Fadhla Mohamed Mutua

# SM3201434

1. Introduction

This study investigates three machine learning parts with simulated data from random numbers. The random number generation comes from the use of numpy np.linspace(-5, 5, train\_points) for linear ridge and kernel ridge (where train\_points is 20 for traning and 1000 for testing) regression and sklearn sklearn.datasets.make\_circles for PCA. The first part compares linear ridge regression (using both Gaussian and polynomial kernels) to assess regression performance under different model assumptions and kernel formulations. The second part compares traditional Principal Component Analysis (PCA) with Kernel PCA (KPCA) to determine how both methods perform in reducing dimensionality on structured data. A third "experimental" section repeats the PCA versus KPCA comparison on a dataset made using sklearn.datasets.make\_classification, highlighting the sensitivity of kernel-based methods.

2. Ridge Regression vs. Kernel Ridge Regression

## 2.1 Objective and Methodology

• Objective: The goal was to determine whether kernelizing the ridge regression approach (with either Gaussian or polynomial kernels) could capture non-linear relationships in data better than standard linear ridge regression. Methodology:

■ Split of training dataset: The dataset consists of 20 training data values in a range from -5 to 5 and whose y value is gotten from the formula:

 $\circ \hspace{0.2cm} y_{train} = (X_{train} + 4) * (X_{train} + 1) * (np. \hspace{0.05cm} cos(X_{train}) - 1) * (X_{train} - 3) + eps$  $\circ$  With  $eps = \text{np.random.normal(0, 1, train_points)}$  an array of random movements (obtained from the normal distribution) from the true value of y

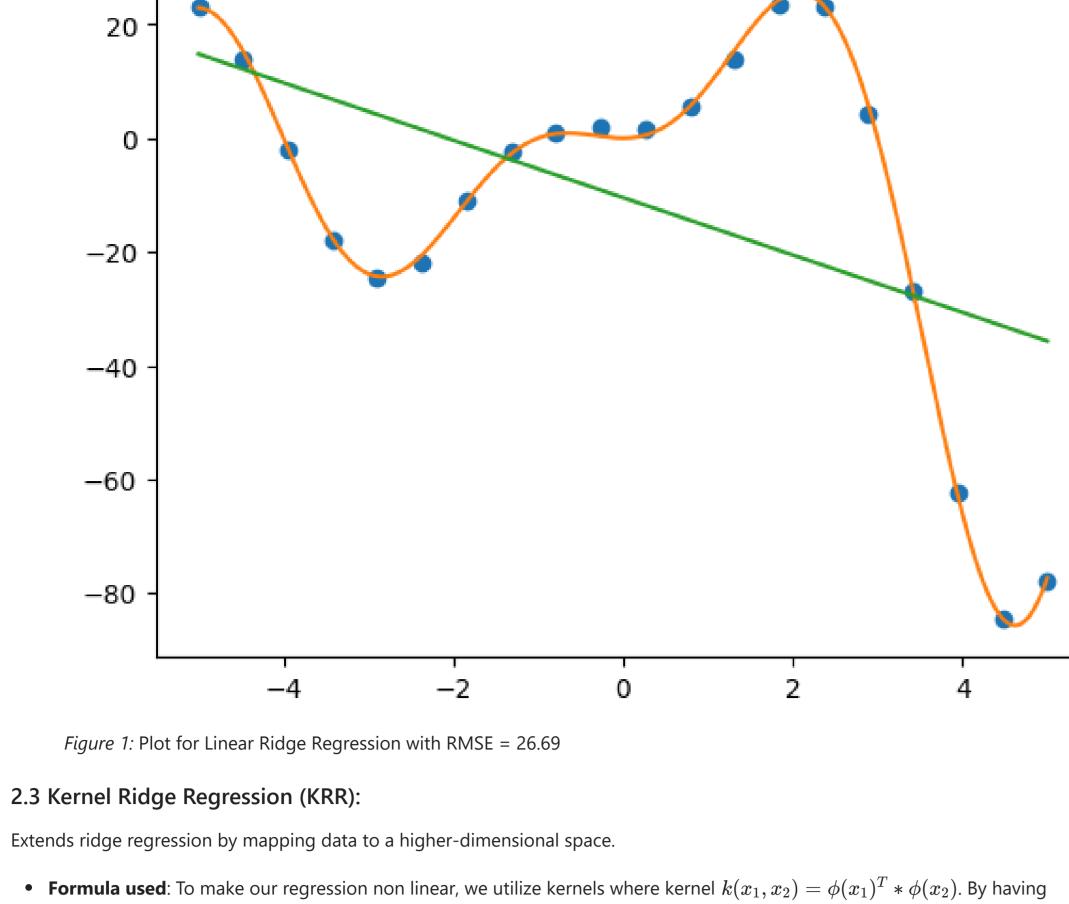
■ Split of testing dataset: The dataset consists of 1000 testing data valus in a range from -5 to 5 and whose y value is gotten from the formula:  $\circ \;\; y_{test} = (X_{test} + 4) * (X_{test} + 1) * (np. cos(X_{test}) - 1) * (X_{test} - 3)$ 2.2 Linear Ridge Regression:

