

# Comparing Performance of Machine Learning Algorithms for Galaxy Classification

F. Korhan YELKENCI<sup>1,3</sup>, E. Kaan ÜLGEN<sup>1</sup>, Sinan ALIŞ<sup>1,2</sup>, Süleyman FİŞEK<sup>1,2</sup>

<sup>1</sup>Istanbul University, Dept. of Astronomy and Space Sciences - Turkey

<sup>2</sup>Istanbul University Observatory Application and Research Center - Turkey

<sup>3</sup>Istanbul University, Dept. of Informatics - Turkey

E-mail: [yelkenci@istanbul.edu.tr](mailto:yelkenci@istanbul.edu.tr)

Web: [cosmology.istanbul.edu.tr](http://cosmology.istanbul.edu.tr)



# Motivation

- Galaxy Data: Deep, wide, big and good quality.
- We don't have enough number of eye to classify all galaxies.
  - Visual classification may have some biases or misclassifications.
  - At high redshift, human eye is ineffective to classify galaxies.
  - We can use parametric features of galaxies to classify (colours, structural parameters, sersic indices, etc.).
  - Try to understand galaxy morphologies.
- Machine learning algorithms are now more efficient and reliable.
  - Comparing performances of ML algorithms on galaxy classifications
  - How accuracy changes by using different parametric features.
  - To see whether accuracy of ML algorithms change with redshift.
- Astronomical data will burst in next years.

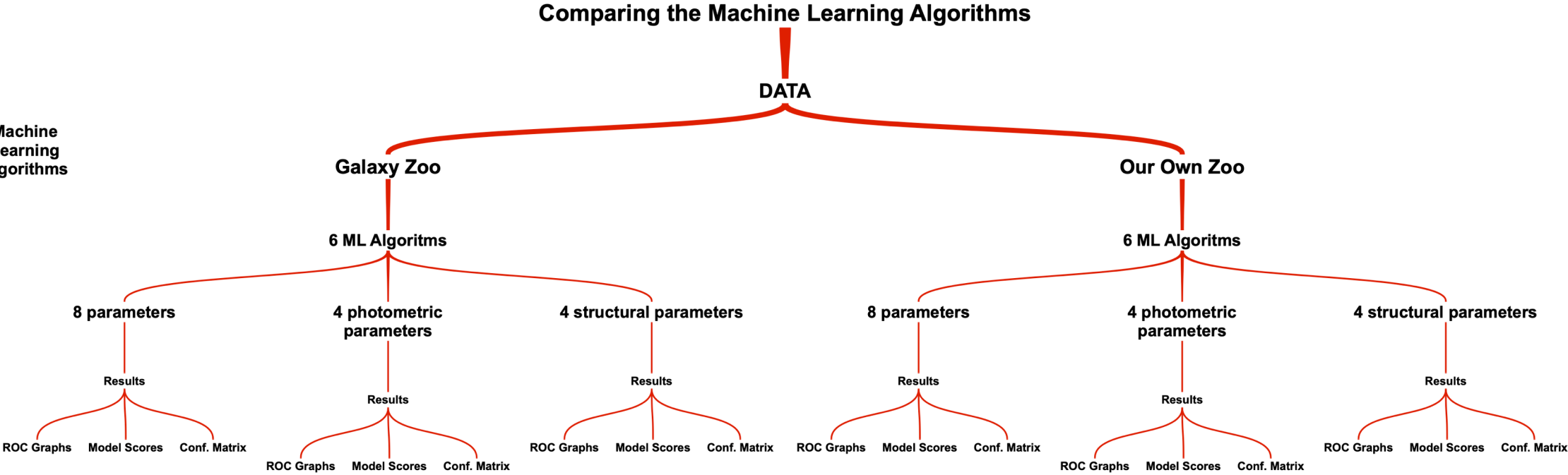
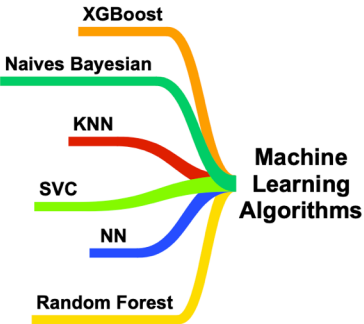
# Background

Statistical learning method	Total sample	Training set	Test set	Number of classes	Dimensions	Accuracy	Reference
SVM	7528	6022 (80 per cent)	1506 (20 per cent)	5	10	75.8 per cent	Results from our work
NN						76.0 per cent	
CT						69.0 per cent	
CTRF						76.2 per cent	
SVM	~1500	500 (33 per cent)	1000 (67 per cent)	2 (early-type, late-type)	12	80 per cent	Huertas-Company et al. (2007)
NN	~1 000 000	~75 000 (7.5 per cent)	~925 000 (92.5 per cent)	3 (early-type, spirals, point sources/artefacts)	12	90 per cent	Banerji et al. (2010)
Oblique CT	5217	~4174 (80 per cent)	~1043 (20 per cent)	5 (E, S0, Sa+Sb, Sc+Sd, Irr)	13	63 per cent	Owens et al. (1996)
Three CT	75 000	67 500 (90 per cent)	7500 (10 per cent)	3 (ellipticals, spirals, unknown)	13	96.2 per cent	Gauci et al. (2010)
algorithms including CTRF							
ConvNet	58 000	47 700 (~82 per cent)	5000 (~9 per cent) 5300 (~9 per cent) used for real-time evaluation during training	5 (probabilities <sup>a</sup> )	Run on images	~99 per cent	Huertas-Company et al. (2015) Dieleman et al. (2015)

Note. <sup>a</sup>Probabilities for each galaxy having a disc or a spheroid, being a point source, having an irregularity or being unclassifiable are the outputs.

Sreejith et al., MNRAS **474**, 5232–5258 (2018)

# Schema of Application



**8 Parametric features:**  $u_i$ ,  $g_r$ ,  $g_i$ ,  $deVAB\_r(b/a)$ ,  $deVRad\_r(\text{eff radius})$ ,  $CI(petroR50\_r/ petroR90\_r)$ ,  $absMagR$ ,  $sersic\_n$   
**4 Structural features:**  $deVAB\_r(b/a)$ ,  $deVRad\_r(\text{eff radius})$ ,  $CI(petroR50\_r/ petroR90\_r)$ ,  $sersic\_n$   
**4 Photometric features:**  $u_i$ ,  $g_r$ ,  $g_i$ ,  $absMagR$  )

# Machine Learning Algorithms

	<u>Unsupervised</u>	<u>Supervised</u>
<u>Continuous</u>	<ul style="list-style-type: none"><li>• Clustering &amp; Dimensionality Reduction<ul style="list-style-type: none"><li>○ SVD</li><li>○ PCA</li><li>○ K-means</li></ul></li></ul>	<ul style="list-style-type: none"><li>• Regression</li><li>• Decision Trees</li><li>• <u>Neural Networks</u></li><li>• <u>Ensemble Methods</u><ul style="list-style-type: none"><li>○ <u>XGBoost</u></li><li>○ <u>Random Forest</u></li></ul></li></ul>
<u>Categorical</u>	<ul style="list-style-type: none"><li>• Association Analysis<ul style="list-style-type: none"><li>○ Apriori</li><li>○ FP-Growth</li></ul></li><li>• Hidden Markov Model</li></ul>	<ul style="list-style-type: none"><li>• Classification<ul style="list-style-type: none"><li>○ KNN</li><li>○ Trees</li><li>○ Logistic Regression</li><li>○ Naive-Bayes</li><li>○ SVM<ul style="list-style-type: none"><li>○ SVC</li></ul></li></ul></li></ul>



## Supervised Learning

- KNN (K-Nearest Neighbour)
- SVM (Support Vector Machine)
  - C-SVC
- Ensemble Methods
  - XGBoost
  - Random Forest
- Gaussian Naive Bayesian
- Neural Network

- SDSS - Galaxy Zoo:

