank** (true multi-class SVM)
- **MPDSVM** (minimal primal-dual SVM)
- **SVMLight**, **SVMLightOneClass**, **SVRLight** (chunking based, thread-parallel non-GPL)
- **ScatterSVM** (multi-class approximation)

Kernels

Popular Kernels

- Gaussian, Polynomial, String Kernels
- Custom Kernel (bring your own)

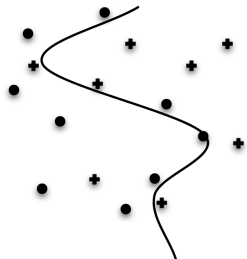
Kernel Normalizers

- SqrtDiagKernelNormalizer (points on hypersphere with $R = 1$)
- ZeroMeanCenterKernelNormalizer (centered in feature space)
- AvgDiagKernelNormalizer, more <http://bit.ly/133aG2V>

Most comprehensive Kernel/SVM Framework

- Over 60 kernels available
- <http://bit.ly/15z1CWY>
- normalizers can be attached to kernels

Kernel based Support Vector Machines



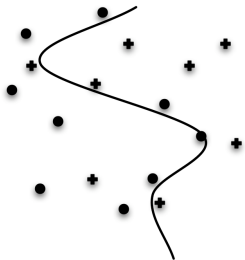
- SVMs learn weights $\alpha \in \mathbb{R}^m$ over training examples in kernel feature space $\Phi : \mathbf{x} \mapsto \mathbb{R}^n$

- Decision function $f(\mathbf{x}) = \text{sign}(\sum_{i=1}^m y_i \alpha_i k(\mathbf{x}, \mathbf{x}_i) + b)$,
with kernel $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$

SVMs rock!

- Kernels - flexible!
- In many applications SVMs define the state-of-the-art!
- But not large scale!

Kernel based Support Vector Machines



- SVMs learn weights $\alpha \in \mathbb{R}^m$ over training examples in kernel feature space $\Phi : \mathbf{x} \mapsto \mathbb{R}^n$
- Decision function $f(\mathbf{x}) = \text{sign}(\sum_{i=1}^m y_i \alpha_i k(\mathbf{x}, \mathbf{x}_i) + b)$,
with kernel $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$

SVMs rock!

- Kernels - flexible!
- In many applications SVMs define the state-of-the-art!
- **But not large scale!**

Kernel based Support Vector Machines are a Dead End

The Curse of Support Vectors

To compute output on all m examples $\mathbf{x}_1, \dots, \mathbf{x}_m$:

$$\forall j = 1, \dots, m : \sum_{i=1}^{m_s} \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}_j) + b$$

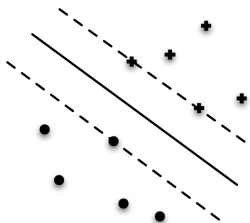
Computational effort:

- All $\mathcal{O}(m_s m T)$, (T time to compute the kernel)
- Effort Scales linearly with $m_s = \mathcal{O}(m) := \#SVs$

⇒ **SVM's in bigO are not faster than standard k-NN.**

⇒ **Kernel Machines are just not large-scale!**

What about Linear SVMs ?



- **Linear** Support Vector Machines learn weights $\mathbf{w} \in \mathbb{R}^n$
- Decision function $\mathbf{f}(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$

Recent Progress in Linear SVM solvers

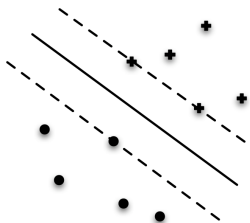
- SGD (Bottou 2007), SGD-QN (Bordes et al., 2009)
- SVM^{perf} (Joachims 2006), liblinear (Fan et al. 2008)
- BMRM (Teo et.al. 2007), OCAS (Franc, Sonnenburg 2009)

⇒ **Linear** training Effort $\mathcal{O}(m)$

⇒ Computing Outputs **Linear** Effort $\mathcal{O}(nm)$

... already linear time but **just linear**

What about Linear SVMs ?