A regularized linear model that imposes penalties on the coefficient sizes to avoid overfitting. • Formula used: Given the formula for linear regression  $Y = X * \omega + b$ , by including the bias term as a feature (in this case all 1), we can express the model as

•  $Y = X_{training\_bias} * \omega$  from which  $ullet \ \omega = (X_{training\_bias}^T * X_{training\_bias} + \lambda * I)^{-1} * X_{training\_bias}^T * Y$ 

lacktriangle The prediction is then formulated as the line  $Y_{predict} = X_{training\_bias} * \omega$ • **Observation**: From the data, we observe that the prediction line  $y_{prediction}$  does not accurately represent the data. In fact, the high RMSE suggests significant underfitting.

Prediction for Linear Ridge regression - RMSE: 26.69 - R2: 0.25



 $\Phi(X) =$ • And replacing all X in the previous linear ridge regression with  $\Phi(X)$  we get:  $\omega = \Phi(X)^T * (\Phi(X)^T * \Phi(X) + \lambda * I)^{-1} * Y$ Then we get

 $Y_{predict} = \phi(x)^T * \omega = \phi(x)^T * \Phi(X)^T * (\Phi(X)^T * \Phi(X) + \lambda * I)^{-1} * Y$ • We can then replace:

 $\Phi(X)^T*\Phi(X)=K$  $\sum_{i=1}^n \phi(x)^T * \phi(x_i) = \sum_{i=1}^n k(x,x_i)$ 

 $(K)_{i,j}=k(x_i,x_j)$  $\phi(x)^T*\Phi(X)^T=\sum_{i=1}^n\phi(x)^T*\phi(x_i)$ 

• From which:

 $lpha = (K + \lambda * I)^{-1} * Y$  $Y_{predict} = \sum_{i=1}^{n} lpha_i * k(x,x_i)$ 

Two kernels were tested:

2.3.1 Gaussian (RBF) Kernel • Formula used: Applying the Gaussian Kernel we have: •  $k(x,x')=e^{-\|x-x'\|^2/(2*\sigma^2)}$  with  $\sigma>0$ 

• **Observation**: By applying a grid search on sigma\_grid = [0.01, 0.1, 1, 5, 10] and lambda\_grid = [0.01, 0.1, 1, 5, 10]

30

25 -

20 RW 15 -

10 -

lambda=5 -- lambda=10

-- lambda=0.01

--- lambda=0.1

→ lambda=1

Figure 2: Plot for RMSE vs Sigma Gauss

ullet We find that the best fit curve has  $\sigma=1$  and  $\lambda=0.01$ 

Best Prediction Gauss - RMSE: 0.80 20 -20 -40 -60 Training data Test data Prediction (param=10, lam=5) -80 MSE=0.64, MAE=0.67, R<sup>2</sup>=1.00 -2 Figure 3: Plot for Kernel Ridge Rigression (Gauss) with RMSE = 0.80

RMSE vs param for each lambda (Gauss kernel)

20

• Formula used: Applying the Polynomial Kernel we have:

2.3.2 Polynomial Kernel

30

25 -

10

 $lacksquare k(x,x') = (x^T x' + 1)^{\sigma}$ 

ullet We find that the best fit curve has  $\sigma=10$  and  $\lambda=5$ 

Figure 4: Plot for RMSE vs Sigma

• **Observation**: By applying a grid search on sigma\_grid = [0.01, 0.1, 1, 5, 10] and lambda\_grid = [0.01, 0.1, 1, 5, 10]

