

Fast low-rank metric learning

Dan-Theodor Oneață



Master of Science
Artificial Intelligence
School of Informatics
University of Edinburgh
2011

Abstract

This doctoral thesis will present the results of my work into the reanimation of lifeless human tissues.

Acknowledgements

Many thanks to my mummy for the numerous packed lunches; and of course to Igor, my faithful lab assistant. DP IM.

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

(Dan-Theodor Oneață)

Table of Contents

1	Introduction	1
2	Background	4
2.1	Theoretical background	4
2.2	Related methods	6
3	Neighbourhood component analysis	8
3.1	General presentation	8
3.2	Practical notes	8
3.3	Class-conditional kernel density estimation interpretation	9
4	Reducing the computational cost	11
4.1	Sub-sampling	11
4.2	Mini-batches	12
4.3	Stochastic learning	15
4.4	Approximate computations	15
4.5	Exact computations	15
4.6	NCA with compact support kernels and background distribution .	16
	Bibliography	19

Chapter 1

Introduction

k Nearest Neighbours (k NN) is one of the oldest and simplest classification methods. It has its origins in an unpublished technical report by ? and since then it became standard textbook material (Russell et al., 1996; Mitchell, 1997; Bishop, 2006). In the last 50 years, k NN was present in most of the machine learning related fields (pattern recognition, statistical classification, data mining, information retrieval, data compression) and it plays a role in many attractive applications (e.g., face recognition, plagiarism detection, vector quantization).

The idea behind k NN is intuitive and straightforward: classify a given point according to a majority vote of its neighbours; the selected class is the one that is the most represented amongst the k nearest neighbours. This is easy to implement and it usually represents a good way to approach new problems or data sets. Despite its simplicity k NN is still a powerful tool, performing surprisingly well in practice, as most of the simple classifiers (Holte, 1993).

Yet there are also other characteristics that make k NN an interesting method. First of all, k NN makes no assumptions about the underlying structure of the data, but lets the data “speak for themselves”. This means that no a priori knowledge is needed beforehand and, more importantly, the accuracy increases with the number of points in the data set (in fact, it approaches Bayes optimality as the cardinality of the training set approaches infinity and k is sufficiently large; Cover and Hart, 1967). Secondly, k NN is able to represent complex functions with non-linear decision boundaries by using only simple local approximations (this behaviour is hardly captured by any parametric method). Lastly, k NN operates in a “lazy” fashion: the training data is just stored and their use is delayed until the testing. The quasi-inexistent training allows to easily add new

training examples.

On the other hand, k NN has some drawbacks that influence both the computational performance and also the quality of its predictions. Firstly, because it has to memorise all the exemplars, the storage requirements are directly proportional with the number of instances in the training set. Furthermore, a serious practical disadvantage is represented by the fact that all the computations are done at testing. The cost for this is also linear in the cardinality of the training set and it is often prohibitive for large data sets. However, more important are the issues related to the accuracy. On one hand, there is not clear how we should choose a dissimilarity metric and how notions as “close”/“far” are defined for our data set. This is non-trivial and usually the standard Euclidean metric is not satisfactory (see Subsection ??, for a more detailed discussion). On the other hand, there have been raised concerns with the usefulness of Nearest Neighbours (NN) methods for high dimensional data (Beyer et al., 1999; Hinneburg et al., 2000). The curse of dimensionality arises for k NN, because the distances become indiscernible for many dimensions; alternatively formulated, for a given distribution the maximum and the minimum distance between points become equal in the limit of dimensions.

All these problems are even more acute nowadays when we have to operate on huge sets of data with many attributes (e.g., images, videos, DNA sequences, etc.). There is an entire literature that tries to come up with possible solutions (some of the most prominent papers are discussed in Section ??). One elegant answer is provided by Goldberger et al. (2004). They proposed a new method, Neighbourhood Component Analysis (NCA), that copes with the drawbacks in a unified manner by learning a low-rank metric. This reduces both the storage and the computational cost, because the algorithm uses the data set projected into a lower subspace, with fewer attributes needed. Also the accuracy is improved, because the label information is used for constructing a proper metric that selects the relevant attributes. However, NCA introduces a consistent training time, because we now have an objective function whose gradient is quadratic in the number of data points that is optimized using iterative methods.