<https://data.galaxyzoo.org>

<https://www.sdss.org>

Lintott et al. 2008, MNRAS, 389, 1179

Lintott et al. 2011, 410, 166



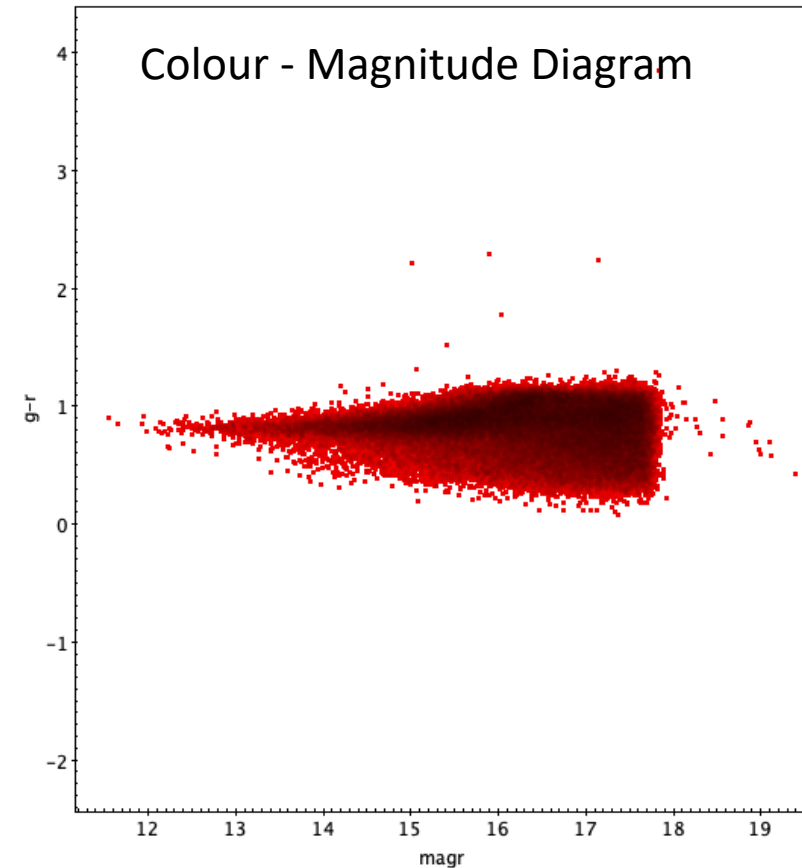
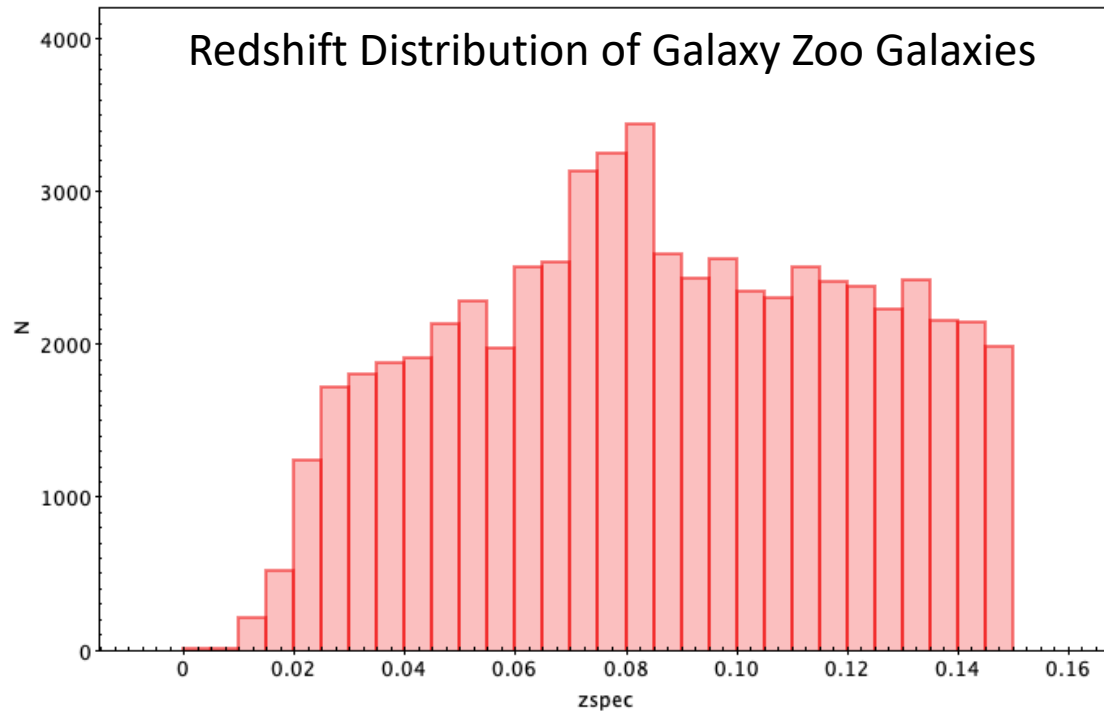
- CFHT-LS

<https://www.cfht.hawaii.edu/Science/CFHTLS/>



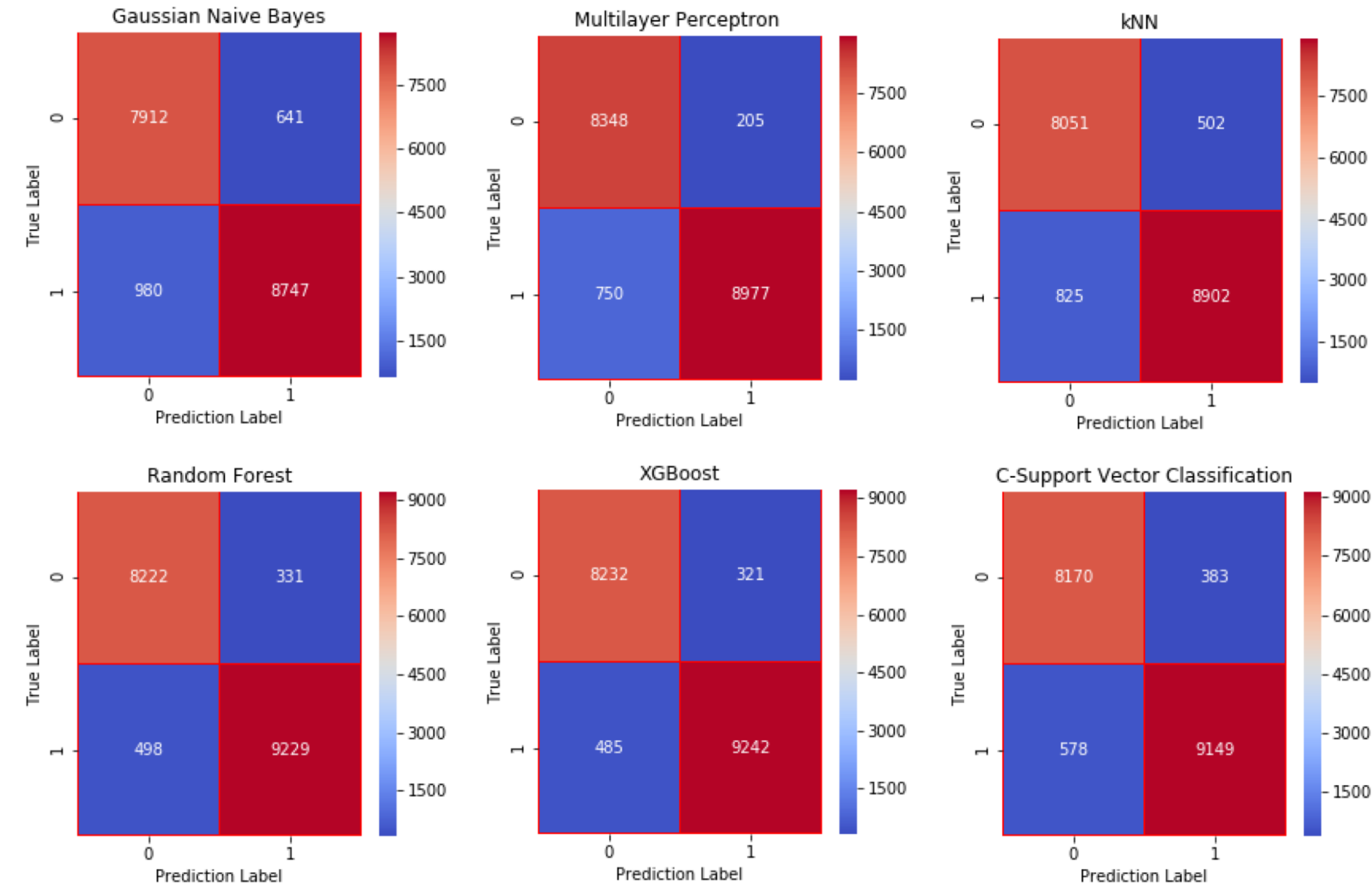
# Galaxy Zoo DATA

- **60932** galaxies → Elliptical: 28591, Spiral: 32341 visually classified in Galaxy Zoo.
- Redshift range:  $0 < z < 0.15$
- with Sersic Index
- Train and Test sets (%70/30 : 42652 / 18280)



