

1) Obtaining and Preparing the Image Data Sets

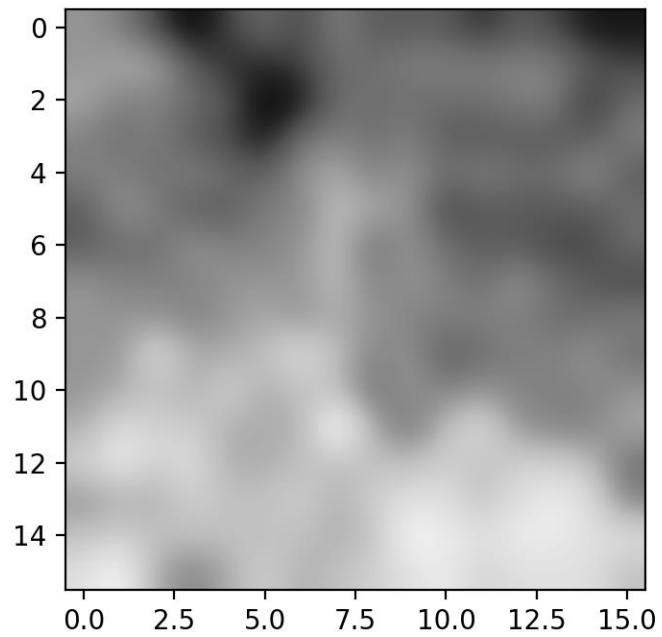


Fig 1. `X_train[3,:]`

2) Assessing Classifiers with Cross Validation

Nearest Centroid Classifier		
Pipeline construction	Parameters	CV score
1. PolynomialFeatures 2. StandardScaler 3. NearestCentroid	Poly_degree = 1	0.7111

K Neighbors Classifier		
Pipeline construction	Parameters	CV score
1. PolynomialFeatures 2. StandardScaler 3. NearestCentroid	Poly_degree = 1 N_neighbors = 1	0.9537

SVC		
Pipeline construction	Parameters	CV score
1. StandardScaler 2. NearestCentroid	Kernel = "rbf" Gamma = 0.16	0.9722

QDA		
Pipeline construction	Parameters	CV score
1. PolynomialFeatures 2. StandardScaler 3. NearestCentroid	Poly_degree = 1 Reg_param= 0.0	0.9593

LDA		
Pipeline construction	Parameters	CV score
1. PolynomialFeatures 2. StandardScaler 3. NearestCentroid	Poly_degree = 2	0.9394

Logistic Regression		
Pipeline construction	Parameters	CV score
1. PolynomialFeatures 2. StandardScaler 3. NearestCentroid	Poly_degree = 2 C = 1.0 solver='liblinear'	0.959

Decision Tree Classifier		
Pipeline construction	Parameters	CV score
1. StandardScaler 2. NearestCentroid	Criterion = "entropy" Splitter = "best" Ccp_alpha = 0.0	0.946

Random Forest Classifier		
Pipeline construction	Parameters	CV score
1. StandardScaler 2. NearestCentroid	N_estimators = 150 Max_features = 4 Max_depth = 60 Criterion = "entropy" Bootstrap = False	0.9832

Note: I used 8x8 images for all classifiers. I switched to 16x16 only for the Random Forest classifier since it looked like the most promising one

3) Choosing an Optimal Classifier

The best classifier is the Random Forest Classifier. The 5-fold cross-validated accuracy on the full training set is 0.9832

4) Optimal Test score

The test score using the best random forest classifier pipeline is 0.9853. The test score is similar to the cross validation mean score of 0.9832 (in fact the test score is slightly higher). Since the classifier appears to generalize well to the test set, I think it is reliable for detecting cracks. The confusion matrix tells us that the false positive and false negative rate is almost the same. This may inform the choice to use the classifier for certain applications where the false negative rate is not acceptable.

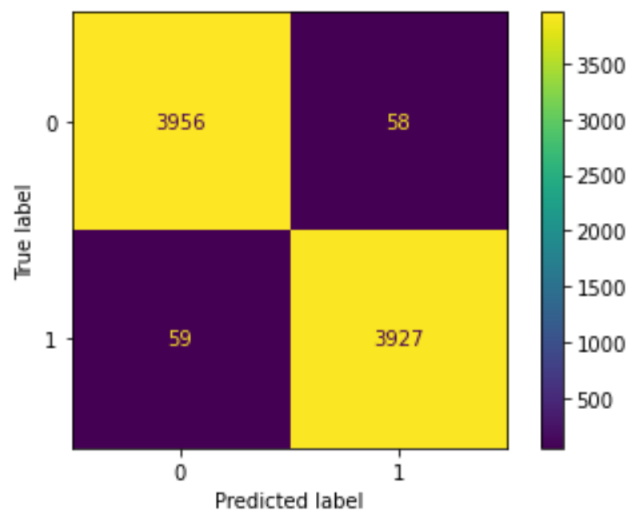


Fig 2. Confusion matrix for the test set