

## **MACHINE LEARNING: CLASSIFICATION**

PREMANANDA INDIC, PH.D.

DEPARTMENT OF ELECTRICAL ENGINEERING

---

**ORS Research Design & Data Analysis Lab**

Office of Research and Scholarship

# ANALYSIS PLATFORM

---



University of Texas at Tyler

[Get Software](#) | [Learn MATLAB](#) | [Teach with MATLAB](#) | [What's New](#)

**MATLAB Access for Everyone at**

**University of Texas at Tyler**

<https://www.mathworks.com/academia/tah-portal/university-of-texas-at-tyler-1108545.html>

# ANALYSIS PLATFORM



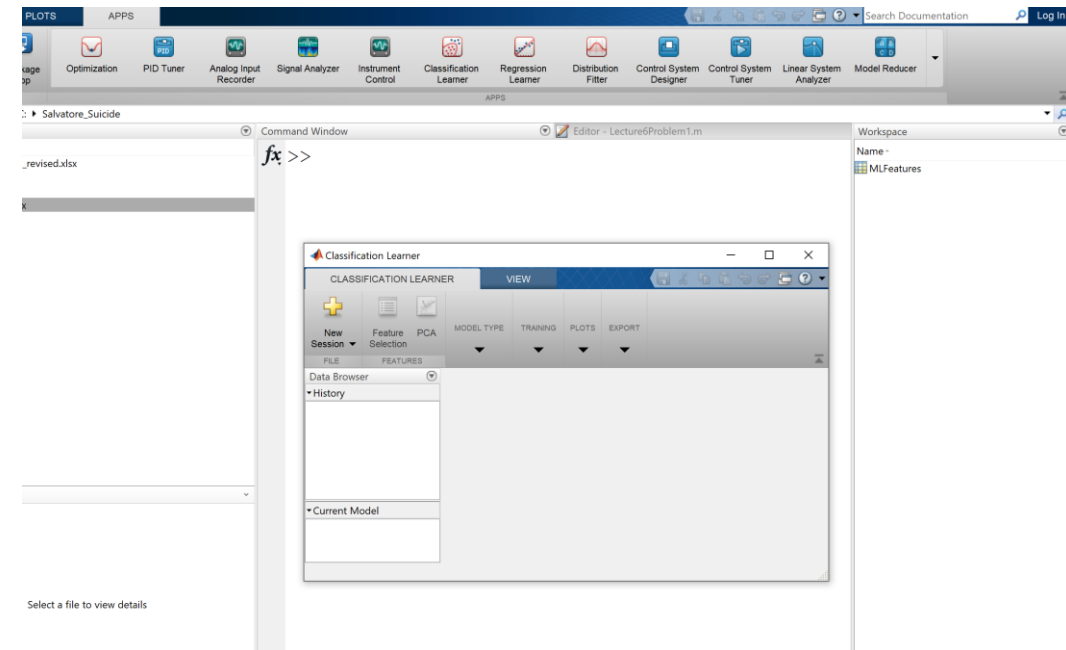
University of Texas at Tyler

[Get Software](#) | [Learn MATLAB](#) | [Teach with MATLAB](#) | [What's New](#)

MATLAB Access for Everyone at

University of Texas at Tyler

<https://www.mathworks.com/academia/tah-portal/university-of-texas-at-tyler-1108545.html>



# OUTLINE

---

➤ INTRODUCTION

➤ DIFFERENT CLASSIFIERS

➤ EXAMPLES

# OUTLINE

---

➤ INTRODUCTION

➤ DIFFERENT CLASSIFIERS

➤ EXAMPLES

# INTRODUCTION

---

## ➤ What is Machine Learning ?

- Machine Learning is a field of study that gives computers the ability to “learn” without being explicitly programmed
  - Prediction
  - Classification

# INTRODUCTION

---

## ➤ What is Machine Learning ?

- Machine Learning is a field of study that gives computers the ability to “learn” without being explicitly programmed
  - Prediction
  - **Classification**

# OUTLINE

---

➤ INTRODUCTION

➤ DIFFERENT CLASSIFIERS

➤ EXAMPLES



# APPROACHES

---

➤ SUPERVISED LEARNING

➤ UNSUPERVISED LEARNING

# APPROACHES

---

## ➤ SUPERVISED LEARNING (Classification / Prediction)

Provide training set with features and solutions

# APPROACHES

---

➤ STANDARD MACHINE LEARNING

➤ ADVANCED MACHINE LEARNING

Based on Artificial Neural Networks (Deep Learning)

# APPROACHES

---

## ➤ CLASSIFICATION

- Logistic Regression
- Support Vector Machine

# APPROACHES

---

## ➤ CLASSIFICATION

- **Logistic Regression**
- Support Vector Machine

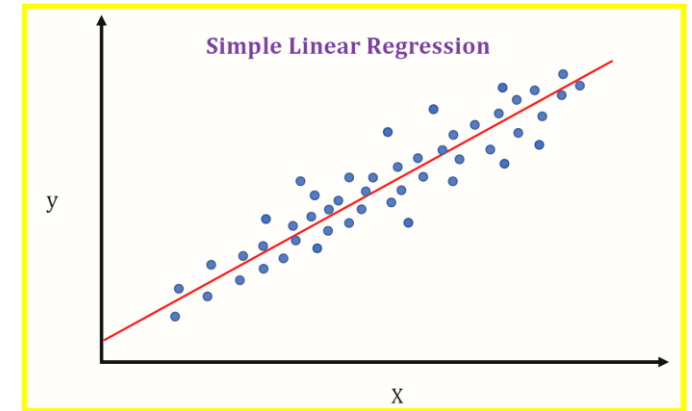
# APPROACHES

---

## ➤ Linear Regression

$$\hat{y}^i = \theta_0 + \theta_1 x_1^i + \theta_2 x_2^i + \dots + \theta_n x_n^i \quad i = 1, 2, \dots, m$$

$$\hat{Y} = \theta^T X$$



<https://medium.datadriveninvestor.com/machine-learning-101-part-1-24835333d38a>

- Gradient Descent by **Louis Augustin Cauchy** in 1847

Cost Function to Minimize

$$J = \left\langle (\hat{y}^i - y^i)^2 \right\rangle = (\hat{Y} - Y)^T (\hat{Y} - Y) = \frac{1}{m} \sum_{i=1}^m (\theta^T X^i - y^i)^2$$

