

Math 521 Classification: Cats and Dogs

Due on Tuesday, May 15, 2018

Kristin Holmbeck and Debbie Tonne

Contents

Theory	2
Introduction	2
Preprocessing	2
Filtering	2
Dimension Reduction	2
Kernel Discriminant Analysis	2
Kernel Trick	4
KDA with the Kernel Trick	4
KDA on concentric circle data	5
KDA on parabolic data	7
Classification	9
Singular Values	9
Results	9
Dogs and Cats	9
Classification Types	9
Code	10

List of Figures

1	Concentric Data	3
2	Figure 1 data with LDA projection vector	6
3	Figure 1 LDA classification	6
4	Figure 1 KDA classification	7
5	Parabolically-separable data	7
6	Figure 5 data with LDA projection vector	8
7	Figure 5 LDA classification	8
8	Figure 5 KDA classification	9

Theory

Introduction

The project we present in this report involves properly classifying two data sets successfully. In this context, the data sets are images of dogs and cats, but the same ideas and algorithms can be successfully applied to other data sets, such as sound waves. Since we are working with images, some preprocessing methods will be explored to add uniformity or variance to the data sets.

Preprocessing

Image *preprocessing* typically involves filtering or computing the Fourier transform of an image prior to analysis. To this end, we will discuss some basics, beginning with filtering.

Filtering

Image filtering uses a *mask* matrix on subsets of an image to perform operations. One filter example is the averaging filter: Given an $m \times n$ mask size, the mask m_a will be

$$m_a = \frac{1}{mn} \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{bmatrix}$$

where the filtering operation involves an $m \times n$ neighborhood around each pixel of the original image matrix A .

Dimension Reduction

Kernel Discriminant Analysis

Before detailing the kernel classification method, we will provide an intuitive example for explaining why one might want to use KDA over LDA. First, consider the toy problem of two concentric circles of data (Figure 1).

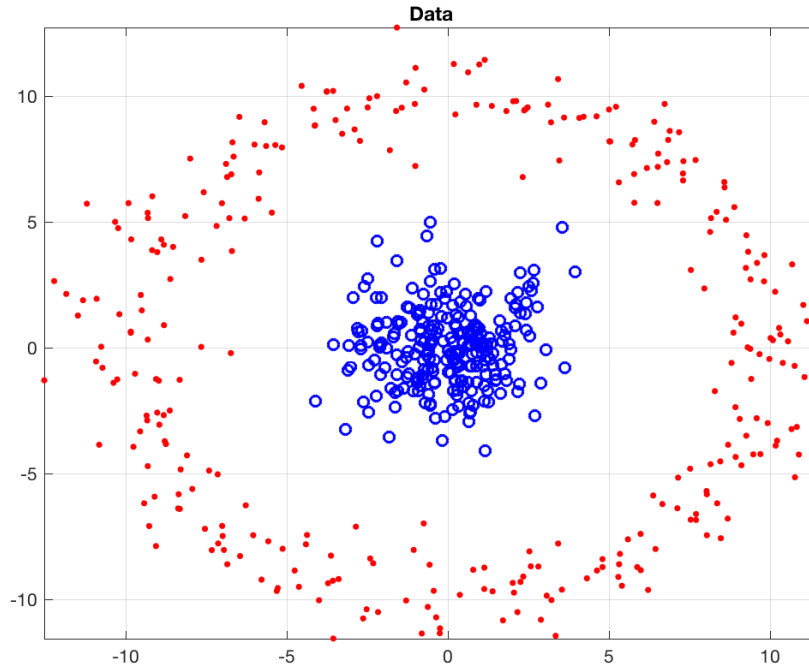


Figure 1: Concentric Data

There linear way to separate the data, but the classes clearly have a defined separation. Kernel Discriminant Analysis (KDA) utilizes the ideas of LDA but with the data set X mapped onto a new *feature* space \mathcal{F} where the data has a linear relationship in \mathcal{F} . For the remainder of this paper, we will present KDA in the context of two classes, although it can be generalized to n classes.