- **Linear** Support Vector Machines learn weights $\mathbf{w} \in \mathbb{R}^n$
- Decision function $\mathbf{f}(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$

Recent Progress in Linear SVM solvers

- SGD (Bottou 2007), SGD-QN (Bordes et al., 2009)
- SVM^{perf} (Joachims 2006), liblinear (Fan et al. 2008)
- BMRM (Teo et.al. 2007), OCAS (Franc, Sonnenburg 2009)

⇒ **Linear** training Effort $\mathcal{O}(m)$

⇒ Computing Outputs **Linear** Effort $\mathcal{O}(nm)$

... already linear time but **just linear**

Support Vector Machines in Shogun II

Linear SVMs

- OnlineLibLinear, LibLinear, MulticlassLibLinear (dual coordinate descent)
- NewtonSVM (newton primal svm)
- OnlineSVMMSGD, SVMMSGD (stochastic gradient descent)
- SGDQN (stochastic gradient descent with quasi-newton steps)
- SVMLin
- SVMOCas, MulticlassOCAS (cutting plane/bundle-method)

Support Vector Machines in Shogun II

Linear SVMs

- **OnlineLibLinear**, **LibLinear**, **MulticlassLibLinear** (dual coordinate descent)
- **NewtonSVM** (newton primal svm)
- **OnlineSVMMSGD**, **SVMMSGD** (stochastic gradient descent)
- **SGDQN** (stochastic gradient descent with quasi-newton steps)
- **SVMLin**
- **SVMOCas**, **MulticlassOCAS** (cutting plane / bundle-method)

Shogun's computational framework for linear SVMs

AIM: Development of a large scale learning framework for SVMs

*"Algorithm [linear SVM solver] improvements do not improve the order of test error convergence. They can simply improve constant factors and therefore **compete evenly with the implementation improvements**. Time spent refining the implementation is time well spent."*

from: Bordes, Bottou, Gallinari: SQD-QN: Careful Quasi-Newton Stochastic Gradient Descent. JMLR 2009.

Towards a computational framework for linear SVMs

Linear SVM solvers like liblinear, SGD, BMRM, OCAS only require two operations to access data:

- (i) dot product between feature vector and the vector \mathbf{w} :

$$r \leftarrow \langle \mathbf{x}, \mathbf{w} \rangle$$

DOT

- (ii) multiplication with scalar $\alpha \in \mathbb{R}$ and addition to vector $\mathbf{v} \in \mathbb{R}^n$

$$\mathbf{v} \leftarrow \alpha \mathbf{x} + \mathbf{v}$$

ADD

COFFIN

COFFIN really is just two simple ideas:

On demand compute...

- ① Features $\Phi(\mathbf{x})$ (only non-zero dims)
 - Non-Linearity Possible
 - Examples: Low Degree Polynomial Kernel, Spectrum Kernel, Weighted Degree Kernel
 - On-the-fly (de)compression
- ② Virtual Examples
 - Incorporating Invariances possible
 - Examples: Image translation, rotation, etc

Needs efficient data structure for \mathbf{w} !

- ..., dense, sorted array, trees, hashes
- fast only when $|\Phi_{\neq 0}(\mathbf{z})| \sim \dim(\mathbf{z})$

COFFIN

COFFIN really is just two simple ideas:

On demand compute...

- ① Features $\Phi(\mathbf{x})$ (only non-zero dims)
 - Non-Linearity Possible
 - Examples: Low Degree Polynomial Kernel, Spectrum Kernel, Weighted Degree Kernel
 - On-the-fly (de)compression
- ② Virtual Examples
 - Incorporating Invariances possible
 - Examples: Image translation, rotation, etc

Needs efficient data structure for \mathbf{w} !

- ..., dense, sorted array, trees, hashes
- fast only when $|\Phi_{\neq 0}(\mathbf{z})| \sim \dim(\mathbf{z})$

Data Structures

Effort of **ADD** and **DOT** for **z** and memory requirement of **w**.

	Dense	Sorted Array	Tree
Add	$\mathcal{O}(\Phi_{\neq 0}(\mathbf{z}))$	$\mathcal{O}(\mathbf{w} _{\neq 0} + \Phi_{\neq 0}(\mathbf{z}))$	$\mathcal{O}(\Phi_{\neq 0}(\mathbf{z}))$ to $\mathcal{O}(K \Phi_{\neq 0}(\mathbf{z}))$
Dot	$\mathcal{O}(\Phi_{\neq 0}(\mathbf{z}))$	$\mathcal{O}(\mathbf{w} _{\neq 0} + \Phi_{\neq 0}(\mathbf{z}))$	$\mathcal{O}(\Phi_{\neq 0}(\mathbf{z}))$ to $\mathcal{O}(K \Phi_{\neq 0}(\mathbf{z}))$
Mem	$\mathcal{O}(n)$	$\mathcal{O}(\sum_{i=1}^m \Phi_{\neq 0}(\mathbf{z}_i))$	$\mathcal{O}(\sum_{i=1}^m \Phi_{\neq 0}(\mathbf{z}_i))$

Sparse data structures have huge overhead!