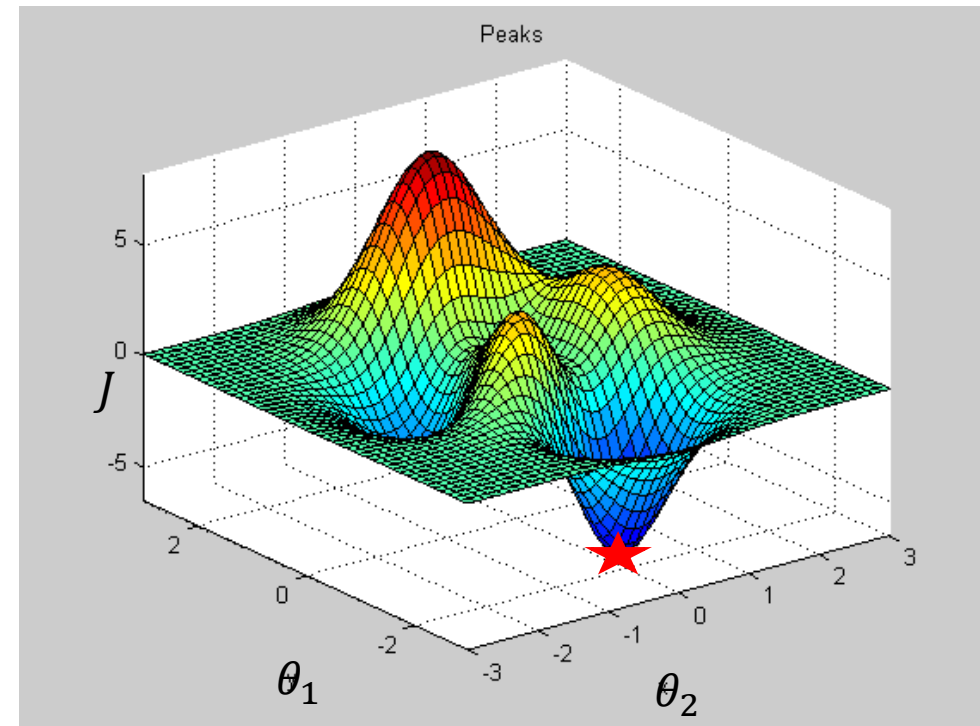
# APPROACHES

---

## ➤ Linear Regression

$$\theta^{k+1} = \theta^k - \gamma \nabla_{\theta} J(\theta)$$

$$\nabla_{\theta} J(\theta) = \frac{2}{m} X^T (X\theta - Y)$$



# APPROACHES

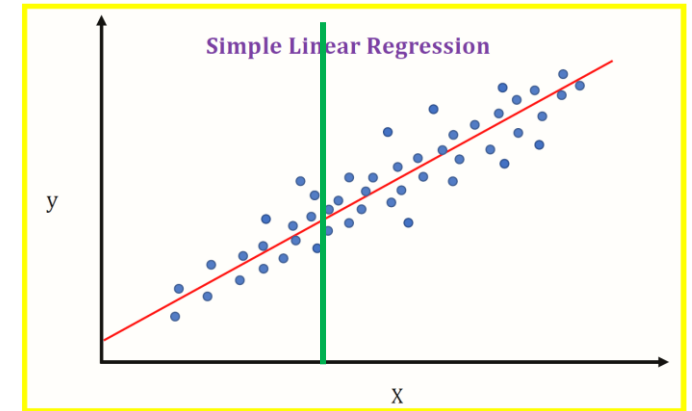
## ➤ Logistic Regression

Two class  $y = 1$  or  $y = 0$

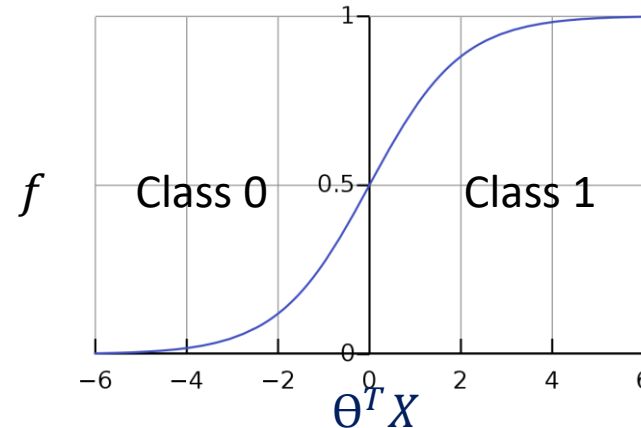
$$\hat{p} = f(\theta^T X) = \frac{1}{1 + e^{-\theta^T X}}$$

$\hat{y} = 1$  if  $\hat{p} < 0.5$ ;  $\hat{y} = 0$  if  $\hat{p} \geq 0.5$

$$J = \frac{1}{m} \sum_{i=1}^m [y^i \log(\hat{p}^i) + (1 - y^i) \log(1 - \hat{p}^i)]$$