Following the LDA method, suppose our two-class data is given by $X = X_1 \cup X_2$ where $X_1 = \{x_i\}_{i=1}^{l_1}$ and $X_2 = \{x_j\}_{j=1}^{l_2}$. Now, let Φ be a nonlinear mapping to some feature space \mathcal{F} , that is, we take a vector $x \in X$ and map it using $\Phi(x) \in \mathcal{F}$. From the LDA algorithm, we need to maximize

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

However, using the new space, we must maximize

$$J(w) = \frac{w^T S_B^\Phi w}{w^T S_W^\Phi w} \quad (1)$$

where

$$S_B^\Phi = (m_1^\Phi - m_2^\Phi)(m_1^\Phi - m_2^\Phi)^T \quad \text{and} \\ S_W^\Phi = \sum_{i=1}^2 \sum_{x \in X_i} (\Phi(x) - m_i^\Phi)(\Phi(x) - m_i^\Phi)^T$$

are the between-class and within-class scatter matrices, respectively, in the \mathcal{F} space, and $m_i = \frac{1}{l_i} \sum_{j=1}^{l_i} \Phi(x_j^{(i)})$, the mean of the i^{th} class. Furthermore, note that we are now finding the projection $w \in \mathcal{F}$.

Kernel Trick

The kernel trick boils down to using a linear classifier to solve a non-linear problem.

A feature map is a map $\Phi : X \rightarrow \mathcal{F}$, where \mathcal{F} is what we call the feature space. Every feature map defines a kernel

$$\kappa(x, y) = \langle \Phi(x), \Phi(y) \rangle$$

where κ is clearly symmetric and positive-definite. In the context of linear algebra, the kernel is the space equivalent to the null space. In the statistical context, the kernel is used as a measure of similarity. In particular, the kernel function κ defines the distribution of similarities of points around a given point x , $\kappa(x, y)$ denotes the similarity of point x with another given point y .

Explicitly computing the mappings of a function $\Phi(x)$ onto \mathcal{F} can become intractable quick. To that end, we instead compute the inner products between the images (images in the linear algebra sense) of all pairs of data in the feature space.

For any x, \hat{x} in X , some kernel functions $\kappa(x, \hat{x})$ can be expressed as an inner product in another space V . In other words, $\kappa : X \times X \rightarrow \mathbb{R}$ and $\Phi : X \rightarrow V$.

$$\kappa(x, \hat{x}) = \langle \Phi(x), \Phi(\hat{x}) \rangle_V$$

KDA with the Kernel Trick

Noting that $w \in \mathcal{F}$, and using the theory of reproducing kernels [3], w lies in the span of all training samples in \mathcal{F} . Hence,

$$w = \sum_{i=1}^l \alpha_i \Phi(x_i^{(j)})$$

Multiplying w^T by the mean,

$$\begin{aligned} w^T m_i^\Phi &= \frac{1}{l_i} \sum_{j=1}^l \alpha_j \Phi^T(x_j) \sum_{k=1}^{l_i} \Phi(x_k^{(i)}) \\ &= \frac{1}{l_i} \sum_{j=1}^l \sum_{k=1}^{l_i} \alpha_j \Phi^T(x_j) \Phi(x_k^{(i)}) \\ &= \frac{1}{l_i} \sum_{j=1}^l \sum_{k=1}^{l_i} \alpha_j \kappa(x_j, x_k^{(i)}) \end{aligned}$$

Defining $(M_i)_j := \frac{1}{l_i} \sum_{k=1}^{l_i} \kappa(x_j, x_k^{(i)})$,

$$w^T m_i = \alpha^T M_i$$

For the details of S_B^Φ ,

$$S_B^\Phi = (m_1^\Phi - m_2^\Phi) (m_1^\Phi - m_2^\Phi)^T = m_1^\Phi m_1^{\Phi T} - m_1^\Phi m_2^{\Phi T} - m_2^\Phi m_1^{\Phi T} + m_2^\Phi m_2^{\Phi T}$$