The main aim of this project is to see whether NCA’s training and testing time can be reduced without significant losses in accuracy. We will try to achieve this by making use of k -d trees (a space partitioning structure; Bentley 1975). These have been successfully applied for speeding up k NN at query time (Friedman et al.,

1977) and we will use them in a similar manner for NCA's testing operations. For training, we will apply k -d trees in a slightly different way (as in Deng and Moore, 1995); as k -d trees organize the space, they can be used to group together near points and quickly compute approximations of the gradient and objective function at each iteration. This idea was also followed by Weinberger and Saul (2009) in accelerating a different distance metric technique, Large Margin Nearest Neighbour (LMNN). As an alternative approach, we could interpret NCA as a class-conditional kernel-density estimate and then use different types of kernels instead of the Gaussian ones; we will focus especially on kernels with compact support, because they do not introduce approximation errors and computations can be done more efficiently than in the previous case. If time permits, it would be interesting to experiment with different tree structures, such as ball trees (Omohundro, 1989; Moore, 2000) (which can perform well in higher dimensional spaces) or dual-trees (Gray and Moore, 2003) (a faster representation of the typical k -d tree). Also we could try this approach on other distance metric learning methods, with different objective functions (e.g., Xing et al., 2003).

Chapter 2

Background

2.1 Theoretical background

A distance metric represents a function or a mapping from a pair of inputs to a scalar proportional to the dissimilarity of the inputs. Additionally, in order to be a proper distance, the given function ought to be non-negative, symmetric and it should respect the triangle inequality. The most common and used metric is the standard Euclidean distance. It often appears and is very useful in many geometrical situations, when distances between points need to be calculated. However, it has two major drawbacks that are problematic especially in machine learning applications. First of all, Euclidean distance is sensitive to scaling. Whilst in mathematical problems this does not constitute an issue, in real situations, we may have features that are measured in different units (e.g., seconds, kilograms, etc.) and we will obtain different distances between our data points if we change the scalings on some axis. The other problem is the fact that it does not take into consideration the correlations in the data structure. It can often happen that more attributes reflect the same information present in the data and, consequently, the distance is strongly influenced by those attributes. Take the example of face recognition: there the pixels from the image background are highly correlated and they usually reflect the same information, i.e., the colour of the background.

The standard notation is $d(\mathbf{x}, \mathbf{y})$ and it represents a function that calculates the distance between two inputs \mathbf{x} and \mathbf{y} (column vectors in a D dimensional space $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$). Essentially, the metric maps a pair from $\mathbb{R}^D \times \mathbb{R}^D$ to a real number scalar.

Formally, $d : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ is a metric if for any $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^D$ the following properties hold ?:

- Non-negativity: $d(\mathbf{x}, \mathbf{y}) \geq 0$.
- Distinguishability: $d(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} = \mathbf{y}$.
- Symmetry: $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$.
- Triangle Inequality: $d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y})$.

Now we can relate the previous definition with the most common and used example: Euclidean distance. This is defined as:

$$d_E(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}, \text{ with } \mathbf{x}, \mathbf{y} \in \mathbb{R}^D. \quad (2.1)$$

Unfortunately, the Euclidean distance has two major drawbacks that are problematic especially in machine learning applications¹:

- **Sensitivity to variable scaling.** In geometrical situations, all variables are measured in the same units of length; for some of the real data variability is stronger on certain dimensions than on others. Naturally, we try to compensate the acute difference; so we want components with high variability to receive less weight than those with low variability. We can obtain this by simply rescaling the components with corresponding values $(s_i)_{i=1, \overline{D}}$:

$$\begin{aligned} d(\mathbf{x}, \mathbf{y}) &= \sqrt{\left(\frac{x_1 - y_1}{s_1}\right)^2 + \dots + \left(\frac{x_D - y_D}{s_D}\right)^2} = \\ &= \sqrt{(\mathbf{x} - \mathbf{y})^T S^{-1} (\mathbf{x} - \mathbf{y})} \end{aligned} \quad (2.2)$$

where $S^{-1} = \text{diag}(s_1^2, \dots, s_D^2)$.

- **Invariance to correlated variables.** Euclidean distance does not take into account the correlations in the data structure. For face images, for example, the pixels in the background, usually have the same colour, especially if they are close to each other; if these pixels differ strongly for other picture of the same person, we will be heavily penalized, because their number is large and the distance will be influenced by them. Therefore, we should not take into account all these variables since they reflect

¹<http://matlabdatamining.blogspot.com/2006/11/mahalanobis-distance.html>

the same information (i.e., the color of the background). Intuitively, we want our distance metric to reflect the correlation in the data. This can be easily achieved by replacing S in Equation (2.2) with the covariance matrix of the data.

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T S^{-1} (\mathbf{x} - \mathbf{y})}, \text{ with } S = \text{cov}(\mathcal{X}) \quad (2.3)$$

where \mathcal{X} is the dataset $\mathcal{X} = (\mathbf{x}_i)_{i=1, \dots, N}$.

Equation (2.3) is the standard definition of the Mahalanobis distance. Since we do not care for any variance in the data, but just for that is representative for our particular task (as discussed in Section ??), we will consider different matrices S that are suitable for the given query.