RMSE vs param for each lambda (poly kernel)

param

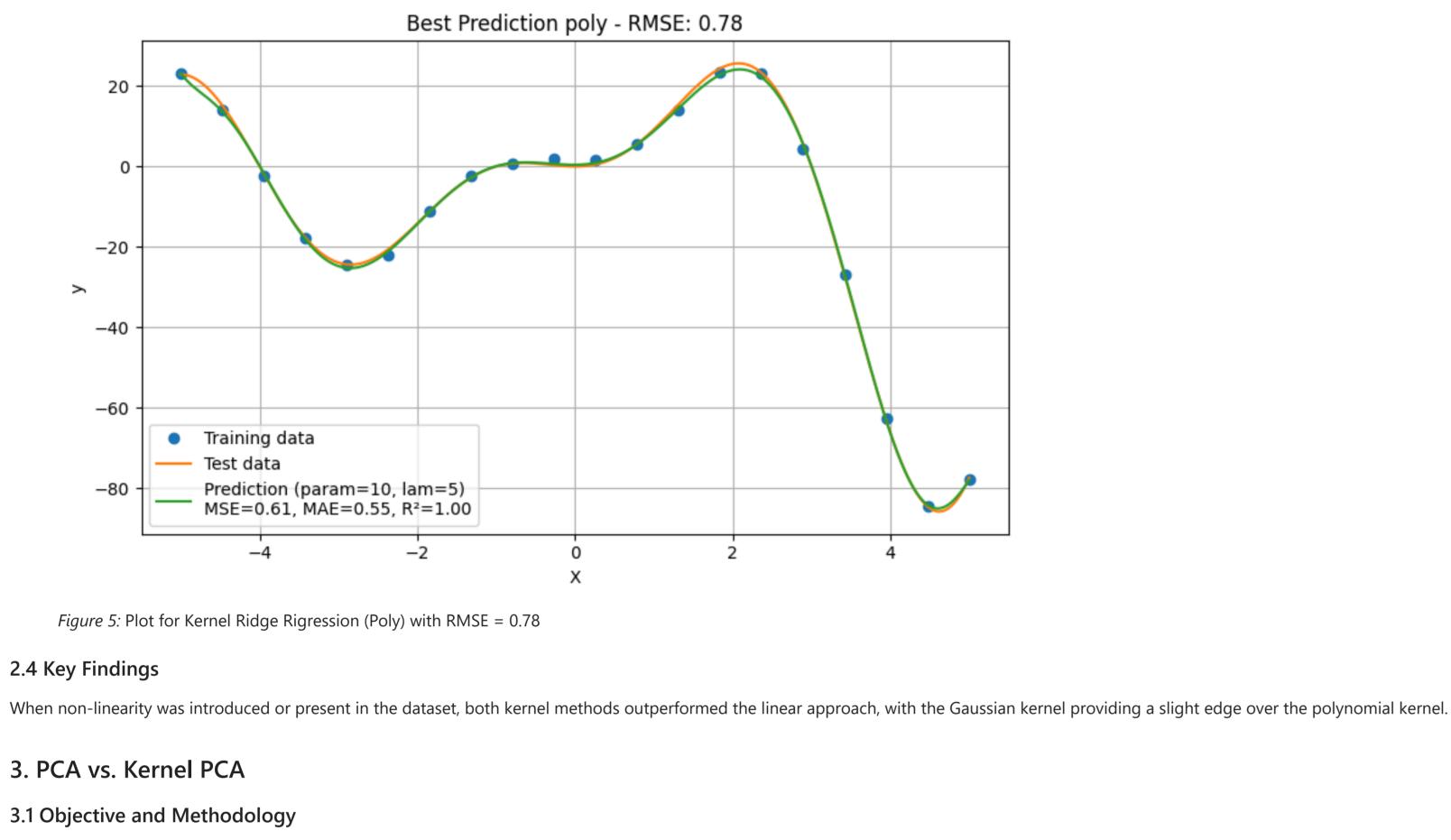
→ lambda=0.01

lambda=0.1

→ lambda=1

→ lambda=5

→ lambda=10



A linear dimensionality reduction technique that projects data onto a subspace spanned by the principal components (directions of maximal variance).

### • **Objective**: To compare the effectiveness of standard PCA against KPCA in terms of capturing the variance and revealing underlying structure in the data for linear SVM. Methodology: • Gets the dataset from sklearn.datasets.make\_circles with 1000 samples and divides it into 20% training and 80% testing with 15% of them being noise. For consistency the sklearn seed is set to 0

1.0

KPCA Performance:

0.6

0.4

0.2

confirmed by the classification report accuracy.