# Galaxy Zoo DATA: Results for all 8 parameters

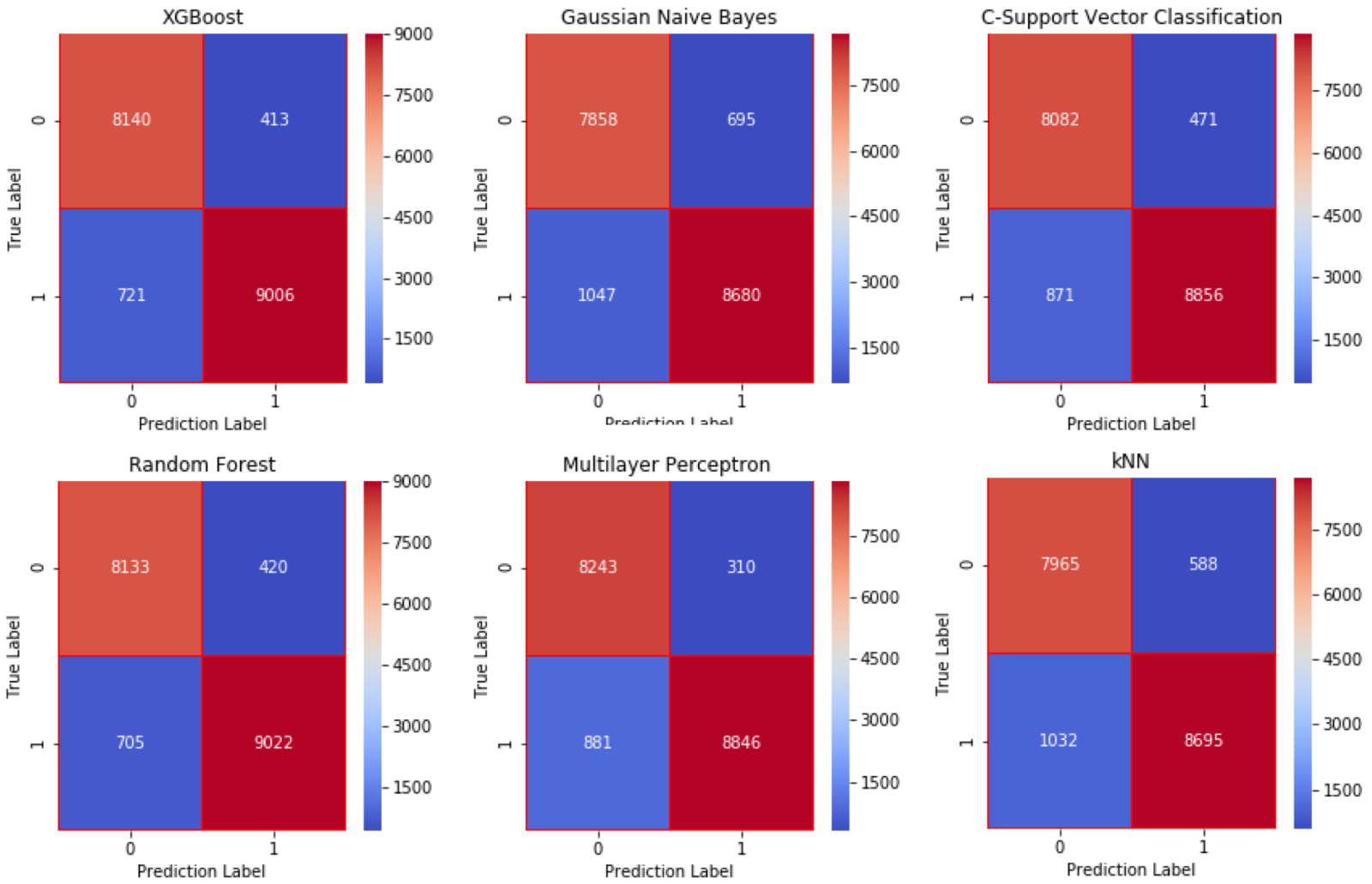
(ui, gr, gi, deVAB\_r(b/a), deVRad\_r(eff radius), CI (petroR50\_r/petroR90\_r), absMagR , sersic\_n)



Model	Score
XGBoost	0.955383
Random Forest	0.954773
NN	0.949405
SVC	0.947247
KNN	0.928444
Naive Bayes	0.912900

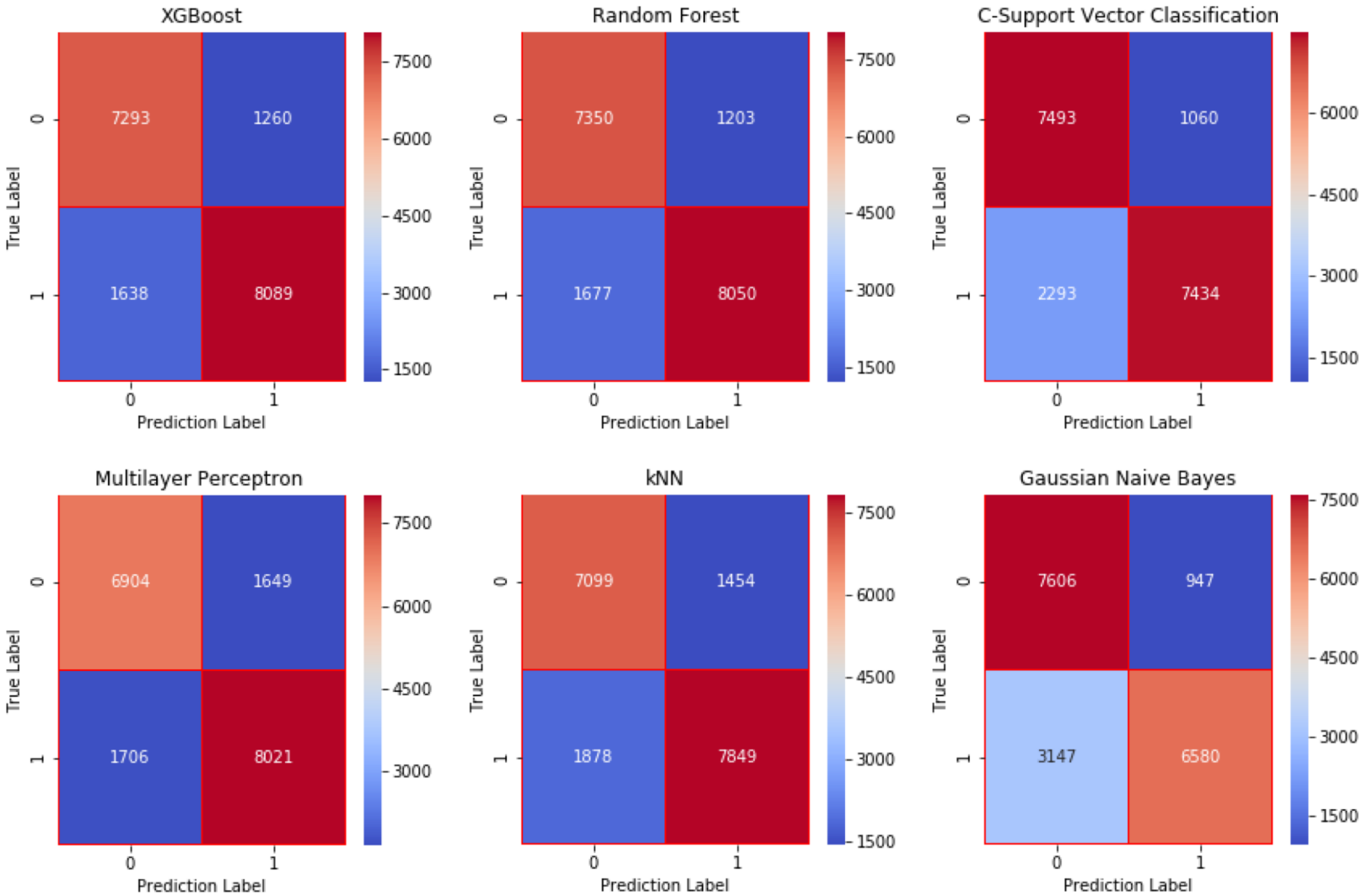


# Galaxy Zoo DATA: Results for 4 structural parameters (deVAB\_r(b/a), deVRad\_r(eff radius), CI (petroR50\_r/petroR90\_r), sersic\_n)