Hashing to the Rescue (Shi et al (2009))

- Always use *dense* **w** with “compressed index”
- Hash function $h(J) \mapsto 1, \dots, 2^\gamma$,
- $(\hat{\Phi}(\mathbf{z}))_j = \sum_{i \in J; h(i)=j} (\Phi(\mathbf{z}))_i$

Data Structures

Effort of **ADD** and **DOT** for **z** and memory requirement of **w**.

	Dense	Sorted Array	Tree
Add	$\mathcal{O}(\Phi_{\neq 0}(\mathbf{z}))$	$\mathcal{O}(\mathbf{w} _{\neq 0} + \Phi_{\neq 0}(\mathbf{z}))$	$\mathcal{O}(\Phi_{\neq 0}(\mathbf{z}))$ to $\mathcal{O}(K \Phi_{\neq 0}(\mathbf{z}))$
Dot	$\mathcal{O}(\Phi_{\neq 0}(\mathbf{z}))$	$\mathcal{O}(\mathbf{w} _{\neq 0} + \Phi_{\neq 0}(\mathbf{z}))$	$\mathcal{O}(\Phi_{\neq 0}(\mathbf{z}))$ to $\mathcal{O}(K \Phi_{\neq 0}(\mathbf{z}))$
Mem	$\mathcal{O}(n)$	$\mathcal{O}(\sum_{i=1}^m \Phi_{\neq 0}(\mathbf{z}_i))$	$\mathcal{O}(\sum_{i=1}^m \Phi_{\neq 0}(\mathbf{z}_i))$

Sparse data structures have huge overhead!

Hashing to the Rescue (Shi et al (2009))

- Always use *dense* **w** with “compressed index”
- Hash function $h(J) \mapsto 1, \dots, 2^\gamma$,
- $(\hat{\Phi}(\mathbf{z}))_j = \sum_{i \in J; h(i)=j} (\Phi(\mathbf{z}))_i$

COFFIN is implemented in Shogun's DotFeatures

DotFeatures

- Dense-, SparseFeatures
- PolyFeatures, BinnedDotFeatures
- CombinedDotFeatures
- ExplicitSpecFeatures, ImplicitWeightedSpecFeatures
- HashedDocDotFeatures

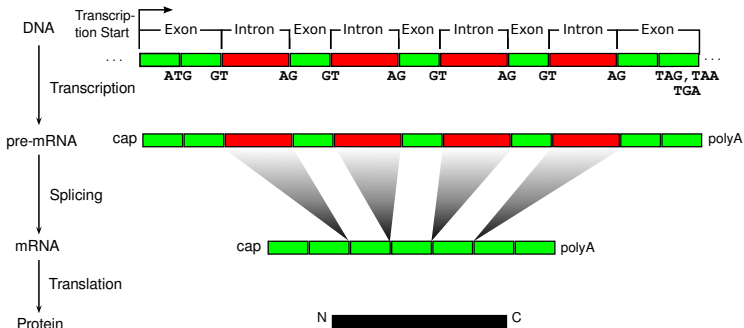
You can

- Stack together features (dense, sparse,...)
- Provides abstraction to data representation

Use it!

Implemented for all linear SVMs

Splice Site Predictions



Application to Human Acceptor Splice Site Prediction

Splice Site Prediction

Discriminate true signal positions against all other positions

≈ 150 nucleotides window around dimer

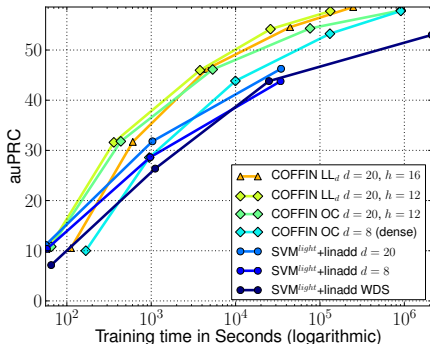
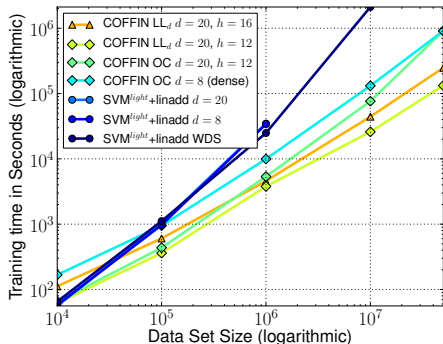
CT...GTCGTA...GAAGCTAGGAGCGC...ACGCGT...GA

