

Large-Scale Similarity and Distance Metric Learning

Aurélien Bellet
Télécom ParisTech

Joint work with K. Liu, Y. Shi and F. Sha (USC), S. Cléménçon and
I. Colin (Télécom ParisTech)

Séminaire Criteo
March 31, 2015

A bit about me

- ▶ PhD in Computer Science (Dec 2012)
 - ▶ Université Jean Monnet, Saint-Etienne
 - ▶ Advisors: Marc Sebban, Amaury Habrard

- ▶ Postdoc in 2013–2014 (18 months)
 - ▶ University of Southern California, Los Angeles
 - ▶ Working with Fei Sha

- ▶ Joined Télécom ParisTech in October
 - ▶ Chaire “Machine Learning for Big Data”
 - ▶ Working with Stéphan Cléménçon

Outline of the talk

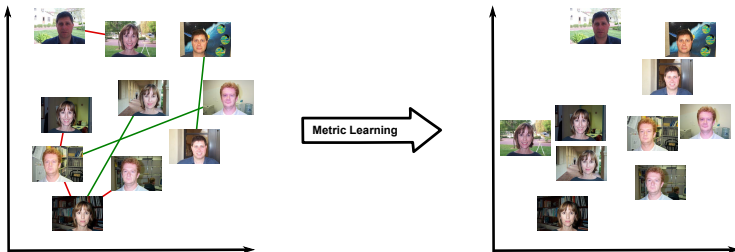
1. Introduction to Metric Learning
2. Learning (Infinitely) Many Local Metrics
3. Similarity Learning for High-Dimensional Sparse Data
4. Scaling-up ERM for Metric Learning

Introduction

Metric learning

Motivation

- ▶ How to assign a **similarity score** to a pair of objects?
 - ▶ Basic component of many algorithms: k -NN, clustering, kernel methods, ranking, dimensionality reduction, data visualization
 - ▶ Crucial to performance of the above
 - ▶ Obviously **problem-dependent**
- ▶ The machine learning way: let's learn it from data!
 - ▶ We'll need some examples of similar / dissimilar things
 - ▶ Learn to approximate the **latent notion of similarity**
 - ▶ Can be thought of as representation learning



Metric learning

Basic recipe

1. Pick a parametric form of distance or similarity function

- ▶ (Squared) Mahalanobis distance

$$d_{\mathbf{M}}^2(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}') \text{ with } \mathbf{M} \in \mathbb{R}^{d \times d} \text{ symmetric PSD}$$

- ▶ Bilinear similarity

$$S_{\mathbf{M}}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{M} \mathbf{x}' \text{ with } \mathbf{M} \in \mathbb{R}^{d \times d}$$

2. Collect similarity judgments on data pairs/triplets

- ▶ \mathbf{x}_i and \mathbf{x}_j are similar (or dissimilar)
- ▶ \mathbf{x}_i is more similar to \mathbf{x}_j than to \mathbf{x}_k

3. Estimate parameters such that metric best satisfies them

- ▶ Convex optimization and the like

Metric learning

A famous algorithm: LMNN [Weinberger et al., 2005]

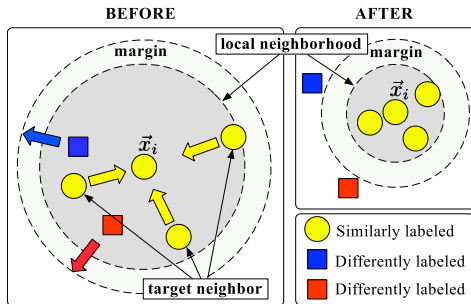
Large Margin Nearest Neighbors

$$\min_{\mathbf{M} \in \mathbb{S}_+^d, \xi \geq 0} (1 - \mu) \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_j) + \mu \sum_{i,j,k} \xi_{ijk}$$

$$\text{s.t. } d_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_k) - d_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \xi_{ijk} \quad \forall (\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in \mathcal{R},$$

$\mathcal{S} = \{(\mathbf{x}_i, \mathbf{x}_j) : y_i = y_j \text{ and } \mathbf{x}_j \text{ belongs to the } k\text{-neighborhood of } \mathbf{x}_i\}$

$\mathcal{R} = \{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) : (\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}, y_i \neq y_k\}$



Metric learning

Large-scale challenges

- ▶ How to efficiently learn **multiple metrics**?
 - ▶ Multi-task learning
 - ▶ Local linear metrics