3.2 PCA

 PCA Performance: ■ Using sklearn.decomposition.PCA, we perform dimensionality reduction on X\_train and visualize the result. Although the dataset is correctly labeled, the projected shape reveals two concentric circles; one enclosed within the other. This circular structure highlights a key limitation: the data is not linearly separable in the reduced space, which significantly reduces the effectiveness of a linear SVM in classifying the points correctly as proven with the classification report on accuracy: **PCA Projection** 

0.30

0.97

0.63

0.63

■ Using sklearn.decomposition.KernelPCA with gamma=5, we perform dimensionality reduction on X\_train and visualize the result. The projected data is well-separated and nearly linearly separable, which significantly improves the performance of a linear SVM, as

1.00

0.99

0.99

0.99

0.99

0.99

125

125

250

250

250

0.72

0.63

0.59

0.59

125

125

250

250

250

2nd princip -0.5-1.0-1.00.0 1.0 1.5 -0.50.5 1st principal component Figure 6: Plot for PCA projection Precision Recall F1-score Support Class 0 0.90 Class 1 0.58 Accuracy 0.74 Macro avg Weighted avg 0.74 3.2 KPCA

2nd principal component o °o

A non-linear extension of PCA that first applies a kernel transformation. The Gaussian or polynomial kernel can help uncover non-linear manifolds hidden in the data.

Kernel PCA Projection (RBF Kernel)

0 -0.40 -0.6-0.4 -0.2 0.2 0.6 0.0 0.4 1st principal component Figure 7: Plot for KPCA projection Precision Recall F1-score Support Class 0 0.98 Class 1 0.98 1.00 Accuracy 0.99 0.99 Macro avg 0.99 Weighted avg 0.99 3.3 Key Findings In scenarios where the data had a non-linear structure, KPCA delivered more insightful lower-dimensional representations. 4. PCA vs. Kernel PCA Part 2 4.1 Objective and Methodology • Methodology: ■ We generate a dataset using sklearn.datasets.make\_classification with 1000 samples and 15% label noise. The data is split into 80% testing and 20% training, with random\_state=0 set for consistency. 4.2 PCA PCA Performance: Using sklearn.decomposition.PCA, we perform dimensionality reduction on X\_train and visualize the result. We observe that the dataset becomes almost linearly separable, as confirmed by the SVM classification report. Kernel PCA Projection (RBF Kernel) 0

odwo	• %			8 8 8	806	0					
cipal c		<b>%</b> % %			0 0						
ind prin				• • •							
7						3. g					
		•	3			•••					
-2			300	•							
				•							
			o								
-4											
				•							
	-4	-2	0 1st princip	) eal component	2	4					
Figure 8: Plo	ot for PCA projectior	١									
							Precision	Recall	F1-score	Support	
						Class 0	0.98	1.00	0.99	124	
									0.55	127	
						Class 1	1.00	0.98	0.99	126	
						Accuracy		0.98	0.99	126 250	
						Accuracy Macro avg	0.99	0.98	0.99 0.99 0.99	<ul><li>126</li><li>250</li><li>250</li></ul>	
4.2. DCA						Accuracy	0.99	0.98	0.99 0.99 0.99	126 250	
						Accuracy Macro avg	0.99	0.98	0.99 0.99 0.99	<ul><li>126</li><li>250</li><li>250</li></ul>	
PCA Performan		tion.KernelPC	with gamma=0.0	1. we perform di	mensionality	Accuracy  Macro avg  Weighted avg	0.99	0.98 0.99 0.99	0.99 0.99 0.99	126 250 250 250	ne dataset becomes nearly linearly separable. However, the classification accuracy does not improve compa
<ul><li>PCA Performan</li><li>Using sk</li></ul>						Accuracy  Macro avg  Weighted avg	0.99	0.98 0.99 0.99	0.99 0.99 0.99	126 250 250 250	ne dataset becomes nearly linearly separable. However, the classification accuracy does not improve compa
<ul><li>Using sk</li></ul>	klearn.decomposit	ghtly decreases,		e SVM classificat	ion report.	Accuracy  Macro avg  Weighted avg	0.99	0.98 0.99 0.99	0.99 0.99 0.99	126 250 250 250	he dataset becomes nearly linearly separable. However, the classification accuracy does not improve compa

