

LARGE-SCALE KERNEL METHODS

MDI 341, MS BIG DATA, TÉLÉCOM PARISTECH

Aurélien Bellet

Inria Lille

March 20, 2017

1. Reminder on Kernel Methods
2. Scalability Issues
3. Random Kernel Features
4. Nyström Approximation
5. Conclusion

REMINDER ON KERNEL METHODS

- Training observations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X}$
- For now, we assume $\mathcal{X} \subset \mathbb{R}^p$
- Training labels $y_1, \dots, y_n \in \{-1, 1\}$
- Linear classifier with parameters $\mathbf{w} \in \mathbb{R}^p, b \in \mathbb{R}$:

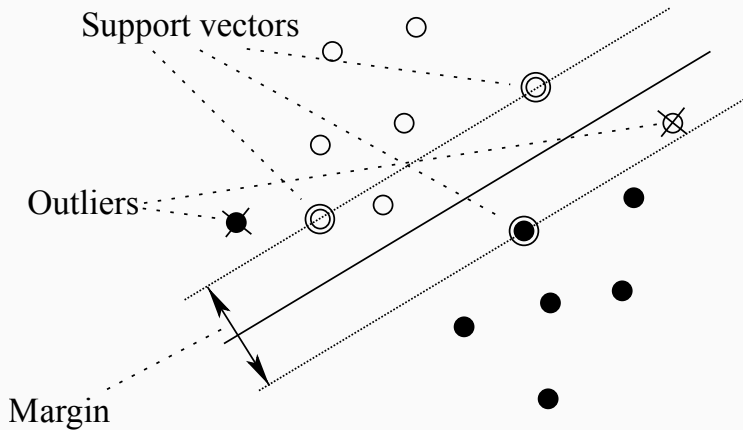
$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

SVM: primal formulation

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^p, b \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x} + b) \geq 1 - \xi_i \quad i = 1, \dots, n \\ & \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

- Key principle: margin maximization
- Convex optimization problem

SUPPORT VECTOR MACHINES: PRIMAL

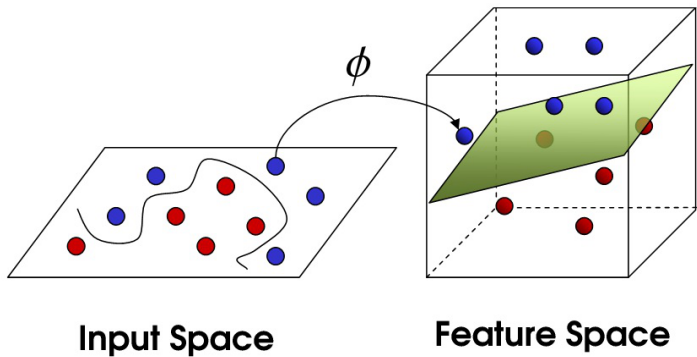


SVM: Lagrange dual formulation

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C \quad i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

- Also convex
- We have $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$ and thus

$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}^T \mathbf{x}_i \right)$$



Definition (Kernel function)

A symmetric function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a *kernel* if there exists a mapping function $\phi : \mathcal{X} \rightarrow \mathbb{H}$ from the instance space \mathcal{X} to a Hilbert space \mathbb{H} such that K can be written as an inner product in \mathbb{H} :

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle .$$

Equivalently, K is a *kernel* if it is symmetric positive semi-definite (PSD), i.e.,

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

for all finite sequences of $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X}$ and $c_1, \dots, c_n \in \mathbb{R}$.

Kernel SVM formulation

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C \quad i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

- We have $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \phi(\mathbf{x}_i)$ and thus

$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i) \right)$$

- Kernels between vectors
 - Linear kernel: $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$
 - Polynomial kernel: $K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d$
 - Gaussian RBF kernel: $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|_2^2)$
 - Laplace RBF kernel: $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|_1)$
- Kernels on structured data
 - Many string, tree and graph kernels [Gärtner, 2003]
 - Frameworks to design structured kernels: convolution kernels [Haussler, 1999], mapping kernels [Shin and Kuboyama, 2008]
- Use closure properties to build new kernels from existing ones
 - Sums, positive combinations, etc

- Kernels allow to obtain **nonlinear variants** for many linear machine learning algorithms
- A few examples
 - SVM
 - Ridge regression
 - PCA
 - CCA
 - K-Means
- Today's topic: **scalability of kernel methods**
 - Identify general issues
 - Study general solutions