- ▶ How to efficiently learn a metric for **high-dimensional data**?
 - ▶ Full $d \times d$ matrix not feasible anymore

- ▶ How to deal with **large datasets**?
 - ▶ Number of terms in the objective can grow as $O(n^2)$ or $O(n^3)$

Learning (Infinitely) Many Local Metrics

[Shi et al., 2014]

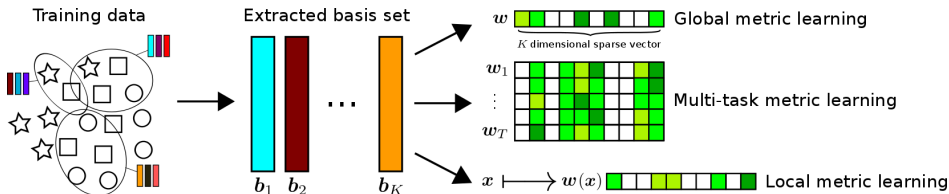
Learning (infinitely) many local metrics

Main idea

- ▶ Assume we are given a **basis dictionary** $B = \{\mathbf{b}_k\}_{k=1}^K$ in \mathbb{R}^d
- ▶ We will learn $d_w^2(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^T \mathbf{M}(\mathbf{x} - \mathbf{x}')$ where

$$\mathbf{M} = \sum_{k=1}^K w_k \mathbf{b}_k \mathbf{b}_k^T, \quad \mathbf{w} \geq 0$$

- ▶ Use sparsity-inducing regularizer to do **basis selection**



Learning (infinitely) many local metrics

Global and multi-task convex formulations

$$\mathcal{C} = \{ \mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^- : d_{\mathbf{w}}^2(\mathbf{x}_i, \mathbf{x}_i^+) \text{ should be smaller than } d_{\mathbf{w}}^2(\mathbf{x}_i, \mathbf{x}_i^-) \}_{i=1}^C$$

SCML-Global

$$\min_{\mathbf{w} \in \mathbb{R}_+^K} \frac{1}{C} \sum_{i=1}^C [1 + d_{\mathbf{w}}^2(\mathbf{x}_i, \mathbf{x}_i^+) - d_{\mathbf{w}}^2(\mathbf{x}_i, \mathbf{x}_i^-)]_+ + \beta \|\mathbf{w}\|_1,$$

where $[\cdot]_+ = \max(0, \cdot)$

mt-SCML : extension to T tasks

$$\min_{\mathbf{W} \in \mathbb{R}_+^{T \times K}} \sum_{t=1}^T \frac{1}{C_t} \sum_{i=1}^{C_t} [1 + d_{\mathbf{w}_t}^2(\mathbf{x}_i, \mathbf{x}_i^+) - d_{\mathbf{w}_t}^2(\mathbf{x}_i, \mathbf{x}_i^-)]_+ + \beta \|\mathbf{W}\|_{2,1}$$

- ▶ $\|\mathbf{W}\|_{2,1}$: induce column-wise sparsity
- ▶ Metrics are task-specific but regularized to share features

Learning (infinitely) many local metrics

Local metric learning

- ▶ Learn multiple local metrics : more flexibility
 - ▶ One metric per region or class
 - ▶ One metric per training instance
- ▶ Many issues
 - ▶ How to define the local regions?
 - ▶ How to generalize at test time?
 - ▶ How to avoid a growing number of parameters?
 - ▶ How to do avoid overfitting?
- ▶ Proposed approach : learn a metric for each point in space

Learning (infinitely) many local metrics

Local nonconvex formulation

- ▶ We use the following parameterization :

$$d_w(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^T \left(\sum_{k=1}^K w_k(\mathbf{x}) \mathbf{b}_k \mathbf{b}_k^T \right) (\mathbf{x} - \mathbf{x}'),$$

- ▶ $w_k(\mathbf{x}) = (\mathbf{a}_k^T \phi(\mathbf{x}) + c_k)^2$: weight of basis k
- ▶ $\phi(\mathbf{x}) \in \mathbb{R}^{d'}$: (nonlinear) projection of \mathbf{x} (e.g., kernel PCA)
- ▶ $\mathbf{a}_k \in \mathbb{R}^{d'}$ and $c_k \in \mathbb{R}$: parameters to learn for each basis