However there are some conditions that S should respect. First of all, we note that $d(\mathbf{x}, 0) = \|\mathbf{x}\| = \sqrt{\mathbf{x}^T S^{-1} \mathbf{x}}$ which imposes $\mathbf{x}^T S^{-1} \mathbf{x} \geq 0, \forall \mathbf{x}$ —this means that S^{-1} must be a positive semi-definite matrix.

We note that we can use Cholesky decomposition and write $S^{-1} = L^T L$, which gives $\|\mathbf{x}\| = \sqrt{(L\mathbf{x})^T (L\mathbf{x})}$. L can be viewed as a linear transformation of the data. Also, using L we can define parametrize our problem by defining a family of metrics over the input space. So the problem of finding a suitable distance metric is virtually equivalent to the problem of finding a good linear transformation and we will discuss these two variants interchangeably.

2.2 Related methods

Xing et al. (2003) proposed learning a Mahalanobis-like distance metric for clustering in a semi-supervised context. There are specified pairs of similar data points and the algorithm tries to find that linear projection that minimizes the distance between these pairs (without collapsing the whole data set to a single point). This approach is formulated as a convex optimization with constraints which can be solved using an iterative procedure, such as Newton-Raphson. The solution is local optima free, but this also means that the method assumes that the data set is uni-modal. This is a strong assumption and it does not coagulate well with the more liberal k NN. Apart from this, the training time is also a problem as the iterative process is costly for large data sets.

Relevant Component Analysis (RCA; Bar-Hillel et al., 2003; ?) is another method that makes use of the labels of the classes in a semi-supervised way to

provide a distance metric. More precisely, RCA uses the so-called chunklet information. A chunklet contains points from the same class, but the class label is not known; data points that belong to the same chunklet have the same label, while data points from different chunklets do not necessarily belong to different classes. The main idea behind RCA is to find a linear transformation which “whitens” the data with respect to the averaged within-chunklet covariance matrix (Weinberger and Saul, 2009). Compared to the method of Xing et al. (2003), RCA has the advantage of presenting a closed form solution, but it has even stricter assumptions about the data: it considers that the data points in each class are normally distributed so they can be described using only second order statistics (Goldberger et al., 2004).

Chapter 3

Neighbourhood component analysis

3.1 General presentation

$$\begin{aligned} f(A) &= \frac{1}{N} \sum_{i=1}^N p_i = \frac{1}{N} \sum_{i=1}^N \sum_{j \in c_i} p_{ij} \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{j \in c_i} \frac{e^{-d_{ij}^2}}{\sum_{k \neq i} e^{-d_{ij}^2}}, \end{aligned} \quad (3.1)$$

where $d_{ij}^2 = \|A\mathbf{x}_i - A\mathbf{x}_j\|^2$.

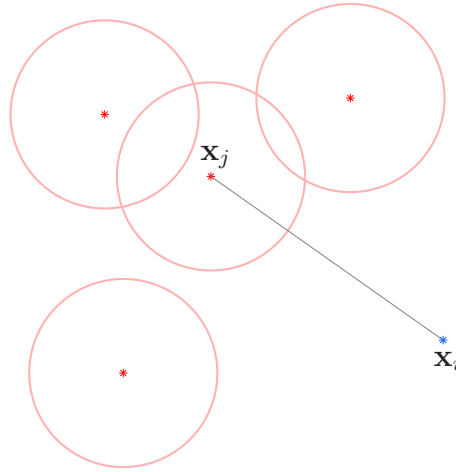
By differentiating with respect to A , we obtain the gradient:

$$\frac{\partial f}{\partial A} = 2A \sum_{i=1}^N \left(p_i \sum_{k=1}^N p_{ik} \mathbf{x}_{ik} \mathbf{x}_{ik}^T - \sum_{j \in c_i} p_{ij} \mathbf{x}_{ij} \mathbf{x}_{ij}^T \right), \quad (3.2)$$

where $\mathbf{x}_{ij} = \mathbf{x}_i - \mathbf{x}_j$.

3.2 Practical notes

- Initialization.
- Optimization methods. Conjugate gradients, gradient descent, bold driver heuristic.
- Numerical issues: log-sum-exp trick, data normalization.
- Regularization: favours 1NN, not sure if optimal thing to do. Might try regularization, as pointed out in (Singh-Miller, 2010). Give equations.
- Further idea: Dimensionality annealing. Regularize each column gradually.



(a) Illustration of the class as a mixture of Gaussians.

Figure 3.1: Formulating NCA as a class-conditional kernel density estimation framework.

3.3 Class-conditional kernel density estimation interpretation

In this section we will present NCA into a class-conditional kernel density estimation framework. This interpretation will allow us to understand what are the assumptions behind NCA. Moreover, this also offers the possibility of altering the model in a suitable way that is efficient for computations. We will see this in the sections 4.4 and 4.5. Similar ideas were previously presented by and , but the following were derived independently and they offer different insights. The following interpretation was inspired by the probabilistic k NN presented by Barber (2011).