<https://medium.datadriveninvestor.com/machine-learning-101-part-1-24835333d38a>





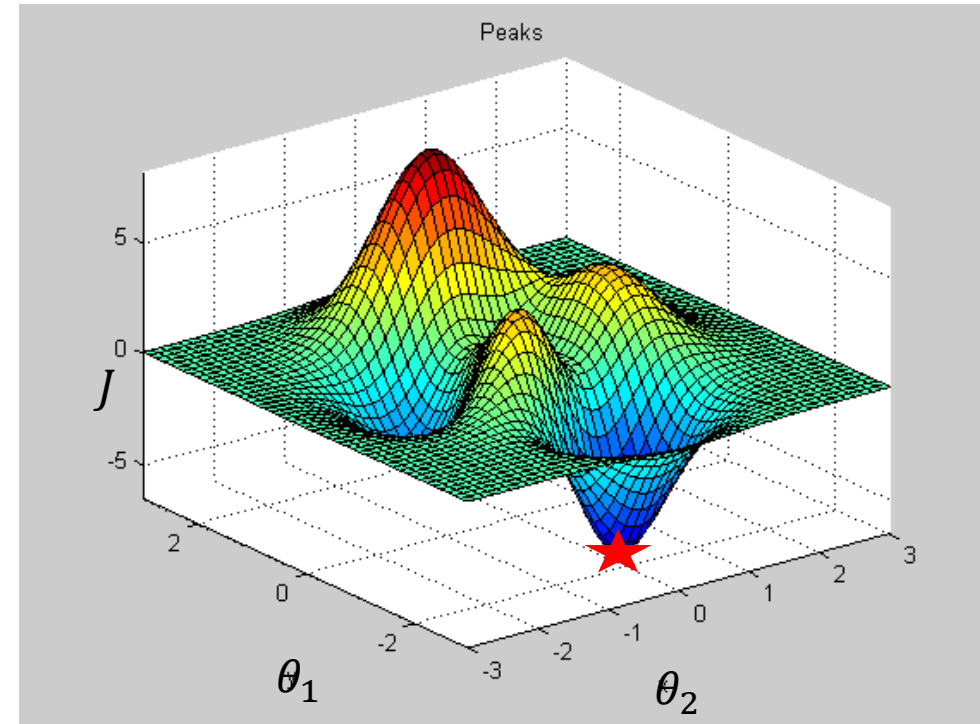
# APPROACHES

---

## ➤ Logistic Regression

$$\theta^{k+1} = \theta^k - \gamma \nabla_{\theta} J(\theta)$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (f(\theta^T X^i) - y^i) x_j^i$$



# APPROACHES

## ➤ Support Vector Machine

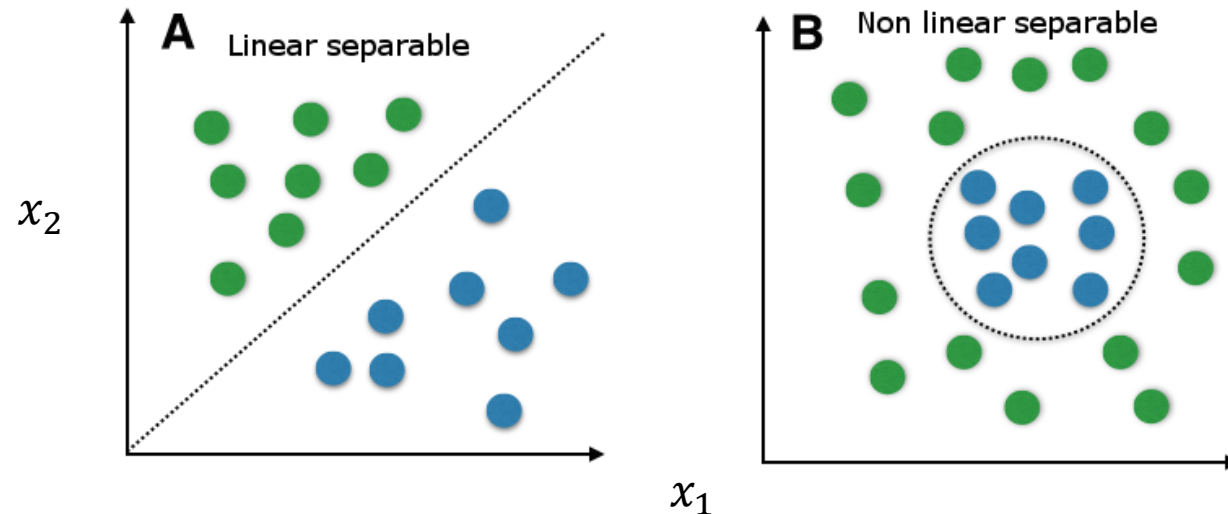
---

$$G(x_j, x_k) = \exp(-\|x_j - x_k\|^2)$$

---

$$G(x_j, x_k) = (1 + x_j'x_k)^q, \text{ where } q \text{ is in the set } \{2,3,\dots\}.$$

$$f(X) = w^T X - b$$

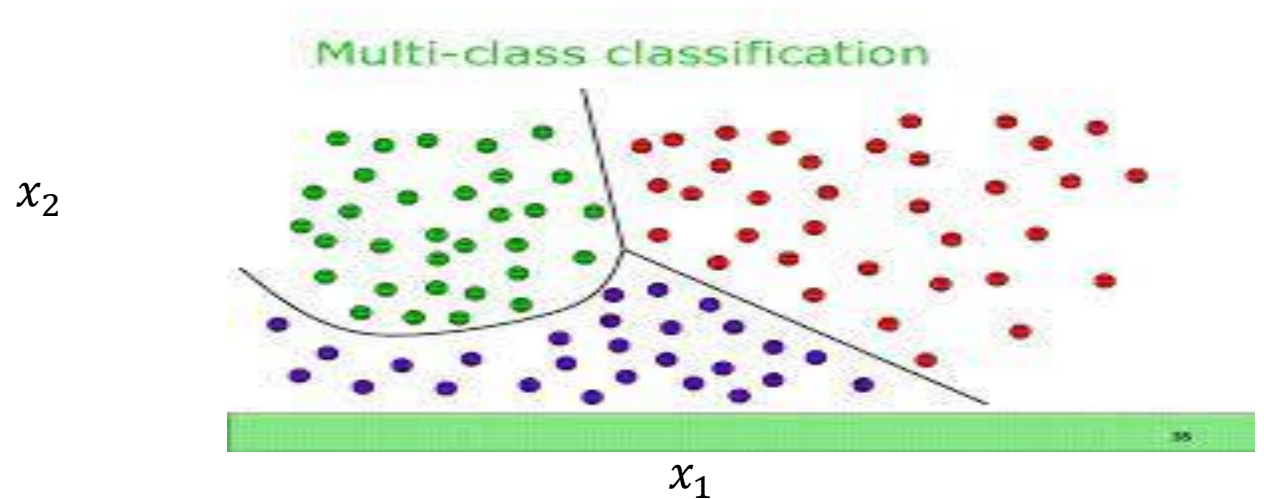


# APPROACHES

---

- SUPERVISED LEARNING (Classification / Prediction)
  - Support Vector Machine (SVM)