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

```
AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG
AAGATTAATAAAAAAAAAACAAATTTTTCATTACAGATATAATAATCTAATT
CACTCCCCAAATCAACGATATTTTAGTTCACCTAACACATCCGTCTGTGCC
TTAATTTCACTTCCACATACTTCCAGATCATCAATCTCCAAAACCAACAC
TTGTTTTAATATTCAATTTTTTACAGTAAGTTGCCAATTCAATGTTCCAC
TACCTAATTATGAAATTAAATTTCAGTGTGCTGATGGAAACGGAGAAGTC
```

- 50 million training examples
- COFFIN with kernels: weighted spectrum and weighted degree (explicit and hashed representation) $\approx 200,000,000$ dims

Applications

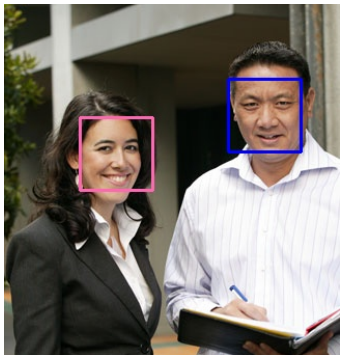


It's fast and works!

- Factor 47 faster on $10 \cdot 10^6$ examples than linadd
- New state-of-the-art results auPRC 58.57% vs. 53.01%

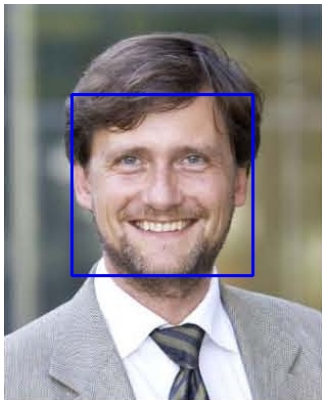
Gender Classification

Distinguish Females from Males solely based on Faces



- learn COFFIN on labelled faces
- virtual examples: translation, rotation, scale
- train ≈ 5 million sample (that would require 50GB) on Vojtechs notebook

Gender Classification I



It's fast and works - again!

- auROC 95.44%
- (vs. auROC 89.57% without VE)

Klaus-Robert Müller is a male!

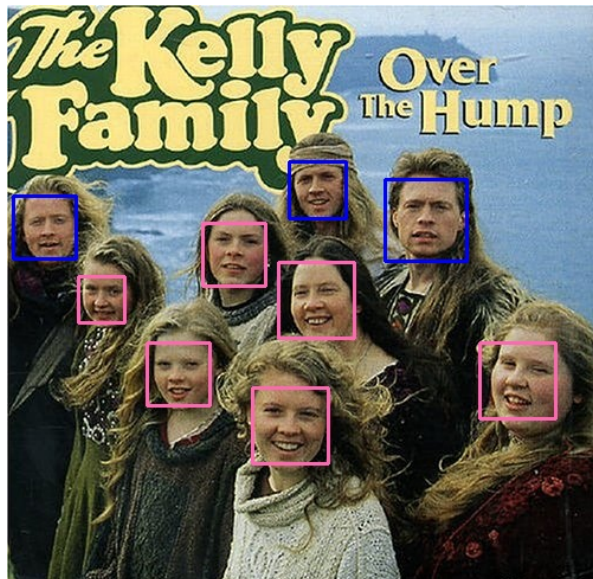
Gender Classification II



Gender Classification III



Gender Classification IV



Summary Support Vector Machines

Kernel SVMs

- Allows non-linear and complex models
- Applicable to mid-size datasets (1e6)
- General and often state-of-the art detectors