We start with the basic assumption that each class can be modelled by a mixture of Gaussians. For each of the N_c data points in class c we consider a Gaussian “bump” centred around it. From a generative perspective, we can view that each point \mathbf{x}_j can generate a point \mathbf{x}_i with a probability given by an isotropic normal distribution with variance σ^2 :

$$p(\mathbf{x}_i|\mathbf{x}_j) = \mathcal{N}(\mathbf{x}_i|\mathbf{x}_j, \sigma^2 \mathbf{I}_D) \quad (3.3)$$

$$= \frac{1}{(2\pi)^{D/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j) \right\}. \quad (3.4)$$

By changing the position of the points through a linear transformation \mathbf{A} , the

probability changes as follows:

$$p(\mathbf{Ax}_i|\mathbf{Ax}_j) \propto \exp \left\{ -\frac{1}{2\sigma^2}(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{A}^T \mathbf{A}(\mathbf{x}_i - \mathbf{x}_j) \right\}. \quad (3.5)$$

We can note that this is similar to the p_{ij} from NCA. Both $p(\mathbf{Ax}_i|\mathbf{Ax}_j)$ and p_{ij} are directly proportional with the same quantity.

Using the mixture of Gaussians assumption, we have that the probability of a point of being generated by class c is equal to the sum of all Gaussians in class c :

$$p(\mathbf{x}_i|c) = \frac{1}{N_c} \sum_{\mathbf{x}_j \in c} p(\mathbf{x}_i|\mathbf{x}_j) \quad (3.6)$$

$$= \frac{1}{N_c} \sum_{\mathbf{x}_j \in c} \mathcal{N}(\mathbf{x}_i|\mathbf{x}_j, \mathbf{I}_D). \quad (3.7)$$

However, we are interested on the inverse probability, given a point \mathbf{x}_i what is the probability of \mathbf{x}_i belonging to class c . We can obtain an expression for $p(c|\mathbf{x}_i)$ using Bayes' theorem:

$$p(c|\mathbf{x}_i) = \frac{p(\mathbf{x}_i|c)p(c)}{p(c|\mathbf{x}_i)} = \frac{p(\mathbf{x}_i|c)p(c)}{\sum_c p(\mathbf{x}_i|c)p(c)}. \quad (3.8)$$

Now if further consider the classes to be equal probable (which might a reasonable assumption if we have no a priori information) we arrive at result that resembles the expression of p_i (see equation):

$$p(c|\mathbf{Ax}_i) = \frac{\frac{1}{N_c} \sum_{\mathbf{x}_j \in c} \exp \left\{ -\frac{1}{2\sigma^2}(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{A}^T \mathbf{A}(\mathbf{x}_i - \mathbf{x}_j) \right\}}{\frac{1}{N_{c'}} \sum_{c'} \sum_{\mathbf{x}_k \in c'} \exp \left\{ -\frac{1}{2\sigma^2}(\mathbf{x}_i - \mathbf{x}_k)^T \mathbf{A}^T \mathbf{A}(\mathbf{x}_i - \mathbf{x}_k) \right\}} \quad (3.9)$$

We are interested in predicting the correct class of the point \mathbf{x}_i . So we try to find that linear transformation \mathbf{A} that maximises the class conditional probability of this point $p(c_i|\mathbf{x}_i)$ to its true class c_i .

$$f(\mathbf{A}) = \sum_i p(c_i|\mathbf{Ax}_i). \quad (3.10)$$

The gradient of the objective function is the following:

$$\frac{\partial f}{\partial \mathbf{A}} = \sum_i \left\{ \frac{\frac{\partial p(\mathbf{Ax}_i|c)}{\partial \mathbf{A}} p(c)}{\sum_c p(\mathbf{x}_i|c)p(c)} - \underbrace{\frac{p(\mathbf{Ax}_i|c)p(c)}{\sum_c p(\mathbf{Ax}_i|c)p(c)}}_{p(c|\mathbf{Ax}_i)} \frac{\sum_c \frac{\partial p(\mathbf{Ax}_i|c)}{\partial \mathbf{A}} p(c)}{\sum_c p(\mathbf{Ax}_i|c)p(c)} \right\}. \quad (3.11)$$

Chapter 4

Reducing the computational cost

- As mentioned in section , evaluating the objective function needs computing all the pairwise distances between the points. Also, the evaluating the gradient is expensive. This is done in $\mathcal{O}(N^2D^2)$ flops. So it is not trivial to successfully use NCA on large data sets.
- This chapter presents a wide palette of ideas and methods that can be applied to speed up the computations. Most of the methods rely on the fact that the learnt metric is low ranked.
- Every method presented can be regarded as an alteration of the original objective function. We basically change our objective function such that the new objective will have a reduced cost.