Used for regression as well as **classification**



# APPROACHES

---

## ➤ SUPERVISED LEARNING (Classification )

- Logistic Regression
- Support Vector Machines
- k-Nearest Neighbors
- Decision Trees and Random Forests

# SECTION 1: Learner App

---

## ➤ Home Value Classification: 9 features to classify high vs low medianHouseValue

longitude: A measure of how far west a house is; a higher value is farther west

latitude: A measure of how far north a house is; a higher value is farther north

housingMedianAge: Median age of a house within a block; a lower number is a newer building

totalRooms: Total number of rooms within a block

totalBedrooms: Total number of bedrooms within a block

population: Total number of people residing within a block

households: Total number of households, a group of people residing within a home unit, for a block

medianIncome: Median income for households within a block of houses (measured in tens of thousands of US Dollars)

**medianHouseValue: Median house value for households within a block (measured in US Dollars)**

oceanProximity: Location of the house w.r.t ocean/sea

<https://www.kaggle.com/camnugent/california-housing-prices>

Demo with N=5000

70% Training Data

30% Test Data

Models Trained:

Logistic Regression

SVM

# SECTION 1: Learner App

---

## ➤ Prediction of House Price Classification Problem

**Confusion Matrix**

<b>True Class</b>	1	True Positive	False Negative	➔ Total Positive
	0	False Positive	True Negative	
		1	0	
		<b>Predicted Class</b>		

$\text{True Positive Rate} = \text{True Positive} / \text{Total Positive}$

$\text{True Negative Rate} = \text{True Negative} / \text{Total Negative} = 1 - \text{False Positive Rate}$

# SECTION 1: Learner App

## DATA IMPORT & CLASSIFICATION LEARNER INITIALIZATION

New Session from Arguments

**Data set**

Data Set Variable: Ttrain (3500x11 table)

Response: hi\_lo\_label (double, 0 .. 1)

**Predictors**

Name	Type	Range
<input checked="" type="checkbox"/> longitude	double	-124.35 .. -114.56
<input checked="" type="checkbox"/> latitude	double	32.57 .. 41.92
<input checked="" type="checkbox"/> housing_median_age	double	2 .. 52
<input checked="" type="checkbox"/> total_rooms	double	25 .. 39320
<input checked="" type="checkbox"/> total_bedrooms	double	3 .. 6210
<input checked="" type="checkbox"/> population	double	13 .. 16305
<input checked="" type="checkbox"/> households	double	5 .. 5258

Add All Remove All

[How to prepare data](#)

**Validation**

Cross-Validation  
Protects against overfitting by partitioning the data set into folds and estimating accuracy on each fold.

Cross-validation folds: 5

Holdout Validation  
Recommended for large data sets.

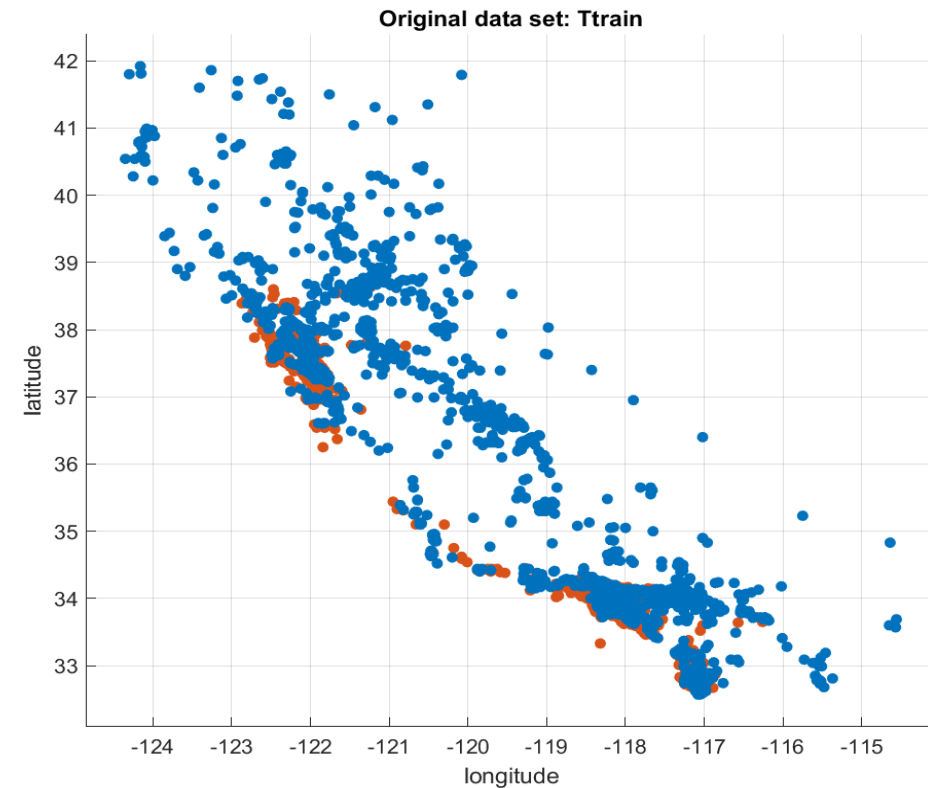
Percent held out: 25

Resubstitution Validation  
No protection against overfitting. The app uses all the data for both training and validation.

[Read about validation](#)

Start Session Cancel

⚠ Response variable is numeric. Distinct values will be interpreted as class labels.