Looking at just one of the terms,

$$\begin{aligned} &= \frac{1}{l_a l_b} \sum_{j=1}^{l_a} \sum_{k=1}^{l_b} \Phi \left(x_j^{(a)} \right) \Phi^T \left(x_k^{(b)} \right) \\ &= \frac{1}{l_a l_b} \sum_{j=1}^{l_a} \sum_{k=1}^{l_b} \kappa \left(x_j^{(a)}, x_k^{(b)} \right) \end{aligned}$$

With this, we can now write

$$S_B^\Phi = (M_1 - M_2)(M_1 - M_2)^T \quad \text{with} \quad (M_i)_j := \frac{1}{l_i} \sum_{k=1}^{l_i} \kappa \left(x_j, x_k^{(i)} \right)$$

where $M_i \in \mathbb{R}^{l \times l}$ (**VERIFY THIS**). Thus, the numerator of (1) can be written as $\alpha^T M \alpha$.

Following the same logic, the denominator $w^T S_W^\Phi w$ can be written as $\alpha^T N \alpha$ where

$$\begin{aligned} N &:= \sum_{j=1}^2 K_j (I - \mathbf{1}_{l_j}) K_j^T \\ (K_j)_{nm} &:= \kappa \left(x_n, x_m^{(j)} \right), K_j \in \mathbb{R}^{l \times l_j} \\ \mathbf{1}_{l_j} &:= \frac{1}{l_j} \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{bmatrix} \in \mathbb{R}^{l_j \times l_j} \end{aligned}$$

Then, (1) can be rewritten as

$$J(\alpha) = \frac{\alpha^T M \alpha}{\alpha^T N \alpha} \tag{2}$$

and a *new* projection of x onto w is given by

$$w \cdot \Phi(x) = \sum_{i=1}^l \alpha_i \kappa(x_i, x) \tag{3}$$

KDA on concentric circle data

Let us revisit the concentric data problem from Figure 1, and compare the classification of LDA to KDA. For LDA, we obtain the failed classification (Figure 2) and the projection vector along with the data in Figure 3.

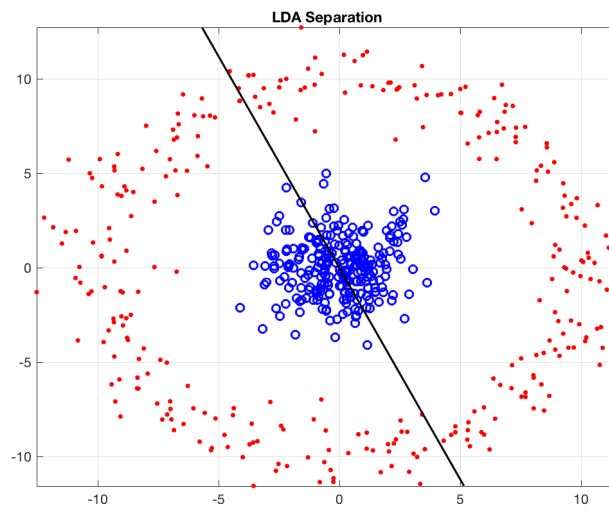


Figure 2: Figure 1 data with LDA projection vector

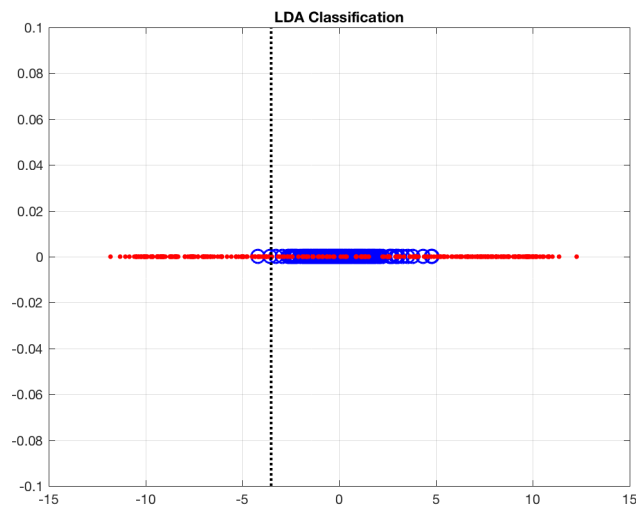


Figure 3: Figure 1 LDA classification

As described above, the kernel conversion is kind of tricky, and thus we cannot plot the equivalent projection vector. However, the successful classification is shown in Figure 4.