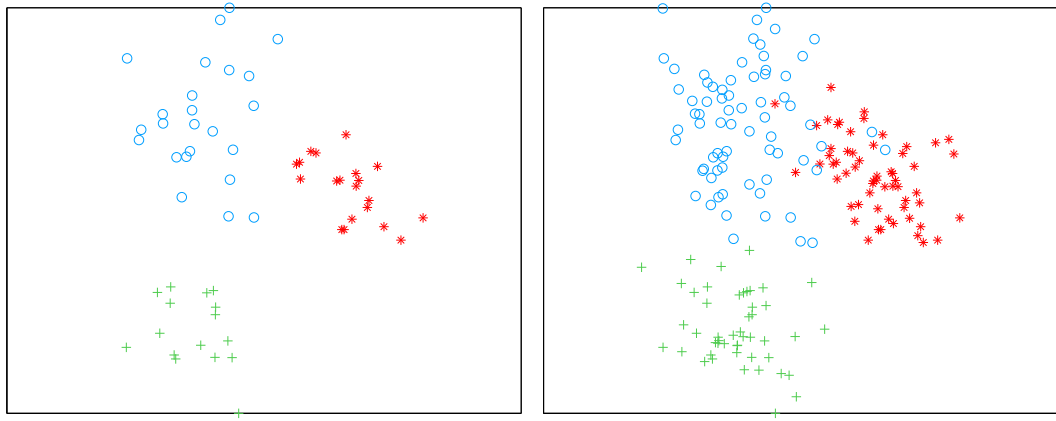
4.1 Sub-sampling

Sub-sampling is the simplest idea that can help speeding up the computations. For the training procedure we use a randomly selected sub-set \mathcal{D}_n of the original data set \mathcal{D} :

$$\mathcal{D}_n = \{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_n}\} \subseteq \mathcal{D}.$$

If n is the size of the sub-set then the cost of the gradient is reduced to $\mathcal{O}(n^2D^2)$. After the projection matrix \mathbf{A} is learnt, it can be applied to the whole data set and all the data points are used for classification.

While easy to implement, this method discards a lot of information available. Also it is affected by the fact the sub-sampled data has a thinner distribution than the real data.



(a) Learnt projection \mathbf{A} on the sub-sampled data set \mathcal{D}_n . (b) The projection \mathbf{A} applied to the whole data set \mathcal{D} .

Figure 4.1: Result of sub-sampling method on wine. There were used one third of the original data set for training, *i.e.*, $n = N/3$. We note that the points that belong to the sub-set \mathcal{D}_n are perfectly separated. But after applying the metric to the whole data there appear different misclassifications. The effects are even more acute if we use smaller sub-sets.

4.2 Mini-batches

The next obvious idea is to use sub-sets in an iterative manner, similar to the stochastic gradient descent method: split the data into mini-batches and train on them successively. Again the cost for one evaluation of the gradient will be $\mathcal{O}(n^2 D^2)$ if the mini-batch consists of n points.

There are different possibilities for splitting the data-set:

1. Random selection. In this case the points are assigned randomly to each mini-batch and after one pass through the whole data set another random allocation is done. As in section 4.1, this suffers from the thin distribution problem. In order to alleviate this and achieve convergence, large-sized mini-batches should be used (similar to Laurens van der Maaten's implementation). The algorithm is similar to Algorithm 4.1, but lines 2 and 3 will be replaced with a simple random selection.
2. Clustering. Constructing mini-batches by clustering ensures that the point density in each mini-batch is conserved. In order to maintain a low computational cost, we consider cheap clustering methods, *e.g.*, farthest point

Algorithm 4.1 Training algorithm using mini-batches formed by clustering

Require: Data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and initial linear transformation \mathbf{A} .

- 1: **repeat**
 - 2: Project each data point using \mathbf{A} : $\mathcal{D}_{\mathbf{A}} = \{\mathbf{A}\mathbf{x}_1, \dots, \mathbf{A}\mathbf{x}_N\}$.
 - 3: Use either algorithm 4.2 or 4.3 on $\mathcal{D}_{\mathbf{A}}$ to split \mathcal{D} into K mini-batches $\mathcal{M}_1, \dots, \mathcal{M}_K$.
 - 4: **for all** \mathcal{M}_i **do**
 - 5: Update parameter: $\mathbf{A} \leftarrow \mathbf{A} - \eta \frac{\partial f(\mathbf{A}, \mathcal{M}_i)}{\partial \mathbf{A}}$.
 - 6: Update learning rate η .
 - 7: **end for**
 - 8: **until** convergence.
-

clustering (FPC; Gonzalez, 1985) and recursive projection clustering (RPC; Chalupka, 2011). Algorithm

FPC gradually selects cluster centres until it reaches the desired number of clusters K . The point which is the farthest away from all the current centres is selected as new centre. The precise algorithm is presented below.