SCALABILITY ISSUES

- All kernel methods rely on the **Gram matrix** $\mathbf{G} \in \mathbb{R}^{n \times n}$

$$\mathbf{G} = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \cdots & K(\mathbf{x}_1, \mathbf{x}_n) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \cdots & K(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}_n, \mathbf{x}_1) & K(\mathbf{x}_n, \mathbf{x}_2) & \cdots & K(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

- Assuming complexity of one kernel evaluation is constant
 - Constructing \mathbf{G} takes $O(n^2)$ time
 - If we need to invert \mathbf{G} (e.g., Kernel Ridge Regression): $O(n^3)$ time
- **Training time of popular algorithms:** $O(n^2)$ or $O(n^3)$
- This is infeasible for large n

- Kernel methods are **nonparametric**: the learned model relies on the training data points
- To process a test point \mathbf{x} , one generally needs to evaluate the kernel between \mathbf{x} and the training points
 - For instance in SVM:

$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i) \right)$$

- Number of support vectors grows linearly with n [Steinwart, 2003]
- **Prediction time of kernel methods: $O(n)$**
- This is slow when n is large

- Today we will study two techniques which can be used to scale up any kernel method
- **Random Kernel Features**: approximate kernel function / map
- **Nyström Approximation**: approximate Gram matrix
- Both methods are very popular and successful in practice
- Still an active area of research

RANDOM KERNEL FEATURES

Kernel SVM : primal formulation

$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 \quad + \quad C \sum_{i=1}^n [y_i(\mathbf{w}^T \phi(\mathbf{x}) + b)]_+$$

where $[a]_+ = \max(0, 1 - a)$ is the hinge loss function

- When $\phi(\mathbf{x})$ is known and finite-dimensional (e.g., linear kernel)
 - Training time linear in n (see e.g., [Shalev-Shwartz et al., 2011])
 - Prediction complexity independent of n
- But $\phi(\mathbf{x})$ is usually unknown and potentially infinite-dimensional
- In this case we can only solve the problem in dual form
 - Training time quadratic or cubic in n
 - Prediction complexity linear in n

- **Idea:** find a **finite-dimensional feature map** $\hat{\phi}(\mathbf{x}) \in \mathbb{R}^c$ such that

$$\langle \hat{\phi}(\mathbf{x}), \hat{\phi}(\mathbf{x}') \rangle \approx K(\mathbf{x}, \mathbf{x}')$$

- We can then solve in primal form to get $\mathbf{w} \in \mathbb{R}^c$ and $b \in \mathbb{R}$
- We can predict using

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \hat{\phi}(\mathbf{x}) + b)$$

- If $c \ll n^2$, training is much faster
- If $c \ll n$, prediction is also much faster

Definition (Shift-invariant kernel)

Let $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ be a positive definite kernel. K is said to be *shift-invariant* if for any $\mathbf{a} \in \mathbb{R}^p$ and any $(\mathbf{x}, \mathbf{x}') \in \mathbb{R}^p \times \mathbb{R}^p$

$$K(\mathbf{x} - \mathbf{a}, \mathbf{x}' - \mathbf{a}) = K(\mathbf{x}, \mathbf{x}').$$

For simplicity we denote $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x} - \mathbf{x}')$.

- Examples of shift-invariant kernels
 - Gaussian RBF kernel: $K(\mathbf{x} - \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|_2^2)$
 - Laplace kernel: $K(\mathbf{x} - \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|_1)$

Theorem (Bochner's theorem, see [\[Rahimi and Recht, 2007\]](#))

A continuous shift-invariant kernel $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x} - \mathbf{x}')$ is positive definite if and only if $K(\boldsymbol{\delta})$ is the Fourier transform of a nonnegative measure. In particular, if K is properly scaled, we have

$$K(\mathbf{x} - \mathbf{x}') = \int_{\mathbb{R}^p} P(\boldsymbol{\omega}) e^{i\boldsymbol{\omega}^T \mathbf{x}} e^{-i\boldsymbol{\omega}^T \mathbf{x}'} d\boldsymbol{\omega},$$

where $P(\boldsymbol{\omega})$ is a real-valued probability density function over \mathbb{R}^p .