# SECTION 1: Learner App

---

## ➤ DATA IMPORT & CLASSIFICATION LEARNER INITIALIZATION

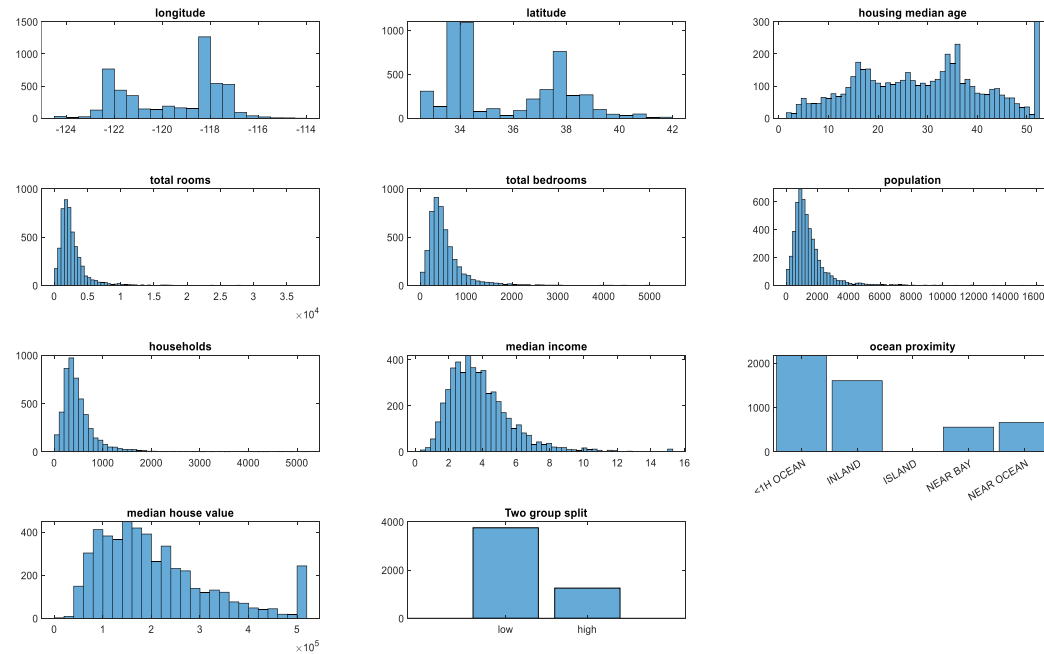
```
classificationLearner(Ttrain, 'hi_lo_label');
```

Demo Learner App in MATLAB - logistic regression and linear SVM



# SECTION 2: Raw Data Analysis

Visualize the data, Summarize variables, data cleaning, pre-processing if needed



207 Missing values, replace with median values

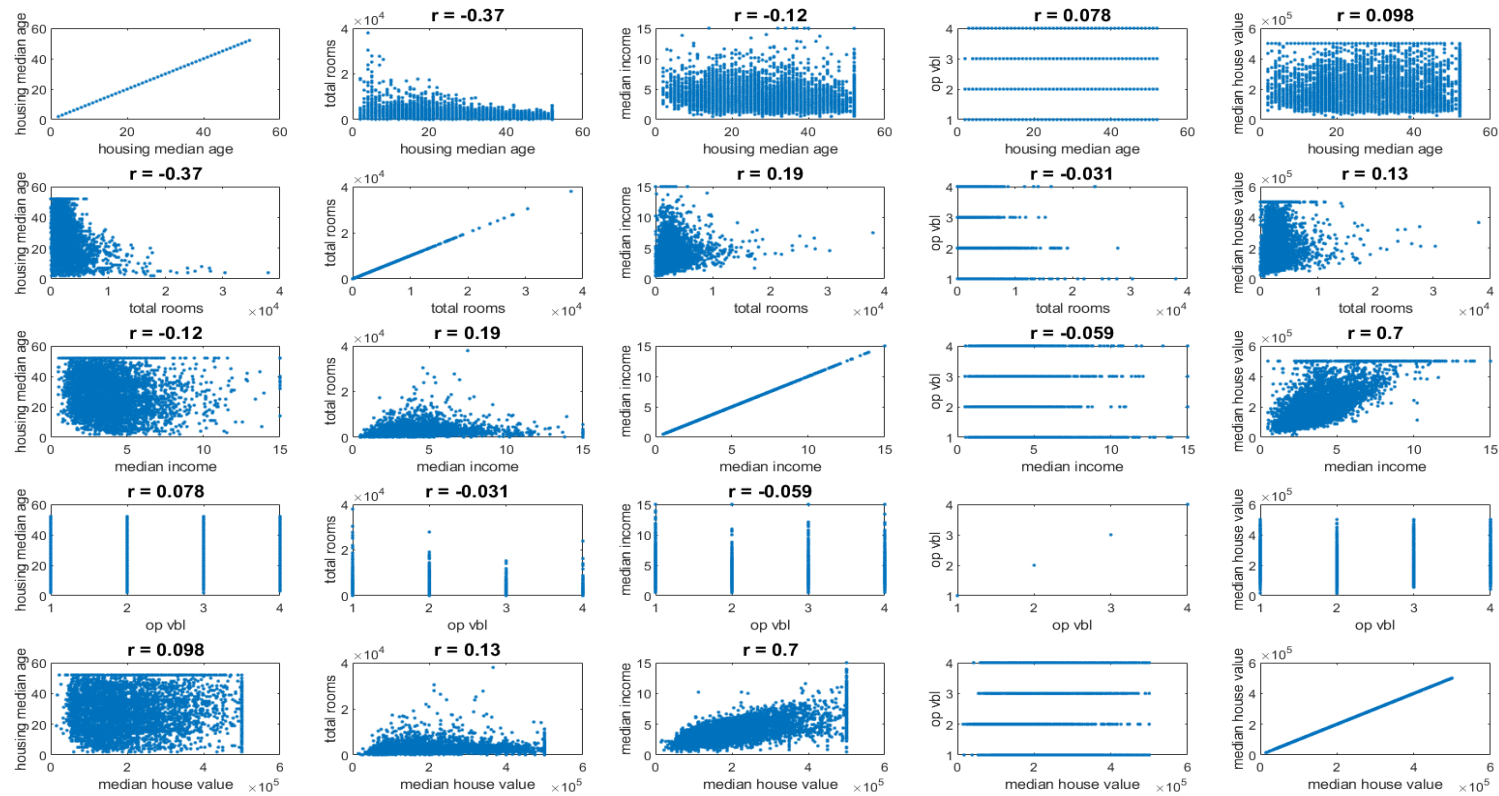
ocean\_proximity: 20636×1 categorical

Values:

<1H OCEAN	9135
INLAND	6550
<b>ISLAND</b>	<b>5</b>
NEAR BAY	2289
NEAR OCEAN	2657

# SECTION 3: Correlation Analysis

FIND VARIABLE CORRELATIONS TO EACH OTHER AND THE MEDIAN HOUSE VALUE



```
[R, pp] = corr(table2array(T1(:, select_vars)));
```

# SECTION 4: Logistic Regression

## SPLIT INTO TRAINING AND TEST DATA AND FIT LOGISTIC REGRESSION MODEL

Estimated Coefficients:

	Estimate	SE	tStat	pValue
(Intercept)	-154.19	14.421	-10.692	1.1065e-26
longitude	-1.7683	0.17448	-10.135	3.8752e-24
latitude	-1.8133	0.18885	-9.6018	7.8546e-22
housing_median_age	0.044239	0.0051484	8.5928	8.4901e-18
total_rooms	0.0003444	9.7387e-05	3.5364	0.00040561
total_bedrooms	0.00080298	0.00084259	0.95299	0.3406
population	-0.0023529	0.00020995	-11.207	3.7737e-29
households	0.0039573	0.00094559	4.185	2.8517e-05
median_income	1.0172	0.053904	18.87	2.0101e-79
ocean_proximity_INLAND	-0.053285	0.24937	-0.21368	0.8308
ocean_proximity_ISLAND	0	0	NaN	NaN
ocean_proximity_NEAR BAY	-0.10616	0.19861	-0.53449	0.593
ocean_proximity_NEAR OCEAN	0.11076	0.15948	0.6945	0.48737

```
mdl = fitglm([Ttrain(:,1:9)  
table(y)], 'Distribution', 'binomial');
```

3500 observations, 3488 error degrees of freedom

Dispersion: 1

Chi^2-statistic vs. constant model: 1.83e+03, p-value = 0

Remove Insignificant features

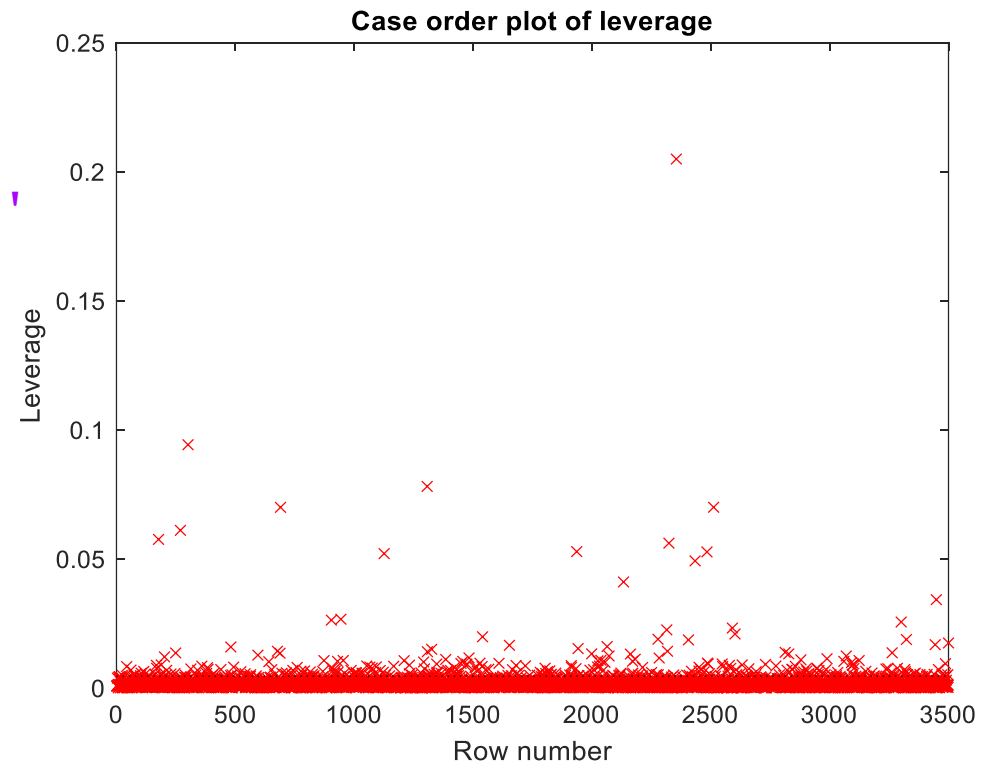
# SECTION 5: Outliers

## DIAGNOSTICS OF MODELS- IDENTIFY OUTLIERS

```
mdl1 = fitglm([Ttrain(:, [1:4 6:8])  
table(y, 'variablenames', {'Hi_lo_label'})], '  
Distribution', 'binomial');
```

```
plotDiagnostics(mdl1, 'leverage')
```