Algorithm 4.2 Farthest point clustering

Require: Data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and K number of clusters.Randomly pick a point that will be the first centre \mathbf{c}_1 .Allocate all the points in the first cluster $\mathcal{M}_1 \leftarrow \mathcal{D}$.**for** $i = 1$ to K **do**

Select the i -th cluster centre \mathbf{c}_i as the point that is farthest away from any cluster centre $\mathbf{c}_1, \dots, \mathbf{c}_{i-1}$.

Move to the cluster \mathcal{M}_i those points that are closer to its centre than to any other cluster centre: $\mathcal{M}_i = \{\mathbf{x} \in \mathcal{D} \mid d(\mathbf{x}; \mathbf{c}_i) < d(\mathbf{x}; \mathbf{c}_j), \forall j \neq i\}$

end for

This computational cost of this method is $\mathcal{O}(NK)$. However, we do not have any control on the number of points in each cluster, so we might end up with very unbalanced clusters. A very uneven split has a couple of obvious drawbacks: too large mini-batches will maintain high cost, while on too small clusters there is not too much to learn.

So, as an alternative, RPC was especially developed to mitigate this problem. It constructs the clusters similarly to how the k -d trees are built (see

section). However instead of splitting the data set across axis aligned direction it chooses the splitting directions randomly (see algorithm 4.3). So, because it uses the median value it will result in similar sized clusters and we can easily control the dimension of each cluster. The complexity of this algorithm is $\mathcal{O}()$.

Algorithm 4.3 Recursive projection clustering

Require: Data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and n size of clusters.

if $n < N$ **then**

New cluster: $i \leftarrow i + 1$.

return current points as a cluster: $\mathcal{M}_i \leftarrow \mathcal{D}$.

else

Randomly select two points \mathbf{x}_j and \mathbf{x}_k from \mathcal{D} .

Project all data points onto the line defined by \mathbf{x}_j and \mathbf{x}_k . (Give equation?)

Select the median value $\tilde{\mathbf{x}}$ from the projected points.

Recurse on the data points above and below $\tilde{\mathbf{x}}$: $\text{RPC}(\mathcal{D}_{>\tilde{\mathbf{x}}})$ and $\text{RPC}(\mathcal{D}_{\leq\tilde{\mathbf{x}}})$.

end if

Note that we are re-clustering in the transformed space after one sweep through the whole data set. There are also other alternatives. For example, we could cluster in the original space periodically or we could cluster only once in the original space. However the proposed variant has the advantage of a good behaviour for a low rank projection matrix \mathbf{A} . Not only that is cheaper, but the clusters resulted in low dimensions by using RPC are closer to the real clusters then applying the same method in a high dimensional space.

Further notes:

- The learning rate can be updated using various heuristics as presented in the section related to optimization.
- The convergence can be tested in various ways: stop when there is not enough momentum in the parameter space, when the function value doesn't vary too much or by using early stopping or a maximum number of iterations. Discuss all of these in practical issues section.

4.3 Stochastic learning

4.4 Approximate computations

4.5 Exact computations

The counterpart of the approximate methods are the exact computations. These are also obtained at some cost. We have to change our model such that it allows efficient exact methods. Again motivated by the rapid decaying squared exponential function, some of the contributions are almost zero. The idea is just to use a compact support kernel instead of the squared exponential kernel. After this is done, there will be a large number of points that will lie outside the support of the kernel and their contribution will be exactly zero instead of a very small value in the previous case.

The kernel needs to be differentiable, so we will use the simplest polynomial compact support kernel that satisfies this requirement. Of course, any other kernel that satisfies this two requirements can be used. This is given by the following expression:

$$k(u) = \begin{cases} \frac{15}{16}(a^2 - u^2)^2 & \text{if } u \in [-a; +a] \\ 0 & \text{otherwise.} \end{cases} \quad (4.1)$$

The expression from equation 4.1 also integrates to one, so it can be used in the kernel density estimation context: $p(\mathbf{x}_i | \mathbf{x}_j) = k(d(\mathbf{x}_i; \mathbf{x}_j))$.

