# Fruit Recognition

Eskil Jarslkog: ejarlsko@ucsd.edu, Richard Wang: rsw014@ucsd.edu Joel Andersson: janderss@ucsd.edu

#### Predicting

Our goal for this project was to get acquainted with practical, hands-on machine learning algorithms. We utimately chose to classify images of fruit with the intention to achieve over 90% accuracy using only Support Vector Machines and Random Forests as classifiers.

Our final result was an accuracy of 97 % using PCA+SVM and 94 % using RF + our own features.

#### Data

The data set we used was retrieved from Kaggle, provided by Horea Muresan and Mihai Oltean from Babes-Bolyai University.

#### Specifications:

- 60 different fruits (labels)
- 100x100 pixel images
- Removed background







#### **Features**

We used two sets of features.

Set 1: Each image was flattened and then performed PCA on to extract a fix number of principal components.





PCA4



PCA 1



- Same fruit multiple times but rotated
- ~30000 training set
- ~10000 validation set

## **Models**

#### **Support Vector Machine**

The SVM we used was a Soft-Margin, linear kernel, multiclass SVM, as we had 60 classes and no guarantee of the classes being linearly separable in any one of our feature spaces.

#### Random Forest

**Our Random Forest classifier** was your prototypical RF; a majority vote over a preordained number of binary classification trees.

minimize	$\mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \xi_i^k$	
subject to	$y_i(\mathbf{x}_i^T\mathbf{w}+b) \ge 1-\xi_i,$	
and	$\xi_i \ge 0;  (i=1,\cdots,m)$	

Minimization problem for soft-margin, binary SVM.



An example of a binary classification tree

Set 2 : We extracted nine features from each image by image processing. Six features derived from RGB mean and variance, three features derived from the shape of the fruit.

### <u>Results</u>

	Training Accuracy	Testing Accuracy
SVM on Set 1	100.00%	96.57%
SVM on Set 2	94.15%	91.40%
RF on Set 1	100.00%	90.21%
RF on Set 2	100.00%	93.68%



Confusion matrix for RF on Set 2

Confusion matrix for SVM on Set 1



Overall, we are content with our results. It is favorably compared with a paper on a variant of this dataset, them achieving 94% with Deep Learning. Our results imply that using PCA with the right parameters almost achieves linear separability.

For the case of comparison between the feature sets, examining the most frequent missclassified classes on each feature set (see confusion matrices to the left) can provide a deeper understanding of what information the two different feature sets encapsulates.

Given more time we would also try more data preprocessing such as FFT to extract features based on texture of the fruit, which could potentially help mitigate the most common missclassification shared by both sets.







Apple

Pomegranate