Linear SVMs

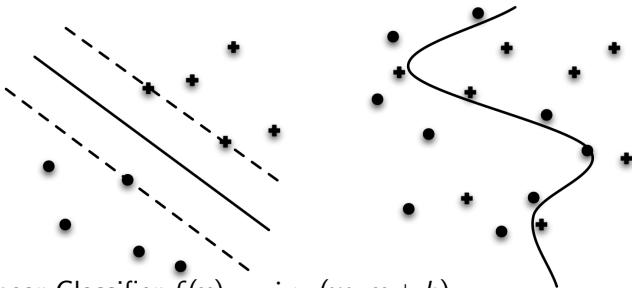
- Allows non-linearity with COFFIN/DotFeatures
- Applicable to huge datasets
- General and often state-of-the art detectors

Datasets, Scripts, Efficient implementation

- Data and Scripts <http://sonnenburgs.de/soeren/coffin>
- Implementation <http://www.shogun-toolbox.org>
- More machine learning software <http://mloss.org>

Multiple Kernel Learning

Given training examples $(\mathbf{x}_i, y_i)_{i=1}^N \in (\mathcal{X}, \{-1, +1\})^N$



- Linear Classifier $f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$
- Kernel Machine (e.g. Support Vector Machine), learn weighting $\alpha \in \mathbb{R}^N$ on training examples in kernel feature space $\Phi(\mathbf{x})$

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

where **Kernel** $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$

Classification using Kernel Machines I

Single Kernel

- Kernel Machine (e.g. Support Vector Machine)
 - learn weighting $\alpha \in \mathbb{R}^N$ on training examples $(\mathbf{x}_i, y_i)_{i=1}^N$ in kernel feature space $\Phi(\mathbf{x})$

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

- where **Kernel** $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$
- still linear classifier $f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$ in kernel feature space, with weighting $\mathbf{w} = \sum_{i=1}^N y_i \alpha_i \Phi(\mathbf{x}_i)$ and examples $\mathbf{x} \mapsto \Phi(\mathbf{x})$

via **kernel: non-linear** discrimination in input space

Classification using Kernel Machines II

Multiple Kernels

- Linear combination of kernels $k(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^M \beta_j k_j(\mathbf{x}, \mathbf{x}')$, $\beta_j \geq 0$. Learn α and β . Resulting classifier:

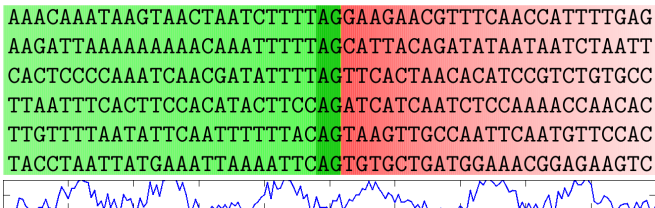
$$f(\mathbf{x}) = \text{sign} \left(\sum_{j=1}^M \beta_j \sum_{i=1}^N y_i \alpha_i k_j(\mathbf{x}, \mathbf{x}_i) + b \right)$$

Learn weighting over training examples α and kernels β

When is Multiple Kernel Learning useful?

Combining Heterogeneous Data

- Consider data from different domains: e.g. Bioinformatics features: DNA-strings, binding energies, conservation, structure,...



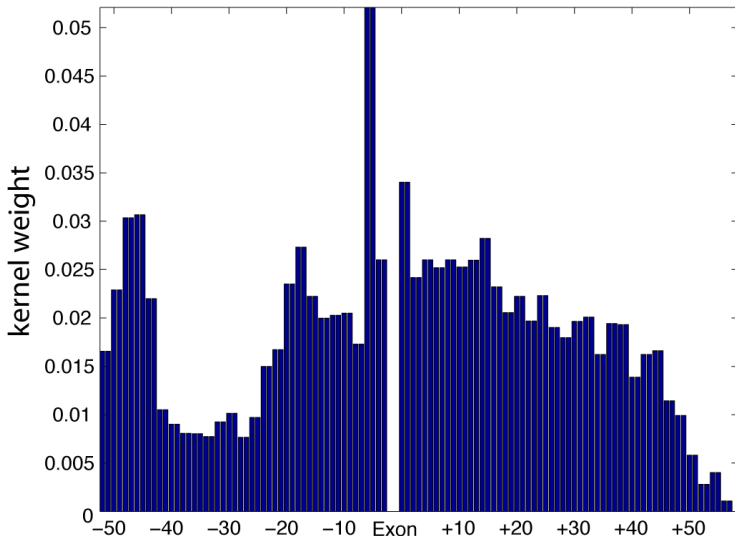
$$k(\mathbf{x}, \mathbf{x}') =$$

$$\beta_1 k_{dna}(\mathbf{x}_{dna}, \mathbf{x}'_{dna}) + \beta_2 k_{nrg}(\mathbf{x}_{nrg}, \mathbf{x}'_{nrg}) + \beta_3 k_{3d}(\mathbf{x}_{3d}, \mathbf{x}'_{3d}) + \dots$$

When is Multiple Kernel Learning useful?