Model	Score
Random Forest	0.938151
XGBoost	0.937682
NN	0.936697
SVC	0.925068
KNN	0.909031
Naive Bayes	0.903521

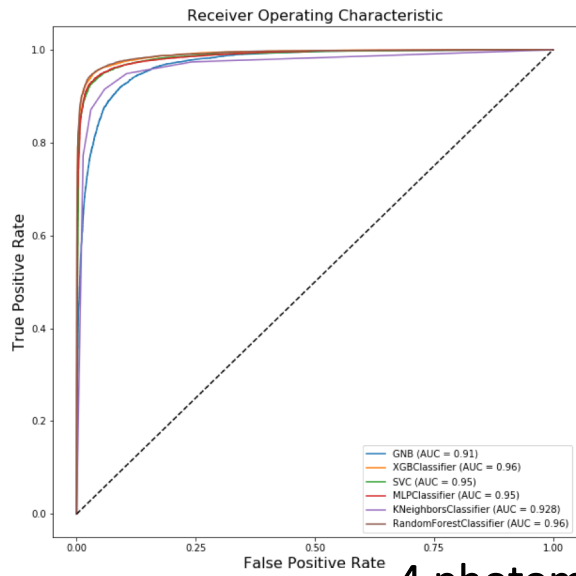
# Galaxy Zoo DATA: Results for 4 photometric parameters (ui, gr, gi, absMagR)



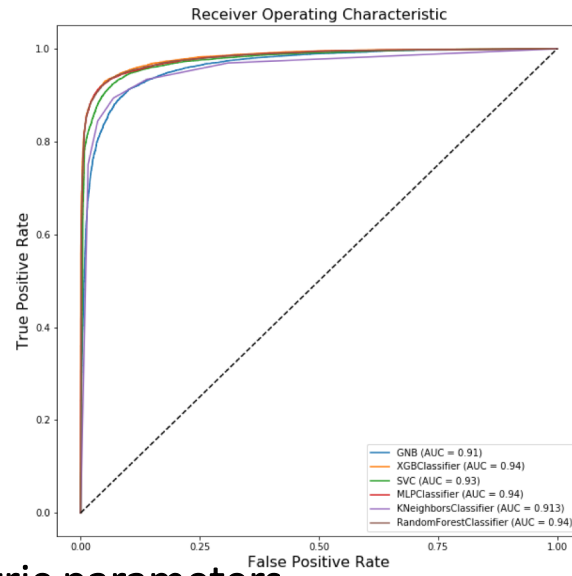
Model	Score
Random Forest	0.840242
XGBoost	0.838484
KNN	0.823057
SVC	0.816000
NN	0.812506
Naive Bayes	0.775908

# Galaxy Zoo DATA: Comparison of ROC graphs of 6 Machine Learning Algorithms

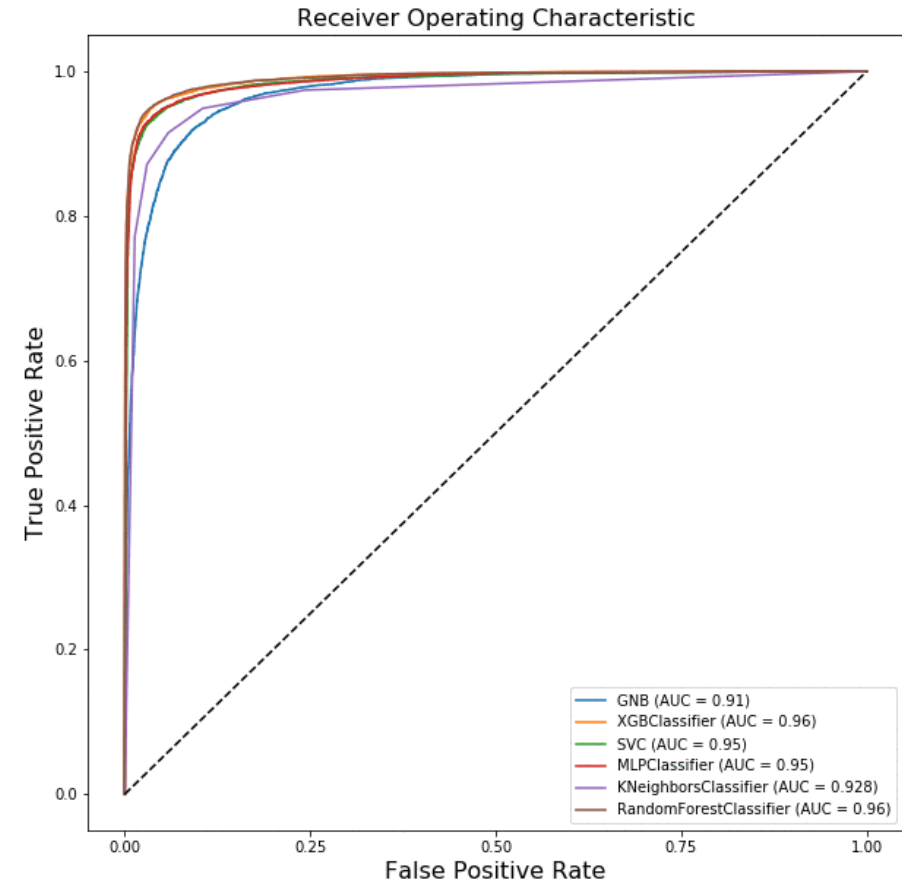
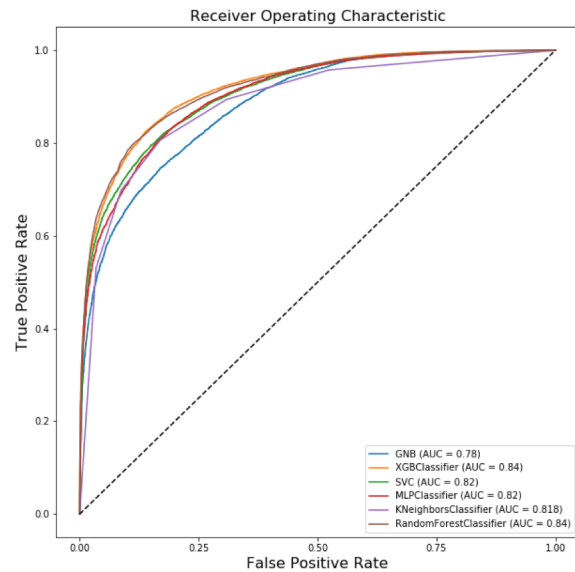
## All 8 parameters



## 4 structural parameters



## 4 photometric parameters



# Impact of 8 Parameters in Random Forest and XGBoost

## Random Forest – Parameters

Feature	Importance
deVAB_r	0.268648
CI	0.241121
sersic_n	0.188087
gi	0.086840
ui	0.083662
gr	0.051916
absMagR	0.039950
deVRad_r	0.039776

## XGBoost - Parameters

Feature	Importance
CI	0.485405
deVAB_r	0.157679
sersic_n	0.106899
ui	0.083637
gi	0.074488
gr	0.038916
deVRad_r	0.034296
absMagR	0.018678