BOCHNER'S THEOREM

Developing the result of Bochner's theorem:

$$\begin{aligned} K(\mathbf{x} - \mathbf{x}') &= \int_{\mathbb{R}^p} P(\boldsymbol{\omega}) e^{i\boldsymbol{\omega}^T \mathbf{x}} e^{-i\boldsymbol{\omega}^T \mathbf{x}'} d\boldsymbol{\omega} \\ &= \int_{\mathbb{R}^p} P(\boldsymbol{\omega}) \cos(\boldsymbol{\omega}^T \mathbf{x} - \boldsymbol{\omega}^T \mathbf{x}') d\boldsymbol{\omega} \end{aligned} \quad (1)$$

$$= \int_{\mathbb{R}^p} P(\boldsymbol{\omega}) \left(\cos(\boldsymbol{\omega}^T \mathbf{x}) \cos(\boldsymbol{\omega}^T \mathbf{x}') + \sin(\boldsymbol{\omega}^T \mathbf{x}) \sin(\boldsymbol{\omega}^T \mathbf{x}') \right) d\boldsymbol{\omega} \quad (2)$$

$$= \int_{\mathbb{R}^p} \int_{b=0}^{2\pi} \frac{P(\boldsymbol{\omega})}{2\pi} 2 \cos(\boldsymbol{\omega}^T \mathbf{x} + b) \cos(\boldsymbol{\omega}^T \mathbf{x}' + b) d\boldsymbol{\omega} db \quad (3)$$

$$= \mathbb{E}_{\boldsymbol{\omega} \sim P, b \sim \mathcal{U}(0, 2\pi)} \left[\sqrt{2} \cos(\boldsymbol{\omega}^T \mathbf{x} + b) \sqrt{2} \cos(\boldsymbol{\omega}^T \mathbf{x}' + b) \right] \quad (4)$$

(1): $K(\mathbf{x} - \mathbf{x}')$, $P(\boldsymbol{\omega}) \in \mathbb{R}$ so we can ignore imaginary part

(2) and (3): use sum of angles formulas

- We have obtained

$$K(\mathbf{x} - \mathbf{x}') = \mathbb{E}_{\boldsymbol{\omega} \sim P, b \sim \mathcal{U}(0, 2\pi)} \left[\sqrt{2} \cos(\boldsymbol{\omega}^T \mathbf{x} + b) \sqrt{2} \cos(\boldsymbol{\omega}^T \mathbf{x}' + b) \right]$$

- K can thus be written as an expectation over $\boldsymbol{\omega}$ drawn from the distribution P
- If we know P for the kernel K of interest, we can approximate K by random sampling

1. Set number of random kernel features c
2. Draw $\omega_1, \dots, \omega_c \sim P(\omega)$ and $b_1, \dots, b_c \sim \mathcal{U}(0, 2\pi)$
3. Map training points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ to their random kernel features $\hat{\phi}(\mathbf{x}_1), \dots, \hat{\phi}(\mathbf{x}_n) \in \mathbb{R}^c$ where

$$\hat{\phi}_j(\mathbf{x}) = \sqrt{\frac{2}{c}} \cos(\omega_j^T \mathbf{x} + b_j), \quad j \in \{1, \dots, c\}$$

4. Train linear model (such as linear SVM) on transformed data $\hat{\phi}(\mathbf{x}_1), \dots, \hat{\phi}(\mathbf{x}_n) \in \mathbb{R}^c$

- For the Gaussian RBF kernel, P is a Gaussian distribution

$$P^{rbf}(\boldsymbol{\omega}) = \mathcal{N}(0, 2\gamma \mathbf{I}_p)$$

- For the Laplace kernel, P is a Cauchy distribution

$$P^{lap}(\boldsymbol{\omega}) = \prod_p \frac{1}{\pi(1 + \omega_p^2/2\gamma)}$$

Theorem ([Rahimi and Recht, 2007])

Let $c \geq 1$ and $\hat{\phi} : \mathbb{R}^p \rightarrow \mathbb{R}^c$ be the feature map obtained by drawing $\omega_1, \dots, \omega_c$ from $P(\omega)$. Then we have with high probability:

$$\sup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}} \left| \langle \hat{\phi}(\mathbf{x}), \hat{\phi}(\mathbf{x}') \rangle - K(\mathbf{x}, \mathbf{x}') \right| \leq O\left(\sqrt{\frac{p}{c}}\right)$$

- Kernel approximation error uniformly decreases in $O(\sqrt{1/c})$
- Can also bound generalization error [Rahimi and Recht, 2008]