SCML-Local

$$\min_{\tilde{\mathbf{A}} \in \mathbb{R}_+^{(d'+1) \times K}} \frac{1}{C} \sum_{i=1}^C [1 + d_w^2(\mathbf{x}_i, \mathbf{x}_i^+) - d_w^2(\mathbf{x}_i, \mathbf{x}_i^-)]_+ + \beta \|\tilde{\mathbf{A}}\|_{2,1}$$

- ▶ $\tilde{\mathbf{A}}$: stack $\mathbf{A} = [\mathbf{a}_1 \dots \mathbf{a}_K]^T$ and \mathbf{c}
- ▶ If $\mathbf{A} = \mathbf{0}$, recover the global formulation

Learning (infinitely) many local metrics

Efficient optimization

- ▶ Nonsmooth objective with large number of terms
 - ▶ Use **stochastic proximal** methods

- ▶ SCML-Global and mt-SCML are **convex**
 - ▶ Regularized Dual Averaging [Xiao, 2010]

- ▶ SCML-Local is **nonconvex**
 - ▶ Forward-backward algorithm [Duchi and Singer, 2009]
 - ▶ Initialize with solution of SCML-Global

Learning (infinitely) many local metrics

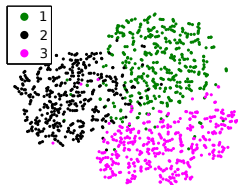
Experiments

Dataset	Euc	LMNN	BoostML	SCML-Global
Vehicle	29.7±0.6	23.5±0.7	19.9±0.6	21.3±0.6
Vowel	11.1±0.4	10.8±0.4	11.4±0.4	10.9±0.5
Segment	5.2±0.2	4.6±0.2	3.8±0.2	4.1±0.2
Letters	14.0±0.2	11.6±0.3	10.8±0.2	9.0±0.2
USPS	10.3±0.2	4.1±0.1	7.1±0.2	4.1±0.1
BBC	8.8±0.3	4.0±0.2	9.3±0.3	3.9±0.2
Avg. rank	3.3	2.0	2.3	1.2

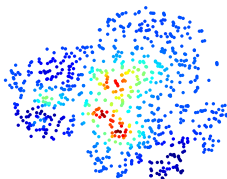
Dataset	M ² LMNN	GLML	PLML	SCML-Local
Vehicle	23.1±0.6	23.4±0.6	22.8±0.7	18.0±0.6
Vowel	6.8±0.3	4.1±0.4	8.3±0.4	6.1±0.4
Segment	3.6±0.2	3.9±0.2	3.9±0.2	3.6±0.2
Letters	9.4±0.3	10.3±0.3	8.3±0.2	8.3±0.2
USPS	4.2±0.7	7.8±0.2	4.1±0.1	3.6±0.1
BBC	4.9±0.4	5.7±0.3	4.3±0.2	4.1±0.2
Avg. rank	2.0	2.7	2.0	1.2

Learning (infinitely) many local metrics

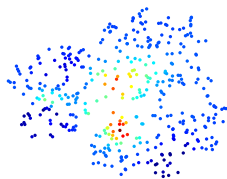
Experiments



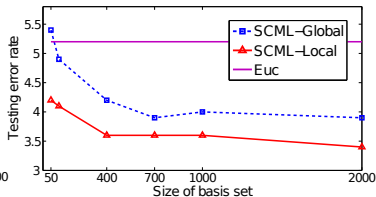
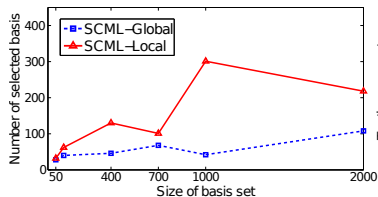
(a) Class membership



(b) Trained metrics



(c) Test metrics



Similarity Learning for High-Dimensional Sparse Data

[Liu et al., 2015]

Similarity learning for high-dimensional sparse data

Problem setting

- ▶ Assume data points are **high-dimensional** ($d > 10^4$) but **D -sparse** (on average) with $D \ll d$
 - ▶ Bags-of-words (text, image), bioinformatics, etc
- ▶ Existing metric learning algorithms **fail**
 - ▶ **Intractable**: training cost $O(d^2)$ to $O(d^3)$, memory $O(d^2)$
 - ▶ Severe **overfitting**
- ▶ Practitioners use **dimensionality reduction** (PCA, RP)
 - ▶ Poor performance in presence of noisy features
 - ▶ Resulting metric difficult to interpret for domain experts
- ▶ Contributions of this work
 - ▶ Learn similarity in original high-dimensional space
 - ▶ Time/memory costs **independent of d**
 - ▶ Explicit **control of similarity complexity**