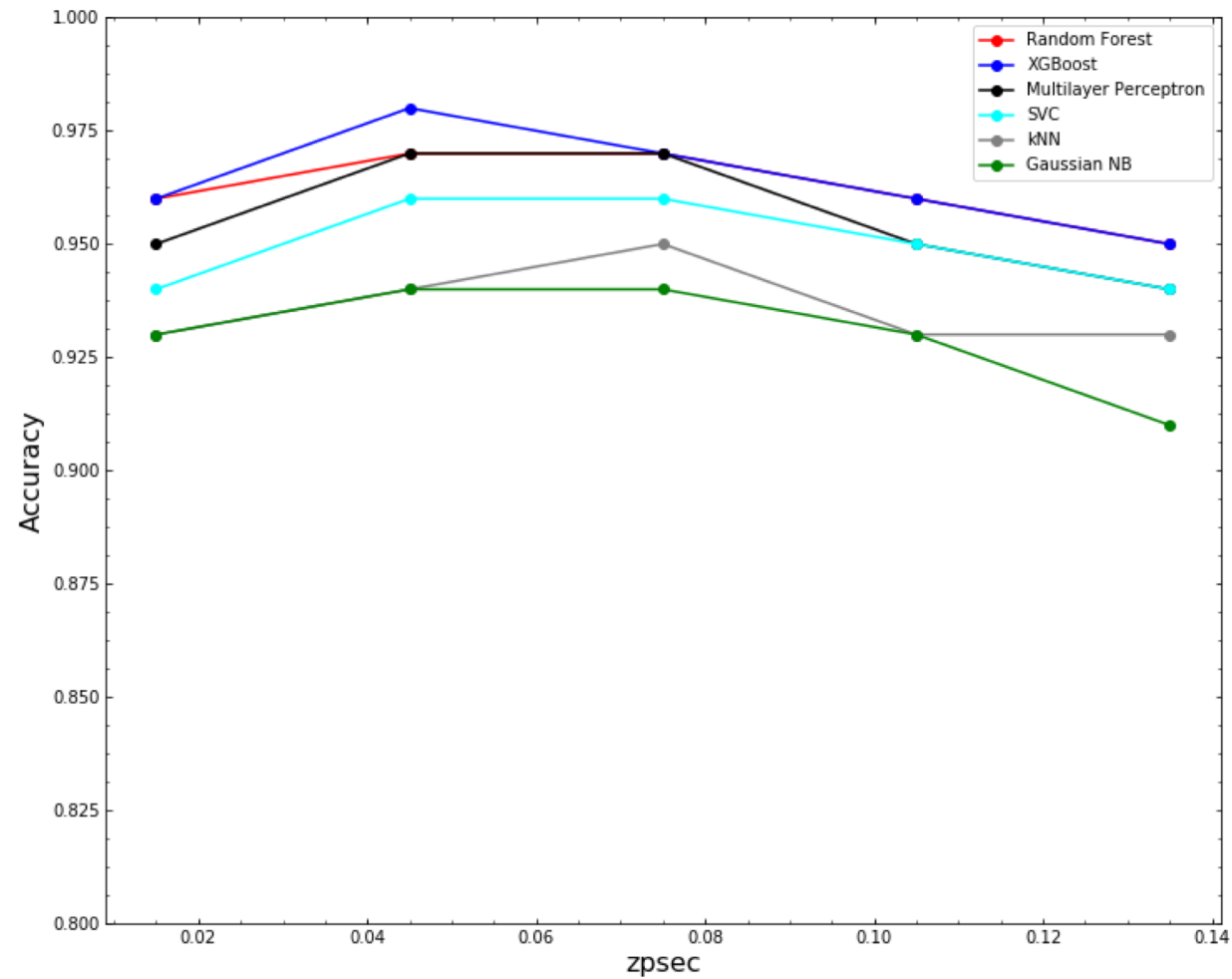
## 4 structural

XGBoost - Parameters	
Feature	Importance
CI	0.551122
deVAB_r	0.219104
sersic_n	0.171079
deVRad_r	0.058695

## 4 photometric

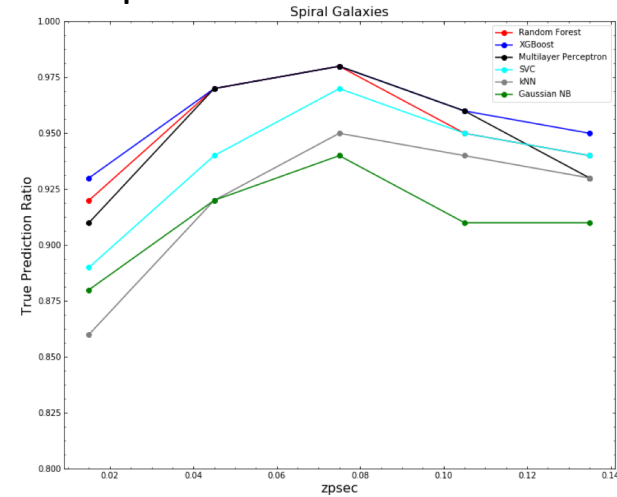
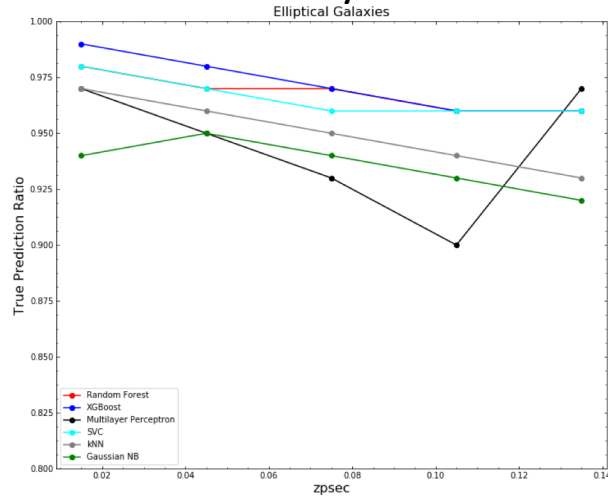
XGBoost - Parameters	
Feature	Importance
ui	0.584734
gr	0.182217
absMagR	0.171967
gi	0.061082

# Galaxy Zoo DATA: Accuracy distribution as a function of redshift

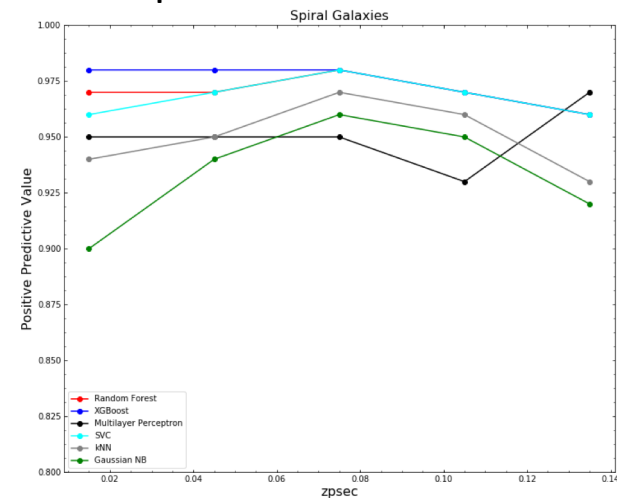
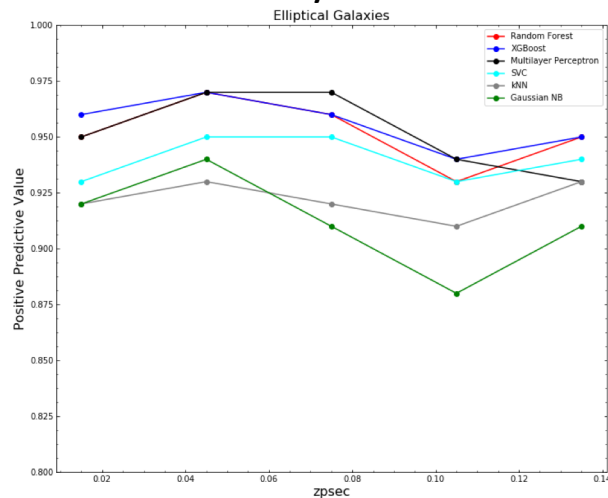