For the simple replacement of the exponential kernel, we have preferred a simplified version of equation 4.1:

$$k(u) = \begin{cases} (1 - u^2)^2 & \text{if } u \in [-1; 1] \\ 0 & \text{otherwise.} \end{cases} \quad (4.2)$$

The width of the kernel's support a will be integrated in the projection matrix \mathbf{A} , similarly as the widths for the Gaussian kernel are absorbed by \mathbf{A} for the classic NCA.

$$p_{ij} = \frac{k(d_{ij})}{\sum_{k \neq i} k(d_{ik})}. \quad (4.3)$$

Objective function will be:

$$f_{\text{CS}}(\mathbf{A}) = \sum_i \sum_{j \in c_i} p_{ij} \quad (4.4)$$

The gradient of the kernel:

$$\frac{\partial}{\partial \mathbf{A}} k(d_{ij}) = \frac{\partial}{\partial \mathbf{A}} [(1 - d_{ij}^2)^2 \cdot \mathbf{I}(|d_{ij}| \leq 1)] \quad (4.5)$$

$$= -4\mathbf{A}(1 - d_{ij}^2)\mathbf{x}_{ij}\mathbf{x}_{ij}^T \cdot \mathbf{I}(|d_{ij}| \leq 1), \quad (4.6)$$

where $d_{ij}^2 = d(\mathbf{Ax}_i; \mathbf{Ax}_j)$, $\mathbf{x}_{ij} = \mathbf{x}_i - \mathbf{x}_j$ and $\mathbf{I}(\cdot)$ is the indicator function that return 1 when its argument is true and 0 when its argument is false.

The gradient of the new objective function:

$$\frac{\partial f_{\text{CS}}}{\partial \mathbf{A}} = 2\mathbf{A} \sum_{i=1}^N \left(p_i \sum_{k=1}^N \frac{p_{ik}}{1 - d_{ik}^2} \mathbf{x}_{ik}\mathbf{x}_{ik}^T - \sum_{j \in c_i} \frac{p_{ij}}{1 - d_{ij}^2} \mathbf{x}_{ij}\mathbf{x}_{ij}^T \right), \quad (4.7)$$

The main concern with this method is what happens when points lie outside the compact support of any other point in the data set. So care must be taken at initialisation. One way would be to initialise with a very small scale \mathbf{A} . However, this means that there is no gain in speed for at least the first iterations. It might be better to use the principal components for initialisation.

- Speedings comes from the fact only few gradient components need to be computed. A nice way of further accelerating this is again via k -d trees. This time they can be used for range search.
- Note regarding the gradient: it was written in this form for simplicity. However, in an implementation it is not recommended to divide each term by $1 - d_{ij}^2$; there might results cases that are not determined. Better is to compute the kernel values in two steps.

4.6 NCA with compact support kernels and background distribution

- An extended variant of the previous scenario. We take care of the points that might remain unallocated by using a background distribution. This extension comes on naturally after the CC-KDE formulation. We can change the assumptions regarding the underlying distribution of a class $p(\mathbf{x}_i|c)$.

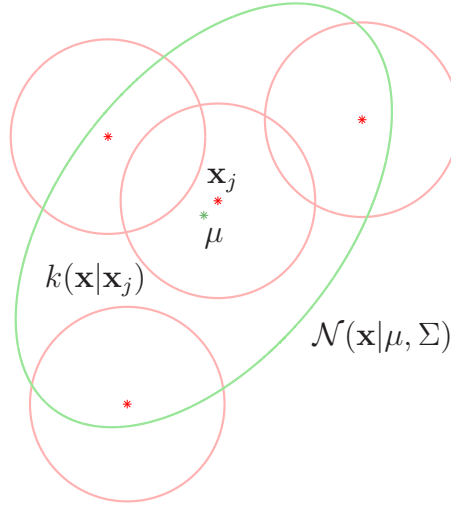


Figure 4.2: Neighbourhood component analysis with compact support kernels and background distribution. Main assumption: each class is a mixture of compact support distributions $k(\mathbf{x}|\mathbf{x}_j)$ plus the normal background distribution $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

So, in this case, we consider that a point can be generated from a class by either a compact distribution from each point $k_{\text{CS}}(\mathbf{x}_i|\mathbf{x}_j)$ or by a global normal distribution that is specific to the whole class $\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$. So, the class-conditional distribution is given by the sum of these components:

$$p(\mathbf{x}_i|c) = \beta \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) + (1 - \beta) \frac{1}{N_c} \sum_{j \in c} k_{\text{CS}}(\mathbf{x}_i|\mathbf{x}_j), \quad (4.8)$$

where β is the mixing coefficient between the background distribution and the compact support model and $\boldsymbol{\mu}_c$ and $\boldsymbol{\Sigma}_c$ are the sample mean and covariance of the class c .

- It is unclear how we can choose the mixing proportion β . By setting $\beta = \frac{1}{N_c+1}$ we will give equal weights to the background distribution as to the compact-support distribution. Setting it too large means that the model will favour convex classes. On one hand, this might diminish NCA's power. One of its strengths lies in the ability to work on non-convex classes. On the other hand, some of the classes This can be fitted as a parameter during the optimization process. The gradient with respect to β can be easily derived and it is easy to evaluate it as the quantities required for this should also be computed for the function evaluation.
- This model can be used by plugging equation 4.8 into the set of equations 3.8, 3.10 and 3.11 from section 3.3.