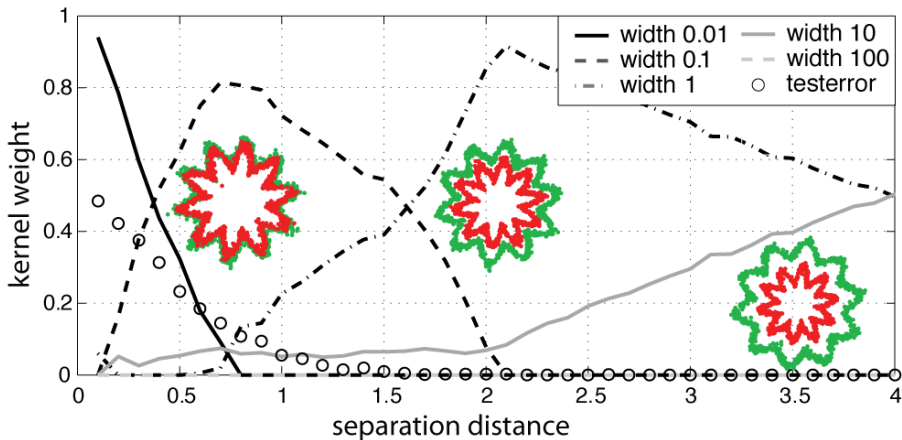
Interpretability

- Bioinformatics: One weight per position in sequence

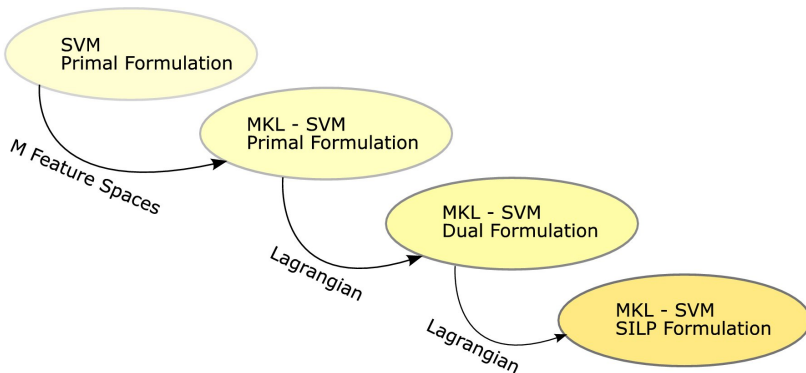


When is Multiple Kernel Learning useful?

Automated Model Selection



Derivation



For details see Sonnenburg, Rätsch, Schäfer, Schölkopf 2006

SVM Primal Formulation

$$\begin{array}{ll}\min & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i \\ \text{w.r.t.} & \mathbf{w} \in \mathbb{R}^D, \boldsymbol{\xi} \in \mathbb{R}_+^N, b \in \mathbb{R} \\ \text{s.t.} & y_i \left(\mathbf{w}^\top \Phi(\mathbf{x}_i) + b \right) \geq 1 - \xi_i, \forall i = 1, \dots, N\end{array}$$

MKL Primal Formulation

$$\begin{aligned}
 \min \quad & \frac{1}{2} \left(\sum_{j=1}^M \beta_j \|\mathbf{w}_j\|_2 \right)^2 + C \sum_{i=1}^N \xi_i \\
 \text{w.r.t.} \quad & \mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_M), \mathbf{w}_j \in \mathbb{R}^{D_j}, \quad \forall j = 1 \dots M \\
 & \beta \in \mathbb{R}_+^M, \xi \in \mathbb{R}_+^N, b \in \mathbb{R} \\
 \text{s.t.} \quad & y_i \left(\sum_{j=1}^M \beta_j \mathbf{w}_j^\top \Phi_j(\mathbf{x}_i) + b \right) \geq 1 - \xi_i, \quad \forall i = 1, \dots, N \\
 & \sum_{j=1}^M \beta_j = 1
 \end{aligned}$$