# Galaxy Zoo Data: Recalls and Precisions as a function of redshift

## Galaxy Zoo Data: Redshift Distribution of Recall for 8 parameters

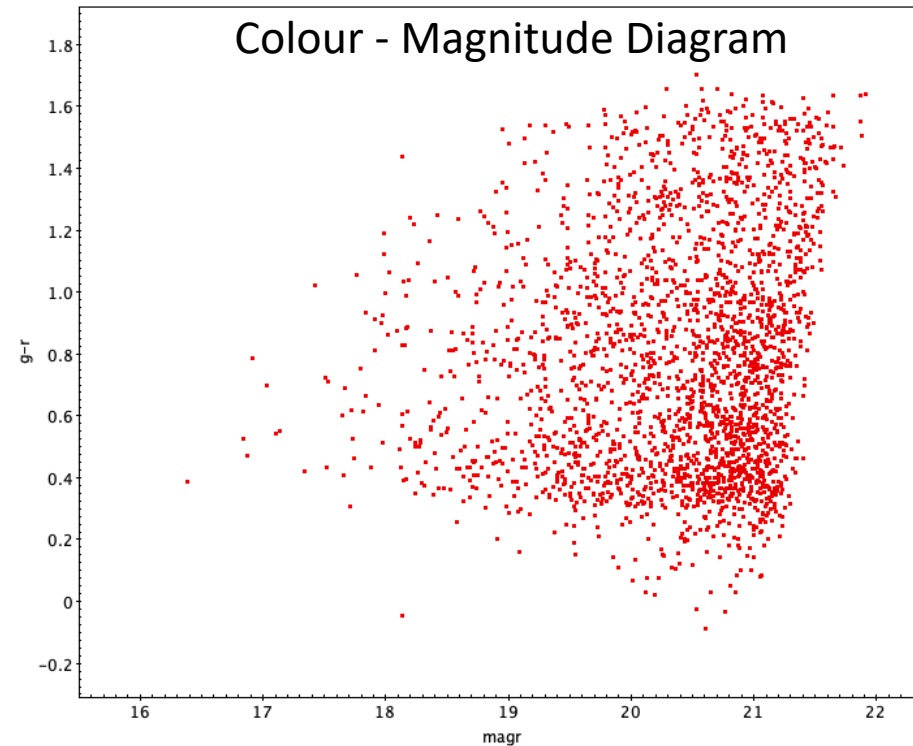
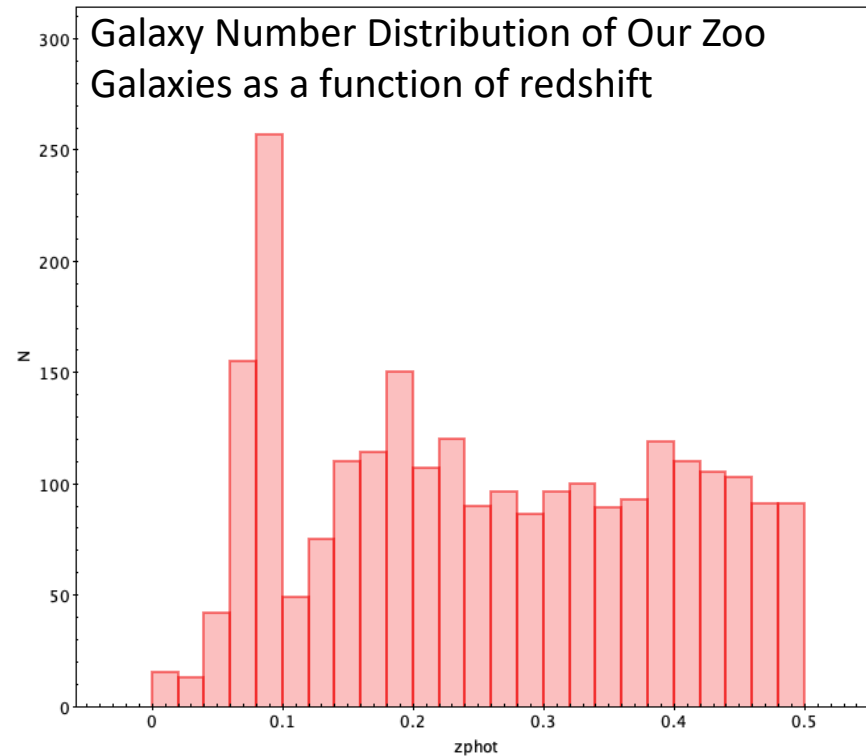


## Galaxy Zoo Data: Redshift Distribution of Precision for 8 parameters



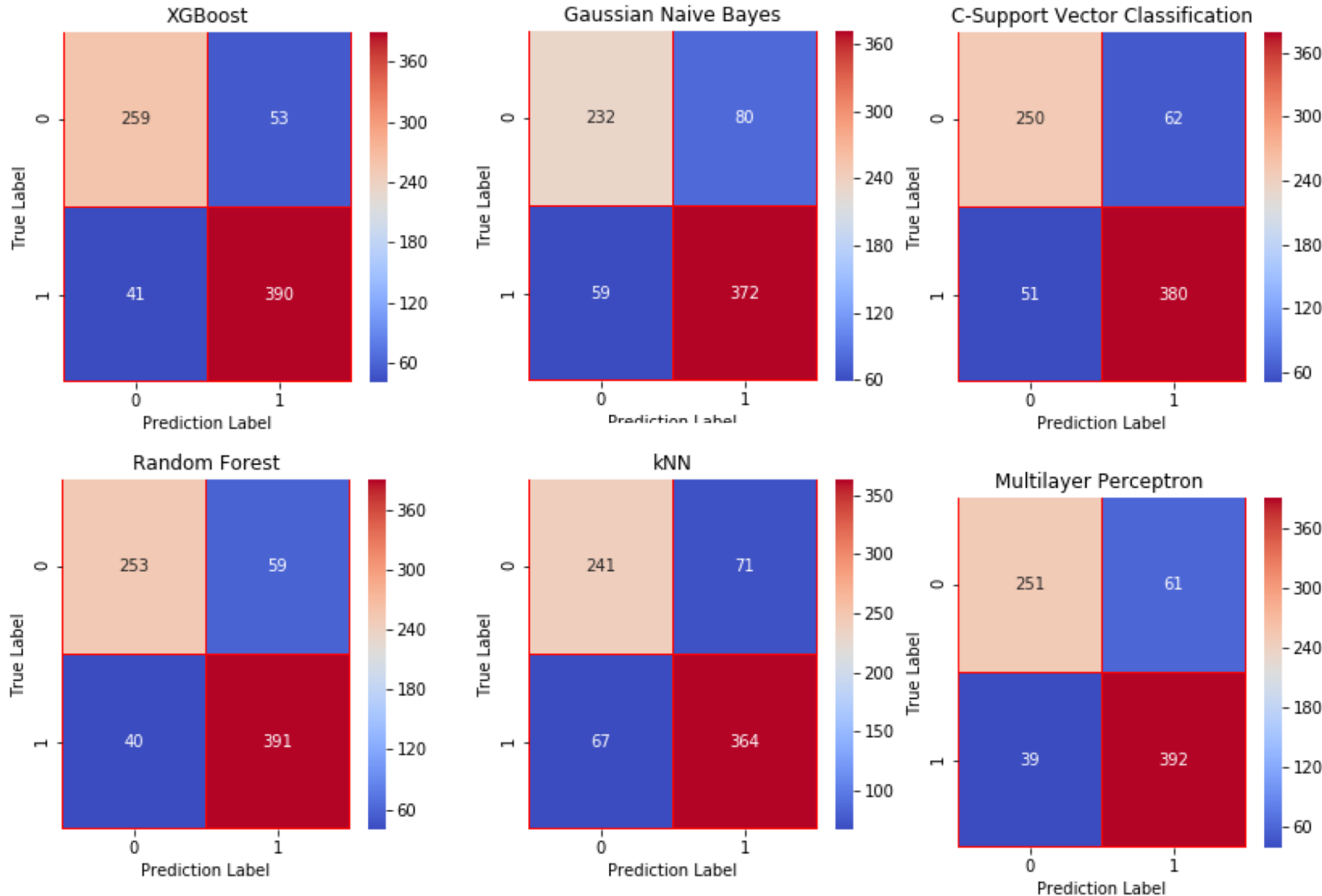
# Our Own Zoo: CFHTLS - W1 catalogue

- **2500** / 180000 galaxies visually classified
- Ellipticals: 1053, Spirals: 1423 -> 2476
- Redshift range:  $0 < z < 0.5$
- Train and Test sets (70/30 : 1733 / 743)