Similarity learning for high-dimensional sparse data

A very large basis set

- ▶ We want to learn a similarity function $S_{\mathbf{M}}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{M} \mathbf{x}'$
- ▶ Given $\lambda > 0$, for any $i, j \in \{1, \dots, d\}$, $i \neq j$ we define

$$\mathbf{P}_{\lambda}^{(ij)} = \begin{pmatrix} \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \lambda & \cdot & \lambda & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \lambda & \cdot & \lambda & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix} \quad \mathbf{N}_{\lambda}^{(ij)} = \begin{pmatrix} \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \lambda & \cdot & -\lambda & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & -\lambda & \cdot & \lambda & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}$$

$$\mathcal{B}_{\lambda} = \bigcup_{ij} \left\{ \mathbf{P}_{\lambda}^{(ij)}, \mathbf{N}_{\lambda}^{(ij)} \right\}$$

$$\mathbf{M} \in \mathcal{D}_{\lambda} = \text{conv}(\mathcal{B}_{\lambda})$$

- ▶ One basis involves only 2 features:

$$S_{\mathbf{P}_{\lambda}^{(ij)}}(\mathbf{x}, \mathbf{x}') = \lambda(x_i x'_i + x_j x'_j + x_i x'_j + x_j x'_i)$$

$$S_{\mathbf{N}_{\lambda}^{(ij)}}(\mathbf{x}, \mathbf{x}') = \lambda(x_i x'_i + x_j x'_j - x_i x'_j - x_j x'_i)$$

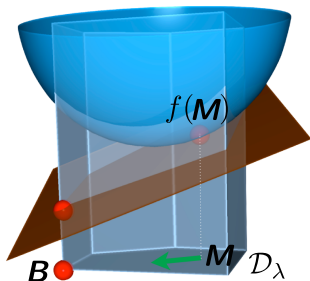
Similarity learning for high-dimensional sparse data

Problem formulation and algorithm

Optimization problem (smoothed hinge loss ℓ)

$$\begin{aligned} \min_{\mathbf{M} \in \mathbb{R}^{d \times d}} \quad & f(\mathbf{M}) = \frac{1}{C} \sum_{i=1}^C \ell \left(1 - \mathbf{x}_i^T \mathbf{M} \mathbf{x}_i^+ + \mathbf{x}_i^T \mathbf{M} \mathbf{x}_i^- \right) \\ \text{s.t.} \quad & \mathbf{M} \in \mathcal{D}_\lambda \end{aligned}$$

- ▶ Use a **Frank-Wolfe algorithm** [Jaggi, 2013] to solve it



Let $\mathbf{M}^{(0)} \in \mathcal{D}_\lambda$

for $k = 0, 1, \dots$ **do**

$$\mathbf{B}^{(k)} = \arg \min_{\mathbf{B} \in \mathcal{B}_\lambda} \left\langle \mathbf{B}, \nabla f(\mathbf{M}^{(k)}) \right\rangle$$

$$\mathbf{M}^{(k+1)} = (1 - \gamma) \mathbf{M}^{(k)} + \gamma \mathbf{B}^{(k)}$$

end for

Similarity learning for high-dimensional sparse data

Convergence and complexity analysis

Convergence

Let $L = \frac{1}{C} \sum_{i=1}^C \|\mathbf{x}_i(\mathbf{x}_i^+ - \mathbf{x}_i^-)^T\|_F^2$. At any iteration $k \geq 1$, the iterate $\mathbf{M}^{(k)} \in \mathcal{D}_\lambda$ of the FW algorithm:

- ▶ has at most rank $k + 1$ with $4(k + 1)$ nonzero entries
 - ▶ uses at most $2(k + 1)$ distinct features
 - ▶ satisfies $f(\mathbf{M}^{(k)}) - f(\mathbf{M}^*) \leq 16L\lambda^2/(k + 2)$
-
- ▶ An optimal basis can be found in $O(CD^2)$ time and memory
 - ▶ An approximately optimal basis can be found in $O(mD^2)$ with $m \ll C$ using a Monte Carlo approximation of the gradient
 - ▶ Or even $O(mD)$ using a heuristic (good results in practice)
 - ▶ Storing $\mathbf{M}^{(k)}$ requires only $O(k)$ memory
 - ▶ Or even the entire sequence $\mathbf{M}^{(0)}, \dots, \mathbf{M}^{(k)}$ at the same cost