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

MKL Dual Formulation

Bach, Lanckriet, Jordan 2004:

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

Properties: equivalent to SVM for $M = 1$

The Semi-Infinite Linear Program I

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

Properties:

- linear, optimize over a convex combination of β
- infinitely many constraints, one for each $0 \leq \alpha \leq C$
- can use standard SVM to identify most violated constraints

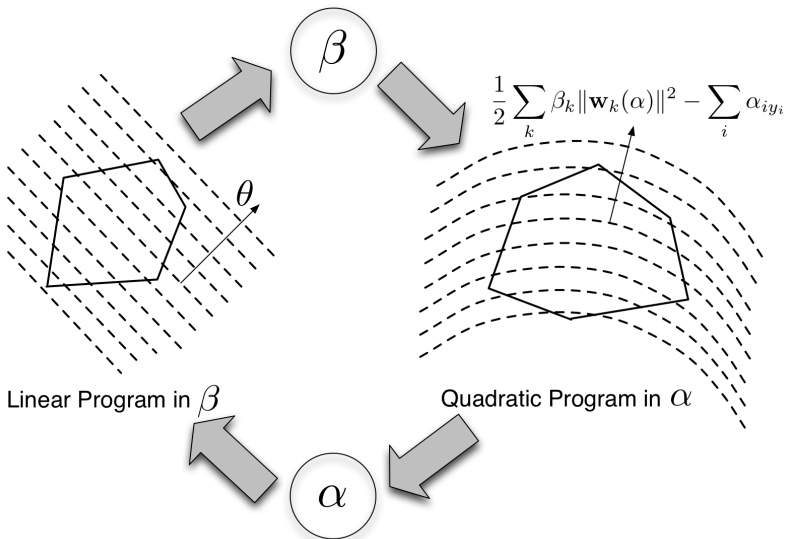
The Semi-Infinite Linear Program II

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

Solving the SILP:

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

Solving the SILP: Column Generation I



Solving the SILP: Column Generation II

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

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

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

- solved by taking set of most violated constraints into account
- most violated constraints given by SVM solution for fixed β

Regression

Primal Formulation:

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

$$\text{w.r.t.} \quad \mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_M), \mathbf{w}_j \in \mathbb{R}^{D_j}, \quad \forall j = 1 \dots M$$

$$\beta \in \mathbb{R}_+^M, \xi \in \mathbb{R}^N, \xi^* \in \mathbb{R}_+^N, b \in \mathbb{R}$$

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

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

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

One Class

Primal Formulation:

$$\begin{aligned}
 \min \quad & \frac{1}{2} \left(\sum_{j=1}^M \beta_j \|\mathbf{w}_j\|_2 \right)^2 + C \sum_{i=1}^N \xi_i - \rho \\
 \text{w.r.t.} \quad & \mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_M), \mathbf{w}_j \in \mathbb{R}^{D_j}, \quad \forall j = 1 \dots M \\
 & \beta \in \mathbb{R}_+^M, \xi \in \mathbb{R}_+^N \\
 \text{s.t.} \quad & y_i \left(\sum_{j=1}^M \beta_j \mathbf{w}_j^\top \Phi_j(\mathbf{x}_i) \right) \geq \rho - \xi_i, \forall i = 1, \dots, N \\
 & \sum_{j=1}^M \beta_j = 1
 \end{aligned}$$

Generalized for arbitrary strictly convex differentiable loss functions (Sonnenburg, Rätsch, Schäfer, Schölkopf 2006)

Multiclass

Primal Formulation (Zien, Ong 2007):