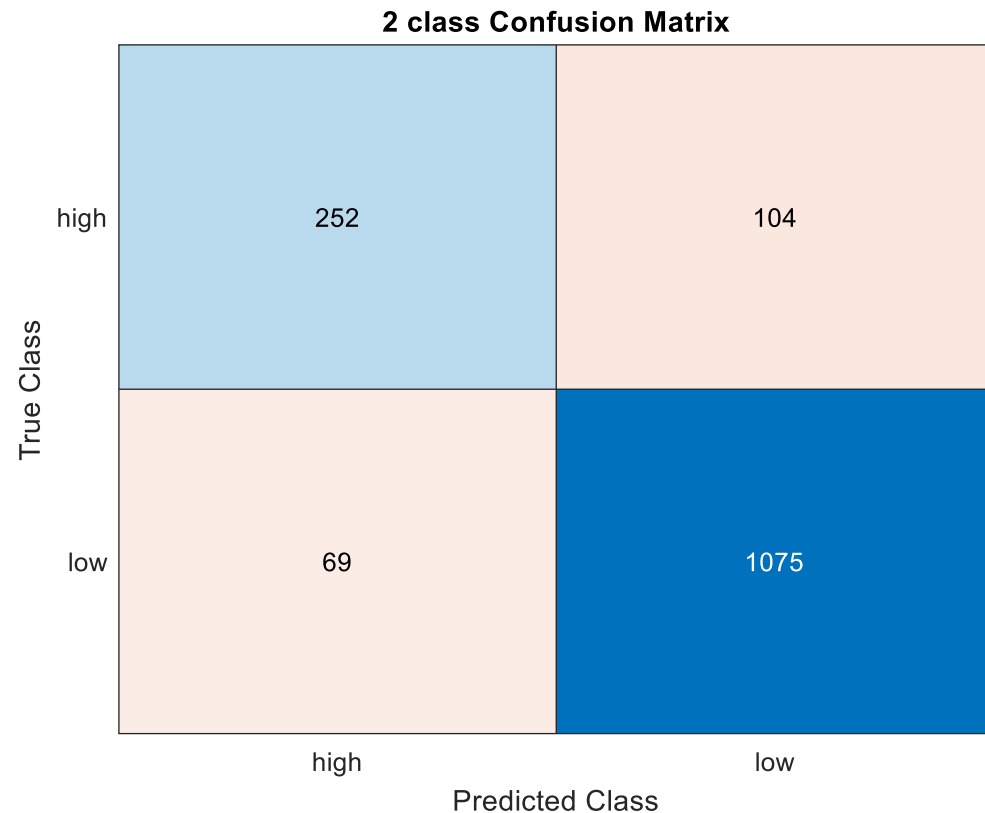
Demo with MATLAB



# SECTION 6: Classification (Clean Data)

---

TEST MODEL FOR TWO CLASS CLASSIFICATION (Logistic Regression)



Test Data N = 1500  
(30% of 5000)

Missing Values  
Insignificant Features  
Outliers

# SECTION 7: SVM Classification

REGULARIZATION OF VARIABLES DONE AUTOMATICALLY, NO NEED TO CHOOSE FEATURES SEPARATELY AS WAS DONE EARLIER FOR LOGISTIC REGRESSION

**SVM - 2 class**

high	255	101
low	69	1075

True Class

71.6%	28.4%
94.0%	6.0%

Test Data N = 1500  
(30% of 5000)

Linear SVM

78.7%	91.4%
21.3%	8.6%

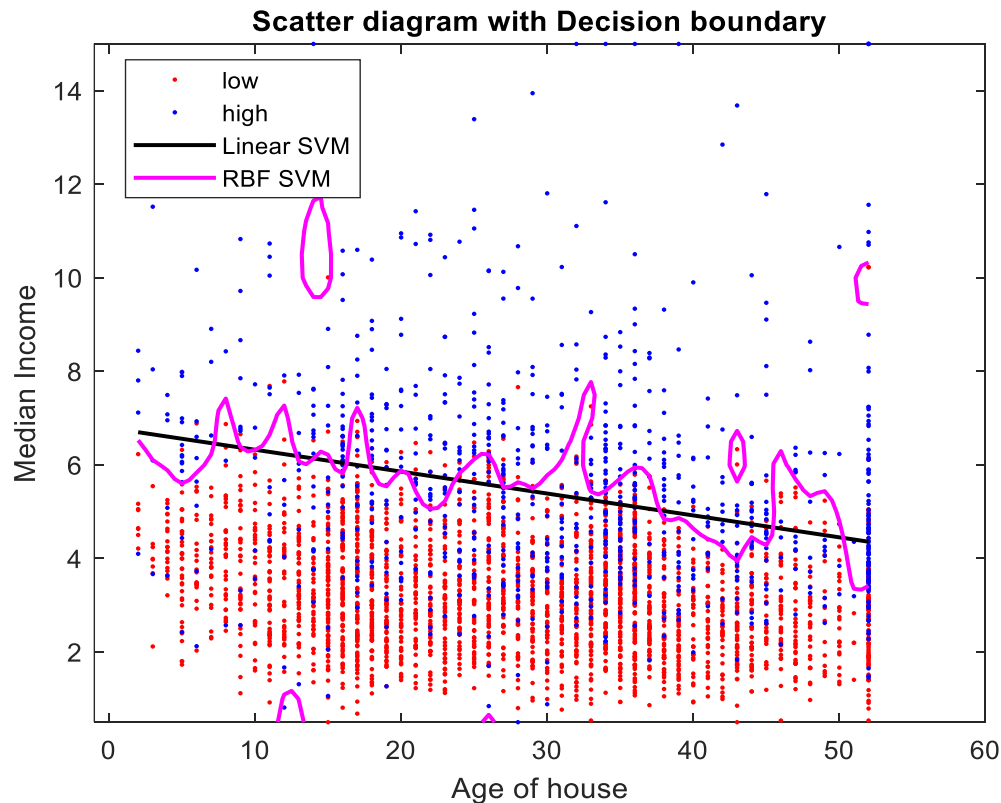
high low  
Predicted Class

```
SVMModel = fitcsvm(Ttrain(:,1:9), y, 'standardize', true);
```

Demo Logistic Regression and SVM binary classification with cleaned up data - PYTHON

# SECTION 8: SVM Classification

## LINEAR vs RADIAL BASIS FUNCTION (RBF) KERNEL



```
fitcsvm([x1 x2],y1);  
fitcsvm([x1 x2],y1,'KernelFunction','rbf');
```