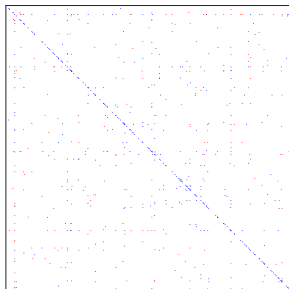
Similarity learning for high-dimensional sparse data

Experiments

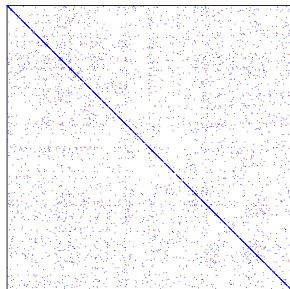
- ▶ K -NN test error on datasets with d up to 10^5

Datasets	IDENTITY	RP+OASIS	PCA+OASIS	DIAG- ℓ_2	DIAG- ℓ_1	HDSL
dexter	20.1	24.0	9.3	8.4	8.4	6.5
dorothea	9.3	11.4	9.9	6.8	6.6	6.5
rcv1.2	6.9	7.0	4.5	3.5	3.7	3.4
rcv1.4	11.2	10.6	6.1	6.2	7.2	5.7

- ▶ Sparsity structure of the matrices



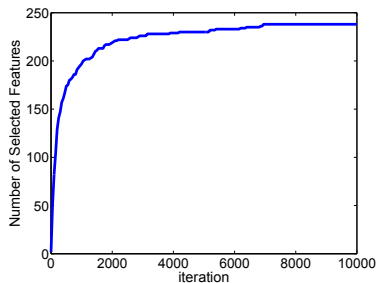
(a) dexter ($20,000 \times 20,000$ matrix, 712 nonzeros)



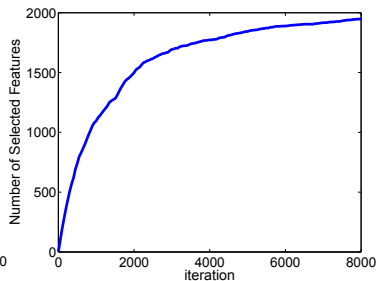
(b) rcv1.4 ($29,992 \times 29,992$ matrix, 5263 nonzeros)

Similarity learning for high-dimensional sparse data

Experiments



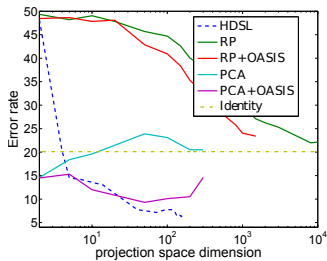
(a) dexter dataset



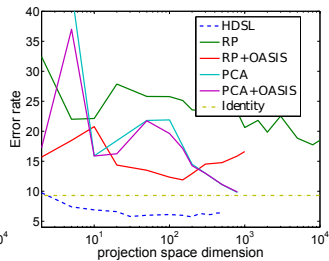
(b) rcv1_4 dataset

Similarity learning for high-dimensional sparse data

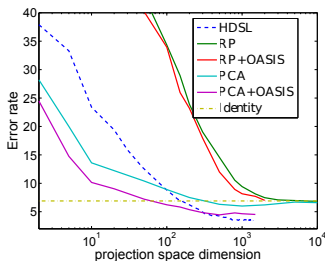
Experiments



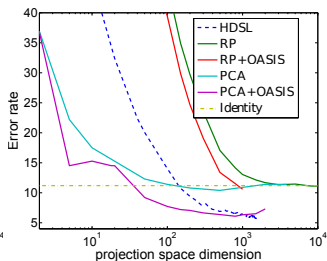
(a) dexter dataset



(b) dorothea dataset



(c) rcv1.2 dataset



(d) rcv1.4 dataset

Scaling-up ERM for Metric Learning

[Cléménçon et al., 2015]