# Own Zoo DATA: Results for all 8 parameters

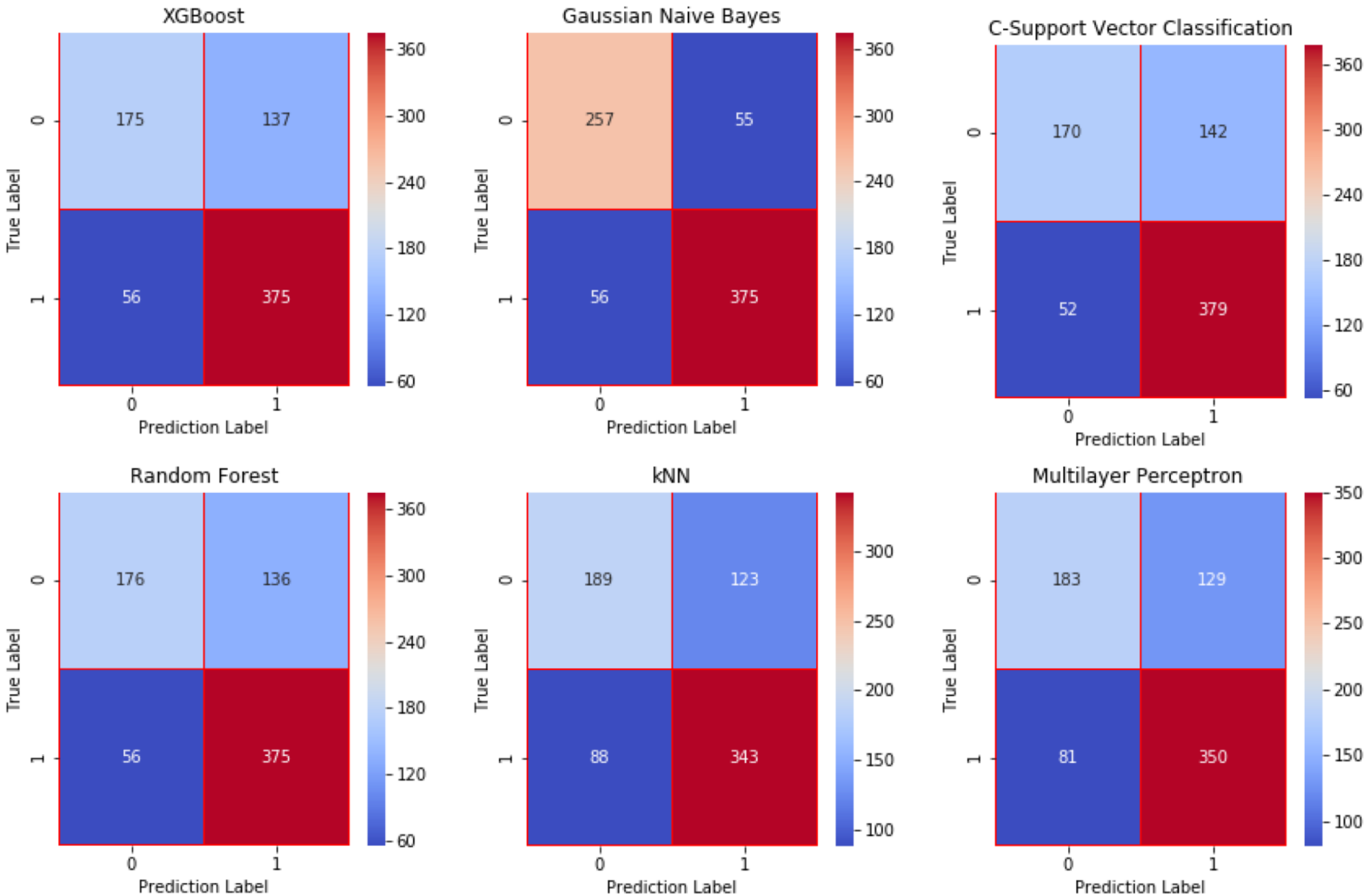
(ui, gr, gi, deVAB\_r(b/a), deVRad\_r(eff radius), CI (petroR50\_r/petroR90\_r), absMagR , sersic\_n)



Model	Score
Random Forest	0.875356
XGBoost	0.874794
NN	0.865559
SVC	0.854586
KNN	0.834990
Naive Bayes	0.812446

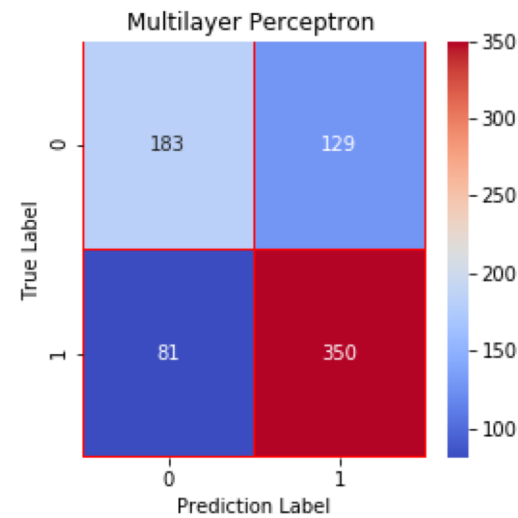
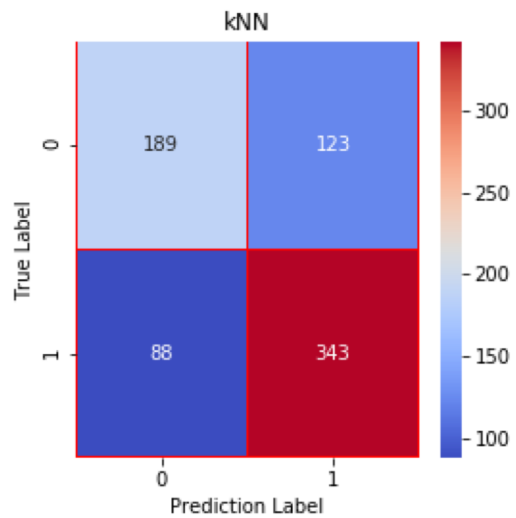
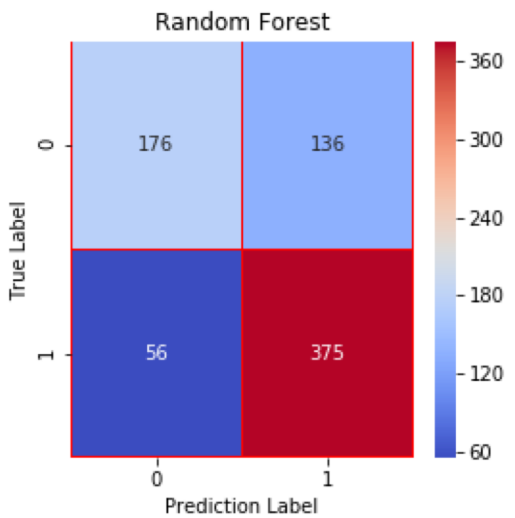
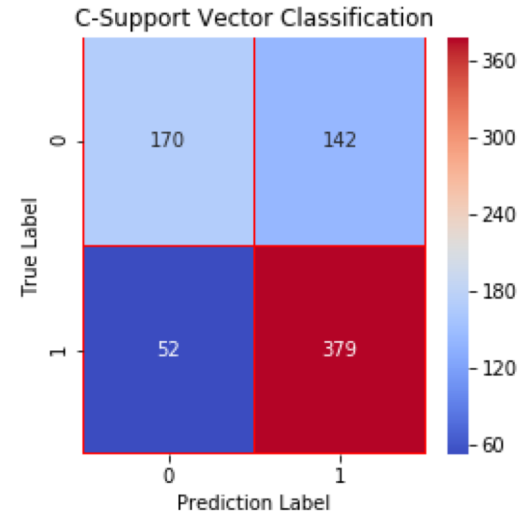
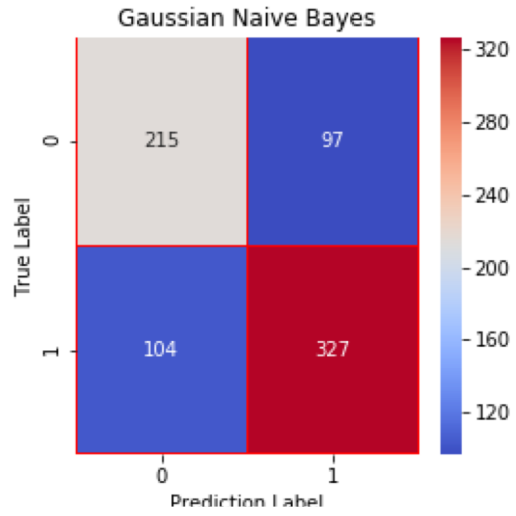
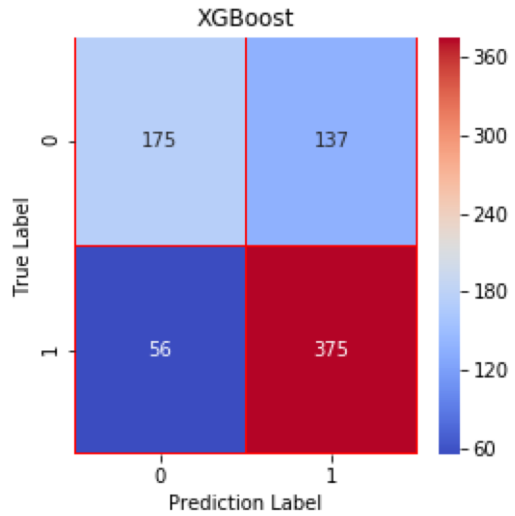