- Example of large-scale application: acoustic models for speech recognition [Lu et al., 2016]
 - Trained on 50 hours of speech with $c = 500K$ random features
 - Performance comparable to deep neural nets
- Random features exist for other kernels, such as dot product kernels (including polynomial kernels) [Kar and Karnick, 2012]
- Techniques to speed up prediction further [Le et al., 2013]
 - From $O(cp)$ to $O(c \log p)$ time
- Easy to combine multiple kernels by stacking their random features [Lu et al., 2014]

NYSTRÖM APPROXIMATION

LOW-RANK APPROXIMATION

- Let us consider the Gram matrix $\mathbf{G} \in \mathbb{R}^{n \times n}$ ($G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$)
- When n is large, could approximate \mathbf{G} with a matrix of rank $k \leq n$
- Consider the spectral decomposition of \mathbf{G}

$$\mathbf{G} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

- $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n]^T \in \mathbb{R}^{n \times n}$ the set of eigenvectors
- $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ the eigenvalues ($\lambda_1 \geq \dots \geq \lambda_n$)
- **Best rank- k approximation** $\mathbf{G}_k \in \mathbb{R}^{n \times n}$ is given by

$$\mathbf{G}_k = \mathbf{U}_k \mathbf{\Lambda}_k \mathbf{U}_k^T$$

where $\mathbf{U}_k = [\mathbf{u}_1, \dots, \mathbf{u}_k]^T \in \mathbb{R}^{n \times k}$ and $\mathbf{\Lambda}_k = \text{diag}(\lambda_1, \dots, \lambda_k)$

NYSTRÖM APPROXIMATION

- In the context of kernel methods this is useless
 - We still need to construct \mathbf{G} : $O(n^2)$ time
 - And compute its k -thresholded spectral decomposition: between $O(n^2)$ and $O(n^3)$ depending on the value of k
- Alternative: Nyström approximation [Drineas and Mahoney, 2005]
 1. Sample a set \mathcal{I} of c indices uniformly in $\{1, \dots, n\}$
 2. Compute $\mathbf{C} \in \mathbb{R}^{n \times c}$ with $\mathbf{C}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$ for $i \in \{1, \dots, n\}$ and $j \in \mathcal{I}$
 3. Form reduced Gram matrix $\tilde{\mathbf{G}} \in \mathbb{R}^{c \times c}$ with $\tilde{\mathbf{G}}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$ for $i, j \in \mathcal{I}$
 4. Compute $\tilde{\mathbf{G}}_k \in \mathbb{R}^{c \times c}$, the best rank- k approximation of $\tilde{\mathbf{G}}$ ($k \leq c$)
 5. Final Gram matrix approximation $\hat{\mathbf{G}} \in \mathbb{R}^{n \times n}$

$$\mathbf{G}_{nys} \in \mathbb{R}^{n \times n} = \mathbf{C} \tilde{\mathbf{G}}_k^+ \mathbf{C}^T$$

where $\tilde{\mathbf{G}}_k^+$ is the pseudoinverse of $\tilde{\mathbf{G}}_k$

- Time complexity is $O(c^3 + nck)$

Theorem ([Drineas and Mahoney, 2005])

Let $\mathbf{G} \in \mathbb{R}^{n \times n}$ be the Gram matrix. Let \mathbf{G}_k be its best rank- k approximation and \mathbf{G}_{nys} be its Nyström approximation. We have with high probability:

$$\|\mathbf{G} - \mathbf{G}_{nys}\|_F \leq \|\mathbf{G} - \mathbf{G}_k\|_F + O\left(\sqrt{\frac{n}{c}}\right)$$

- If c is large enough, \mathbf{G}_{nys} is nearly as good as \mathbf{G}_k
- Some nonuniform sampling techniques have been proposed, but uniform sampling tends to work best in practice
[Kumar et al., 2009]

- \mathbf{G}_{nys} can be used directly to speed up training algorithms
- It can also be used to generate an explicit feature map as in RKF
- Denote $\tilde{\mathbf{G}}_k = \tilde{\mathbf{U}}_k \tilde{\mathbf{\Lambda}}_k \tilde{\mathbf{U}}_k^T$
- Approximate feature map $\hat{\phi}(\mathbf{x}) \in \mathbb{R}^k$ is given by

$$\hat{\phi}(\mathbf{x}) = \mathbf{k}_{\mathbf{x}, \mathcal{I}} \tilde{\mathbf{U}}_k \tilde{\mathbf{\Lambda}}_k^{-1/2}$$

where $\mathbf{k}_{\mathbf{x}, \mathcal{I}} = [K(\mathbf{x}, \mathbf{x}_i)]_{i \in \mathcal{I}}$