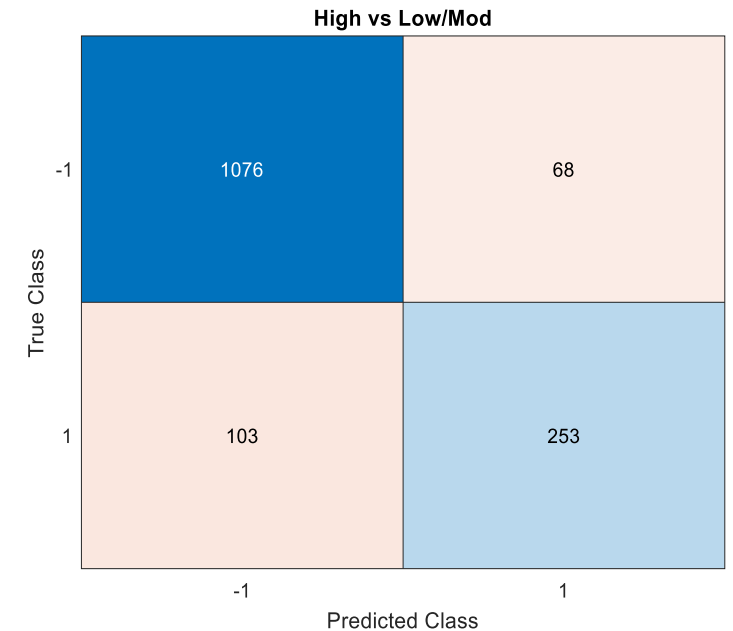
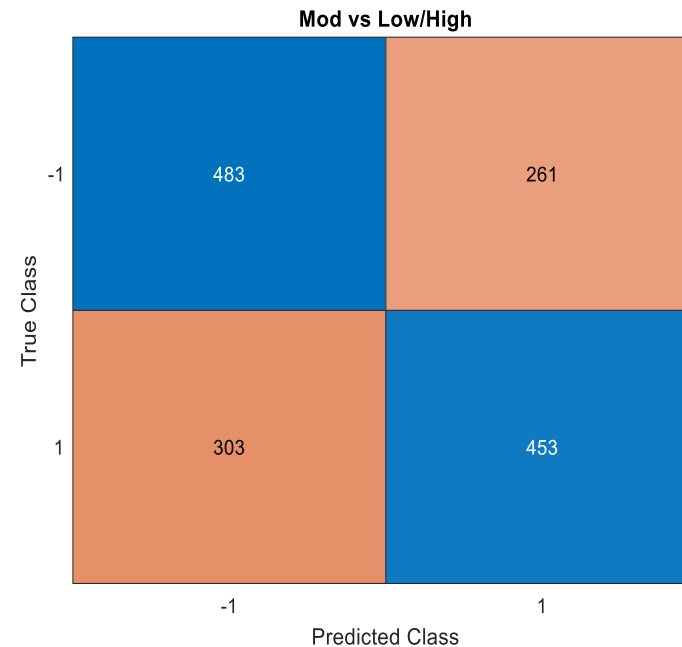
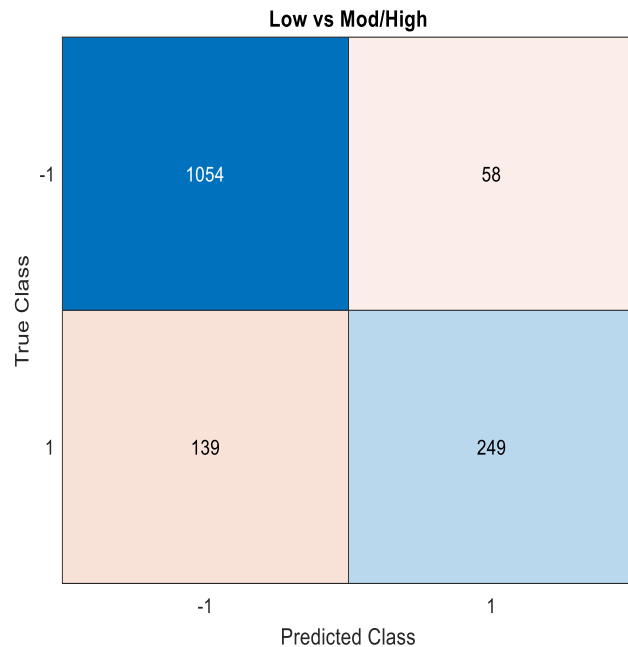
x1: Age of House  
x2: Median Income

Demo SVM decision boundaries  
with MATLAB

# SECTION 9: Multiclassification (SVM)

ONE CLASS vs REST

Also perform one to one class



```
Mdl =  
fitcecoc(Ttrain(:,1:8), y, 'Learners', t, 'Coding', coding, 'ResponseName', responseName, ...  
        'PredictorNames', predictorNames, 'ClassNames', classNames);
```



# SECTION 10: Multiclassification (SVM)

## LOW vs MOD vs HIGH CLASS

```
Mdlp =  
fitcecoc(Ttrain(:,1:8),y,'Learners',t,'FitPosterior',true,...  
'ClassNames',{'low','mod','high'}  
,...  
'Verbose',2);
```

Demo SVM Multi-class  
classification with MATLAB

3 class SVM classification

True Class	Predicted Class				
	high	low	mod		
high	258	2	96	72.5%	27.5%
low	1	257	130	66.2%	33.8%
mod	70	56	630	83.3%	16.7%

# CONCLUSION

---

- Classification divides the data into different groups
- Look at the raw data and understand features in relation to class designation
- Several codes are available to perform classification



SBIR: RAE (Realize, Analyze, Engage) - A digital biomarker based detection and intervention system for stress and cravings during recovery from substance abuse disorders.

*PIs: M. Reinhardt, S. Carreiro, P. Indic*



STARs Award

The University of Texas System  
*P. Indic (PI, UT Tyler)*

**Research Design & Data Analysis Lab**  
Office of Research, Scholarship, and Sponsored Programs



Department of Veterans Affairs

*Design of a wearable sensor system and associated algorithm to track suicidal ideation from movement variability and develop a novel objective marker of suicidal ideation and behavior risk in veterans.*

Clinical Science Research and Development Grant (approved for funding),

*P. Indic (site PI, UT-Tyler)*

*E.G. Smith (Project PI, VA)*

*P. Salvatore (Investigator, Harvard University)*



Pre-Vent

National Institute Of Health Grant

*P. Indic (Analytical Core PI, UT-Tyler)*

*N. Ambal (PI, Univ. of Alabama, Birmingham)*



ViSiON

National Institute Of Health Grant

*P. Indic (Co-Investigator & site PI, UT-Tyler)*

*P. Ramanand (Co-Investigator, UT-Tyler)*

*N. Ambal (PI, Univ. of Alabama, Birmingham)*

# QUESTIONS

---