0.15 0.10

0.89

0.83

0.86

0.86

0.86

0.86

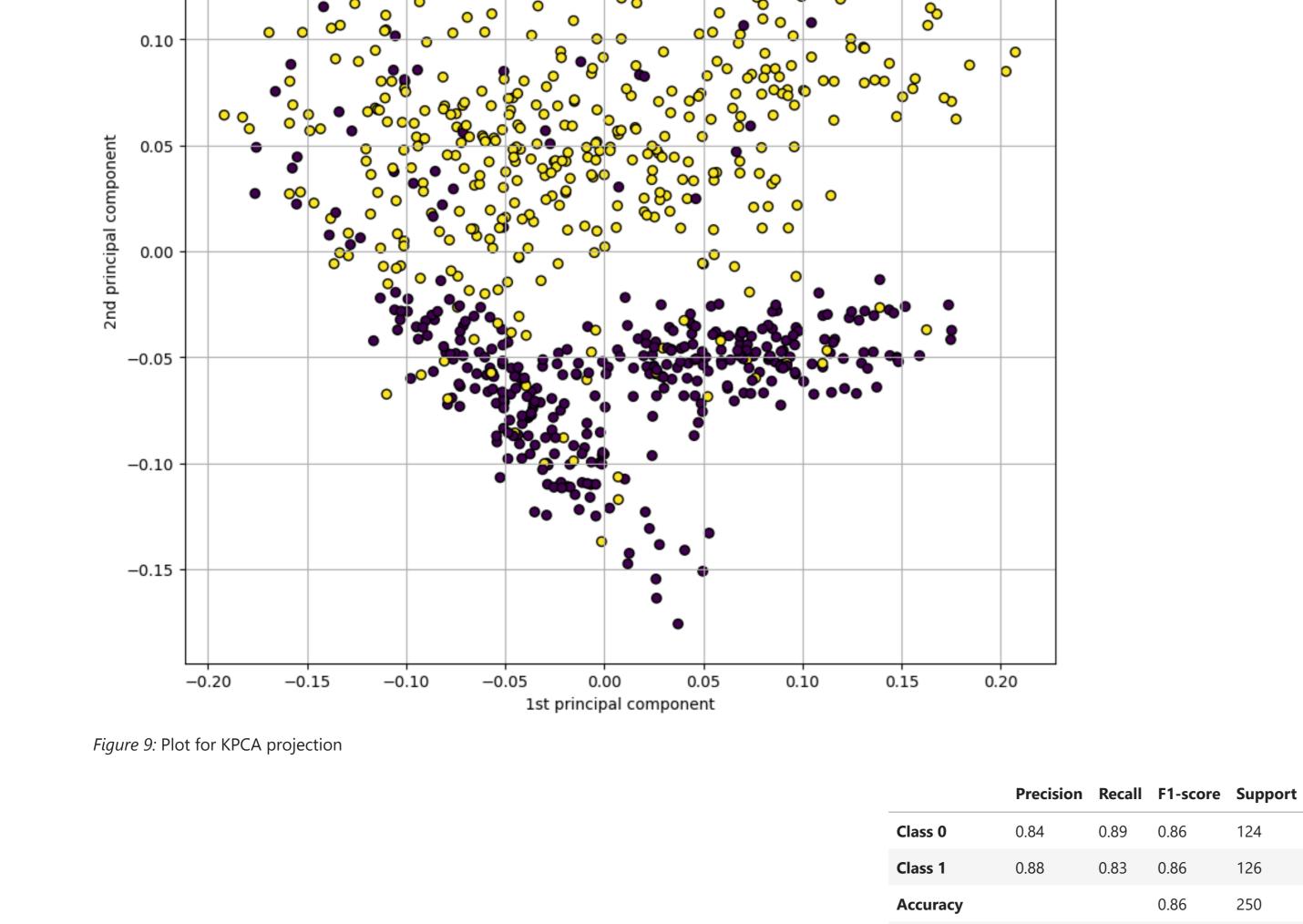
0.86

124

126

250

250



	Weighted avg	0.86	0.86	0.86	250
ngs					
ghlight that while KPCA can be highly effective, its benefits are conditional on the data characteris	stics.				

0

These results high In datasets where the non-linear structure is weak or obscured by noise, traditional PCA's stability and simplicity can make it a more reliable choice.

0.86 Macro avg 4.4 Key Finding