- We give here the gradient for the background distribution with respect to the point positions (the full derivation can be found in the appendix).
- It might be useful to note that projecting the points $\{\mathbf{x}_i\}_{i=1}^N$ into a new space $\{\mathbf{A}\mathbf{x}_i\}_{i=1}^N$ will change the sample mean $\boldsymbol{\mu}_c$ to $\mathbf{A}\boldsymbol{\mu}_c$ and the sample covariance $\boldsymbol{\Sigma}_c$ to $\mathbf{A}\boldsymbol{\Sigma}_c\mathbf{A}^T$. Hence, we have:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{A}} \mathcal{N}(\mathbf{A}\mathbf{x}_i | \mathbf{A}\boldsymbol{\mu}_c, \mathbf{A}\boldsymbol{\Sigma}_c\mathbf{A}^T) &= \mathcal{N}(\mathbf{A}\mathbf{x}_i | \mathbf{A}\boldsymbol{\mu}_c, \mathbf{A}\boldsymbol{\Sigma}_c\mathbf{A}^T) \\ &\times \{ -(\mathbf{A}\boldsymbol{\Sigma}_c\mathbf{A}^T)^{-1} \mathbf{A}\boldsymbol{\Sigma}_c + \mathbf{v}\mathbf{v}^T \mathbf{A}\boldsymbol{\Sigma}_c - \mathbf{v}(\mathbf{x} - \boldsymbol{\mu}_c)^T \}, \end{aligned} \quad (4.9)$$

where $\mathbf{v} = (\mathbf{A}\boldsymbol{\Sigma}_c\mathbf{A}^T)^{-1} \mathbf{A}(\mathbf{x} - \boldsymbol{\mu}_c)$.

Bibliography

- Bar-Hillel, A., Hertz, T., Shental, N., and Weinshall, D. (2003). Learning distance functions using equivalence relations. In *Proceedings of the Twentieth International Conference on Machine Learning*, volume 20, page 11.
- Barber, D. (2011). *Bayesian Reasoning and Machine Learning*. Cambridge University Press. In press.
- Bentley, J. L. (1975). Multidimensional binary search trees used for associative searching. *Commun. ACM*, 18:509–517.
- Beyer, K., Goldstein, J., Ramakrishnan, R., and Shaft, U. (1999). When is “nearest neighbor” meaningful? *Database TheoryICDT99*, pages 217–235.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- Chalupka, K. (2011). Empirical evaluation of Gaussian Process approximation algorithms.
- Cover, T. and Hart, P. (1967). Nearest neighbor pattern classification. 13:21–27.
- Deng, K. and Moore, A. (1995). Multiresolution instance-based learning. In *Proceedings of the Twelfth International Joint Conference on Artificial Intelligence*, pages 1233–1239, San Francisco. Morgan Kaufmann.
- Friedman, J. H., Bentley, J. L., and Finkel, R. A. (1977). An algorithm for finding best matches in logarithmic expected time. *ACM Trans. Math. Softw.*, 3:209–226.
- Goldberger, J., Roweis, S., Hinton, G., and Salakhutdinov, R. (2004). Neighbourhood components analysis. In *Advances in Neural Information Processing Systems*. MIT Press.
- Gonzalez, T. (1985). Clustering to minimize the maximum intercluster distance. *Theoretical Computer Science*, 38:293–306.
- Gray, A. and Moore, A. (2003). Nonparametric density estimation: Toward computational tractability. In *SIAM International Conference on Data Mining*, volume 2003.

- Hinneburg, E., Aggarwal, C., Keim, D., and Hinneburg, A. (2000). What is the nearest neighbor in high dimensional spaces?
- Holte, R. (1993). Very simple classification rules perform well on most commonly used datasets. *Machine learning*, 11(1):63–90.
- Mitchell, T. M. (1997). *Machine Learning*. McGraw-Hill.
- Moore, A. (2000). *The anchors hierarchy: Using the triangle inequality to survive high dimensional data*. Citeseer.
- Omohundro, S. (1989). *Five balltree construction algorithms*.
- Russell, S. J., Norvig, P., Candy, J. F., Malik, J. M., and Edwards, D. D. (1996). *Artificial intelligence: a modern approach*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- Singh-Miller, N. (2010). *Neighborhood Analysis Methods in Acoustic Modeling for Automatic Speech Recognition*. PhD thesis, Massachusetts Institute of Technology.
- Weinberger, K. and Saul, L. (2009). Distance metric learning for large margin nearest neighbor classification. *The Journal of Machine Learning Research*, 10:207–244.
- Xing, E., Ng, A., Jordan, M., and Russell, S. (2003). Distance metric learning with application to clustering with side-information. In *Advances in Neural Information Processing Systems*, pages 521–528.