Scaling-up ERM for metric learning

Empirical Minimization of U -Statistics

- ▶ Given an i.i.d. sample X_1, \dots, X_n , the U -statistic of degree 2 with kernel H is given by

$$U_n(H; \theta) = \frac{2}{n(n-1)} \sum_{i < j} H(X_i, X_j; \theta)$$

- ▶ Can be generalized to higher degrees and multi-samples
- ▶ Empirical Minimization of U -statistic: $\min_{\theta \in \Theta} U_n(H; \theta)$
 - ▶ Applies to metric learning for classification
 - ▶ Also pairwise clustering, ranking, etc
- ▶ Number of terms quickly becomes huge
 - ▶ MNIST dataset : $n = 60000 \rightarrow 2 \times 10^9$ pairs

Scaling-up ERM for metric learning

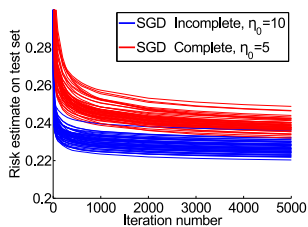
Main idea and results

- ▶ Approximate $U_n(H)$ by an incomplete version $\tilde{U}_B(H)$
 - ▶ Uniformly sample B terms (with replacement)
- ▶ Main result : if $B = O(n)$, the learning rate is preserved
 - ▶ $O(1/\sqrt{n})$ convergence in the general case
 - ▶ Objective function has only $O(n)$ terms v.s. $O(n^2)$ initially
 - ▶ Due to high dependence in the pairs
- ▶ Naive strategy : use complete U -statistic with fewer samples
 - ▶ Uniformly sample $O(\sqrt{n})$ samples
 - ▶ Form all possible pairs : $O(n)$ terms
 - ▶ Learning rate is only $O(1/n^{1/4})$
- ▶ Can use similar ideas in mini-batch SGD
 - ▶ Reduce variance of gradient estimates

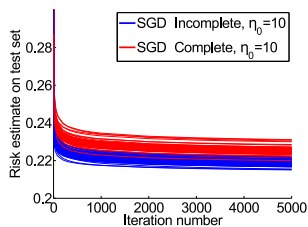
Scaling-up ERM for metric learning

Experiments

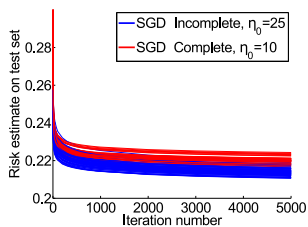
- ▶ MNIST dataset : $n = 60000 \rightarrow 2 \times 10^9$ pairs
- ▶ Mini-batch SGD
 - ▶ Step size tuning
 - ▶ 50 random runs



(a) $B = 10$



(b) $B = 55$



(c) $B = 253$

Conclusion and perspectives

- ▶ Metric learning: useful in a variety of setting
- ▶ Ideas to overcome some key limitations
 - ▶ Large datasets
 - ▶ High dimensionality
 - ▶ Learn many local metrics
- ▶ Details and more results in the papers
- ▶ One interesting challenge : [distributed metric learning](#)
 - ▶ Existing work [Xie and Xing, 2014] too restrictive

References I

- [Clémençon et al., 2015] Clémençon, S., Bellet, A., and Colin, I. (2015).
Scaling-up Empirical Risk Minimization: Optimization of Incomplete U-statistics.
Technical report, arXiv:1501.02629.
- [Duchi and Singer, 2009] Duchi, J. and Singer, Y. (2009).
Efficient Online and Batch Learning Using Forward Backward Splitting.
JMLR, 10:2899–2934.
- [Jaggi, 2013] Jaggi, M. (2013).
Revisiting Frank-Wolfe: Projection-Free Sparse Convex Optimization.
In *ICML*.
- [Liu et al., 2015] Liu, K., Bellet, A., and Sha, F. (2015).
Similarity Learning for High-Dimensional Sparse Data.
In *AISTATS*.
- [Shi et al., 2014] Shi, Y., Bellet, A., and Sha, F. (2014).
Sparse Compositional Metric Learning.
In *AAAI*, pages 2078–2084.
- [Weinberger et al., 2005] Weinberger, K. Q., Blitzer, J., and Saul, L. K. (2005).
Distance Metric Learning for Large Margin Nearest Neighbor Classification.
In *NIPS*, pages 1473–1480.

References II

[Xiao, 2010] Xiao, L. (2010).

Dual Averaging Methods for Regularized Stochastic Learning and Online Optimization.

JMLR, 11:2543–2596.

[Xie and Xing, 2014] Xie, P. and Xing, E. (2014).

Large Scale Distributed Distance Metric Learning.

Technical report, arXiv:1412.5949.