CONCLUSION

NYSTRÖM APPROXIMATION VS. RANDOM KERNEL FEATURES

- Both methods have a $O(1/\sqrt{c})$ convergence rate (c : # of random features for RKF, # of random columns for Nyström)
- Nyström' approximation guarantee is adaptive to the data, while RKF is data-independent
 - In fact, Nyström can achieve $O(1/c)$ convergence when eigengap of \mathbf{G} is large [Yang et al., 2012]
- At equal number of random samples (features/columns)
 - Nyström tends to achieve better performance
 - But RKF are generally cheaper to generate (no spectral decomposition or matrix inversion needed)
- Lacks a solid large-scale comparison

- Kernel methods: general class of nonlinear algorithms
- Training and prediction time scales badly with n
- Two general techniques to make kernel methods scalable:
Random Kernel Features and Nyström approximation
- We can then take advantage of existing fast solvers for linear algorithms

REFERENCES I

- [Drineas and Mahoney, 2005] Drineas, P. and Mahoney, M. W. (2005).
On the nystrom method for approximating a gram matrix for improved kernel-based learning.
Journal of Machine Learning Research, 6:2153–2175.
- [Gärtner, 2003] Gärtner, T. (2003).
A survey of kernels for structured data.
SIGKDD Explorations, 5(1):49–58.
- [Haussler, 1999] Haussler, D. (1999).
Convolution Kernels on Discrete Structure.
Technical Report UCSC-CRL-99-10, University of California at Santa Cruz.
- [Kar and Karnick, 2012] Kar, P. and Karnick, H. (2012).
Random feature maps for dot product kernels.
In AISTATS, pages 583–591.
- [Kumar et al., 2009] Kumar, S., Mohri, M., and Talwalkar, A. (2009).
Sampling techniques for the nystrom method.
In AISTATS, pages 304–311.
- [Le et al., 2013] Le, Q., Sarlós, T., and Smola, A. (2013).
Fastfood — Approximating Kernel Expansions in Loglinear Time.
In ICML.

REFERENCES II

- [Lu et al., 2016] Lu, Z., Guo, D., Bagheri Garakani, A., Liu, K., May, A., Bellet, A., Fan, L., Collins, M., Kingsbury, B., Picheny, M., and Sha, F. (2016).
A Comparison Between Deep Neural Nets and Kernel Acoustic Models for Speech Recognition.
In ICASSP.
- [Lu et al., 2014] Lu, Z., May, A., Liu, K., Bagheri Garakani, A., Guo, D., Bellet, A., Fan, L., Collins, M., Kingsbury, B., Picheny, M., and Sha, F. (2014).
How to Scale Up Kernel Methods to Be As Good As Deep Neural Nets.
Technical report, arXiv:1411.4000.
- [Rahimi and Recht, 2007] Rahimi, A. and Recht, B. (2007).
Random Features for Large-Scale Kernel Machines.
In NIPS.
- [Rahimi and Recht, 2008] Rahimi, A. and Recht, B. (2008).
Weighted Sums of Random Kitchen Sinks: Replacing minimization with randomization in learning.
In NIPS, pages 1313–1320.
- [Shalev-Shwartz et al., 2011] Shalev-Shwartz, S., Singer, Y., Srebro, N., and Cotter, A. (2011).
Pegasos: Primal Estimated sub-GrAdient SOLver for SVM.
Mathematical Programming, 127(1):3–30.

REFERENCES III

- [Shin and Kuboyama, 2008] Shin, K. and Kuboyama, T. (2008).
A generalization of Haussler's convolution kernel: mapping kernel.
In ICML, pages 944–951.
- [Steinwart, 2003] Steinwart, I. (2003).
Sparseness of Support Vector Machines.
Journal of Machine Learning Research, 4:1071–1105.
- [Yang et al., 2012] Yang, T., Li, Y.-F., Mahdavi, M., Jin, R., and Zhou, Z.-H. (2012).
Nyström method vs random fourier features: A theoretical and empirical comparison.
In NIPS, pages 485–493.