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

$$\text{w.r.t.} \quad \mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_M), \mathbf{w}_j \in \mathbb{R}^{k_j}, \quad \forall j = 1 \dots M$$

$$\beta \in \mathbb{R}_+^M, \mathbf{s} \in \mathbb{R}^{N \times c}, \boldsymbol{\xi} \in \mathbb{R}_+^N, \mathbf{b} \in \mathbb{R}$$

$$\text{s.t.} \quad \xi_i = \max_{u \neq y_i} s_{iu}, \quad s_{iu} \geq 0,$$

$$\sum_{j=1}^M \beta_j \mathbf{w}_j^\top (\Phi_j(\mathbf{x}_i, y_i) - \Phi_j(\mathbf{x}_i, u)) + b_{y_i} - b_u \geq 1 - s_{iu},$$

$$\forall i = 1 \dots N, \forall u = 1 \dots c$$

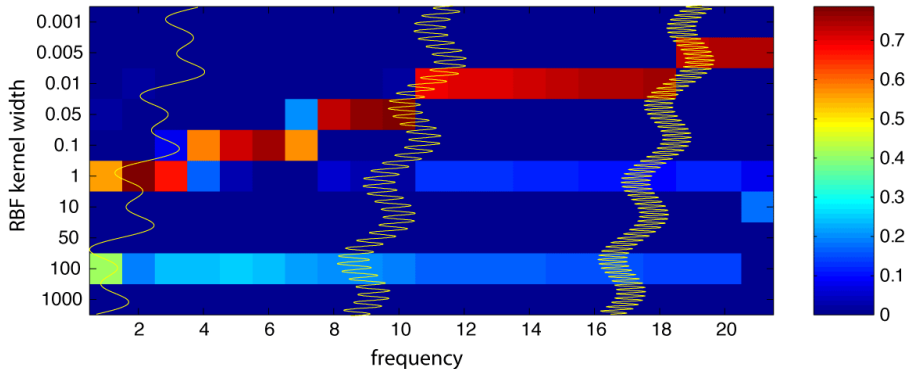
Further Extension

Extension

- ℓ_p norm MKL
- faster algorithms
- convergence bounds
- and many more

⋮

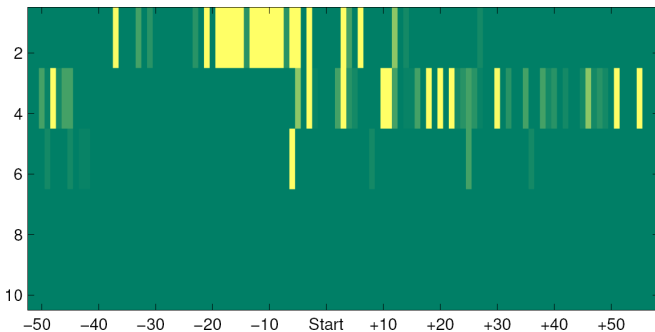
Automated Model Selection - Regression



- $f(x) = \sin(ax) + \sin(bx) + cx$ for varying a
- Support Vector Regression with 10 RBF-Kernels of different width

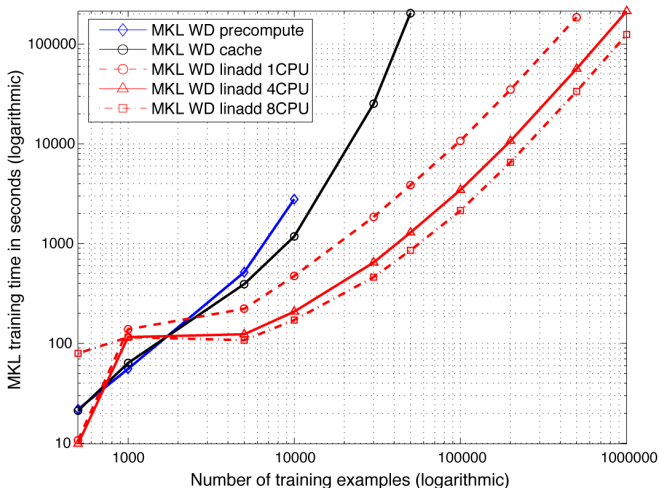
Knowledge discovery

Feature Extraction



- Support Vector Classification on Bioinformatics problem, distinguish “splice sites” from “fake sites” (aligned DNA sequences)
- One weight β_j per position and per sub-sequence length
- Displayed: Learned weights of 500 kernels

A word about speed.



Summary and Outlook

MKL learns convex combination of kernels

- ⇒ allows (to some extent) for automated model selection
- ⇒ allows for interpreting SVM result
- ⇒ matches prior knowledge on real-world bioinformatics problem

- **Simple:** iterative wrapper algorithm around single kernel SVM
- **General:** same technique applicable to a wide range of problems (1-class, 2-class, Multiclass, Regression, ...)
- **Fast:** suitable for large scale problems ($> 100,000$ examples)

Download free source <http://www.shogun-toolbox.org>.