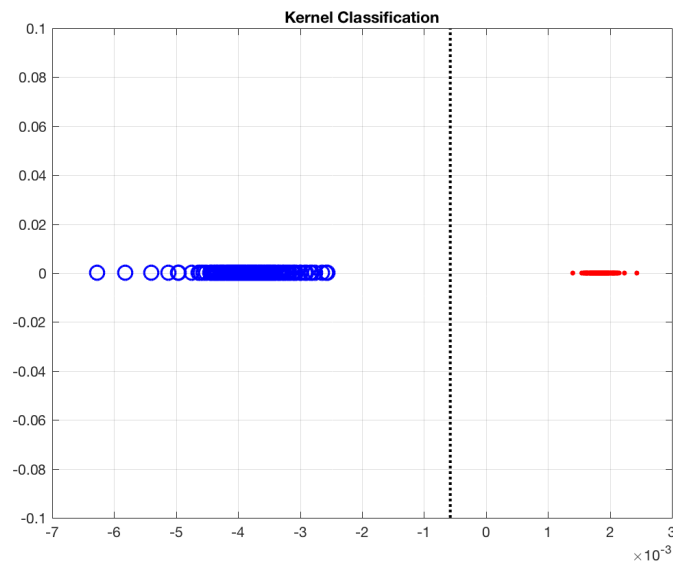


Figure 4: Figure 1 KDA classification

KDA on parabolic data

Another data set is shown in Figure 5. Again, this is easily separable, but it's clear that the separation is nonlinear. Again, for LDA, we obtain the failed classification (Figure 6) and the successful classification in Figure 8.

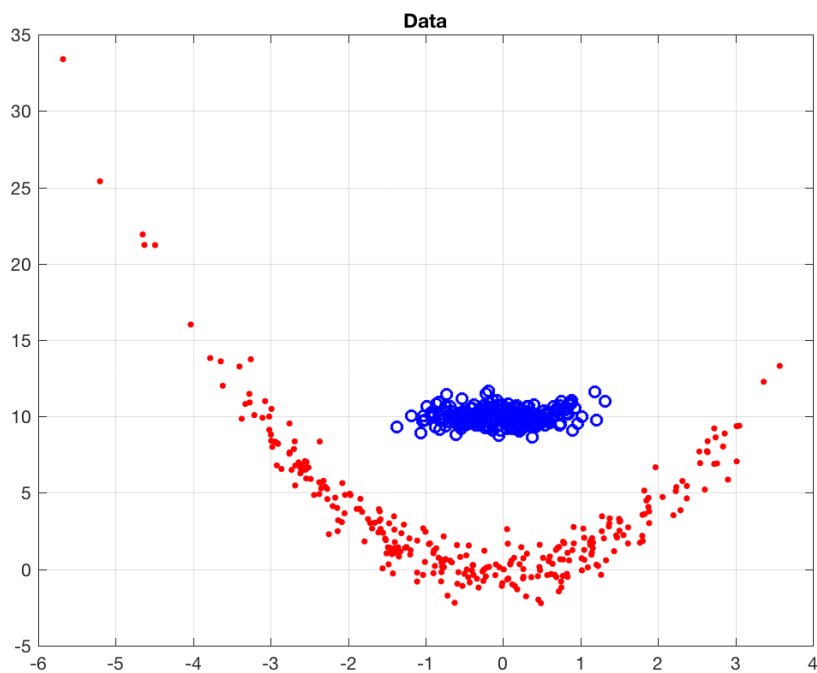


Figure 5: Parabolically-separable data

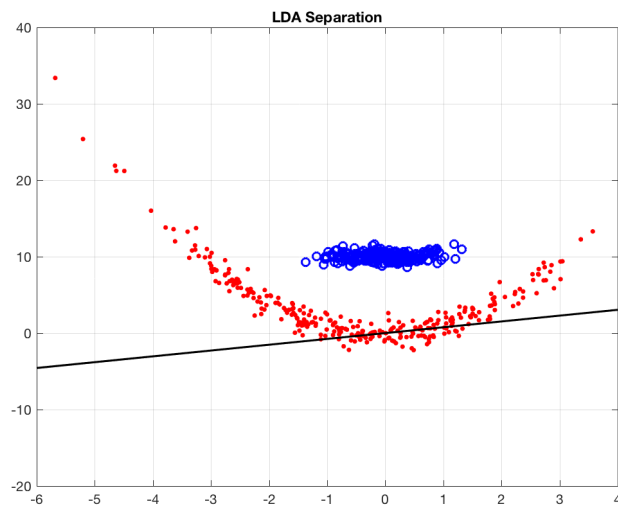


Figure 6: Figure 5 data with LDA projection vector

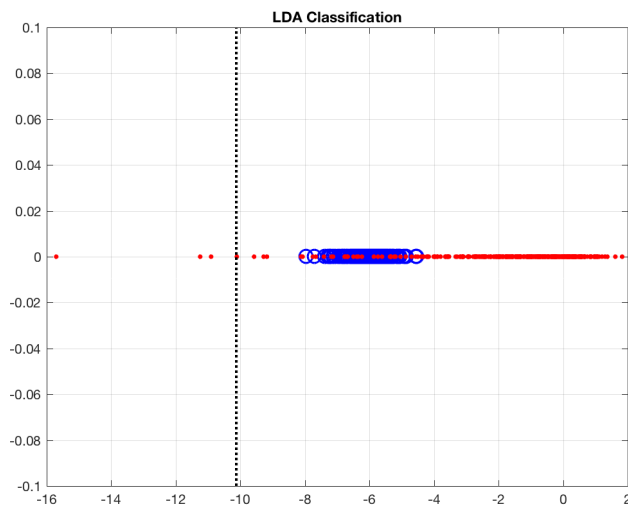


Figure 7: Figure 5 LDA classification

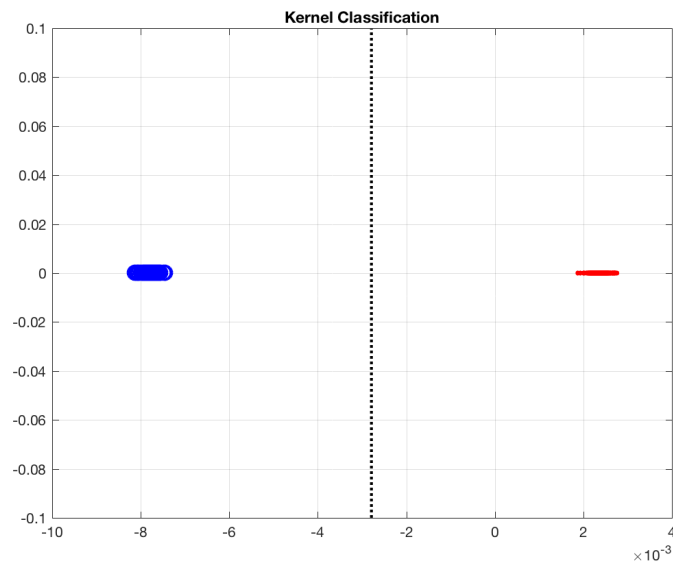


Figure 8: Figure 5 KDA classification

Classification

Singular Values

The Singular Value Decomposition (SVD) is an important first step in classification.

Results

Dogs and Cats

Classification Types

Code

References

- [1] Chang, Jen-Mei. *Matrix Methods for Geometric Data Analysis and Recognition*. 2014.
- [2] S. Mika, G. Ratsch, J. Weston, B. Scholkopf, and K. Muller. Fisher discriminant analysis with kernels. In *Proc. IEEE Neural Networks for Signal Processing Workshop*, pages 4148. IEEE Computer Society Press, 1999.
- [3] Aronszajn, N. “Theory of Reproducing Kernels.” *Transactions of the American Mathematical Society* 68, no. 3 (1950): 337-404. doi:10.2307/1990404.