# Own Zoo DATA: Results for 4 structural parameters (deVAB\_r(b/a), deVRad\_r(eff radius), CI (petroR50\_r/petroR90\_r), sersic\_n)



Model	Score
NN	0.857467
XGBoost	0.856320
Naive Bayes	0.848819
Random Forest	0.847668
SVC	0.841901
KNN	0.830948

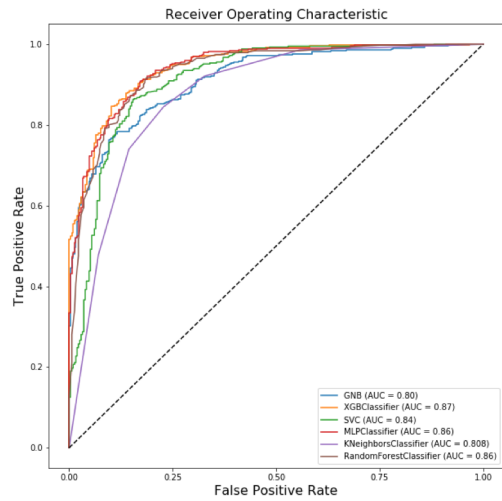
# Own Zoo DATA: Results for 4 photometric parameters (ui, gr, gi, absMagR)



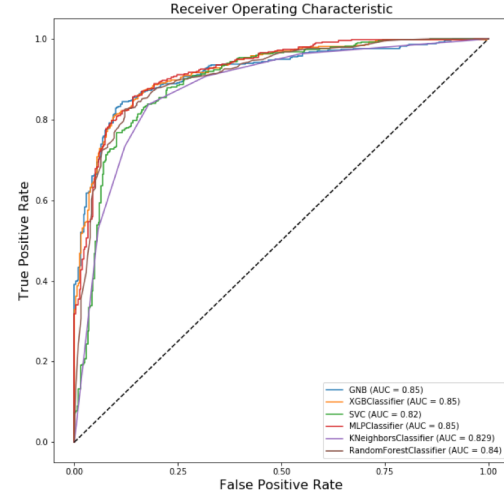
Model	Score
KNN	0.737433
SVC	0.735161
XGBoost	0.730503
Random Forest	0.721863
NN	0.721286
Naive Bayes	0.704537

# Our Own Zoo DATA: Comparing of ROC graphs of 6 Machine Learning Algorithms

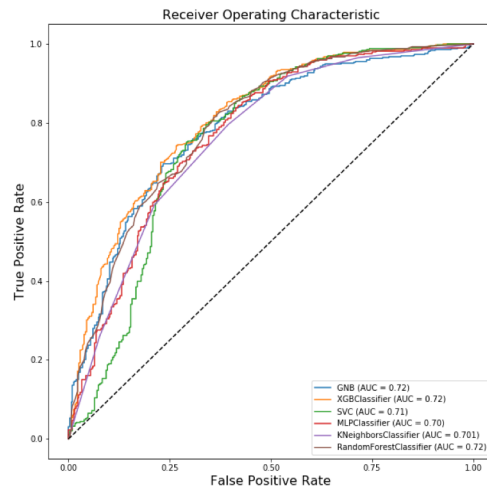
### All 8 parameters



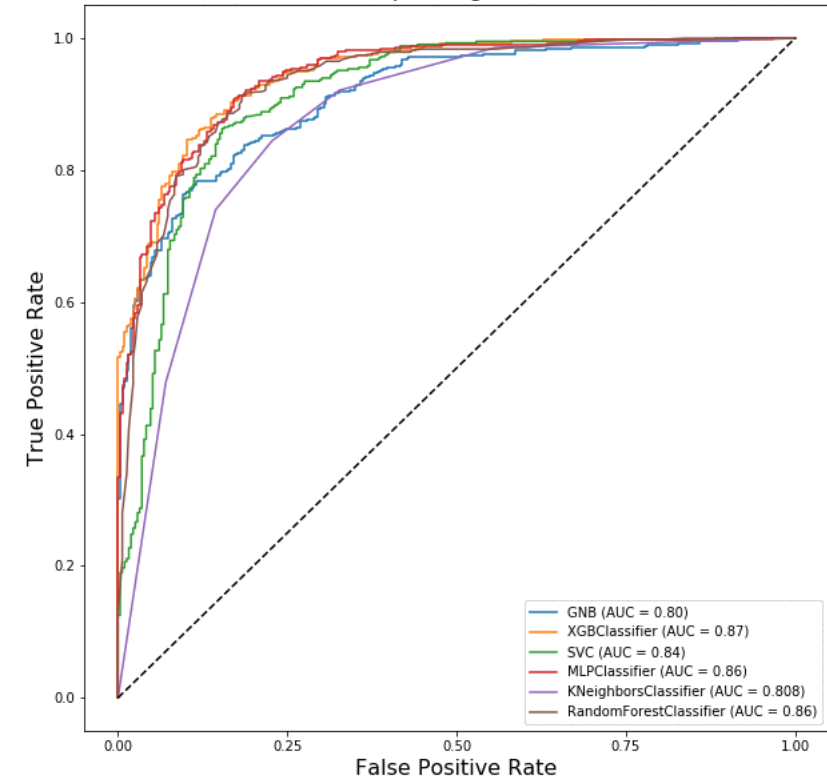
### 4 structural parameters



### 4 photometric parameters



### Receiver Operating Characteristic



# Impact of 8 Parameters in XGBoost and Random Forest

XGBoost - Parameters		Random Forest - Parameters	
Feature	Importance	Feature	Importance
Re	0.310302	Re	0.245362
ui	0.168082	ba	0.213058
ba	0.144646	sersic_r	0.118683
Cl	0.089578	ui	0.117339
sersic_r	0.089413	gr	0.091458
gi	0.079507	Cl	0.082964
gr	0.079207	gi	0.080279
absmagR	0.039266	absmagR	0.050857

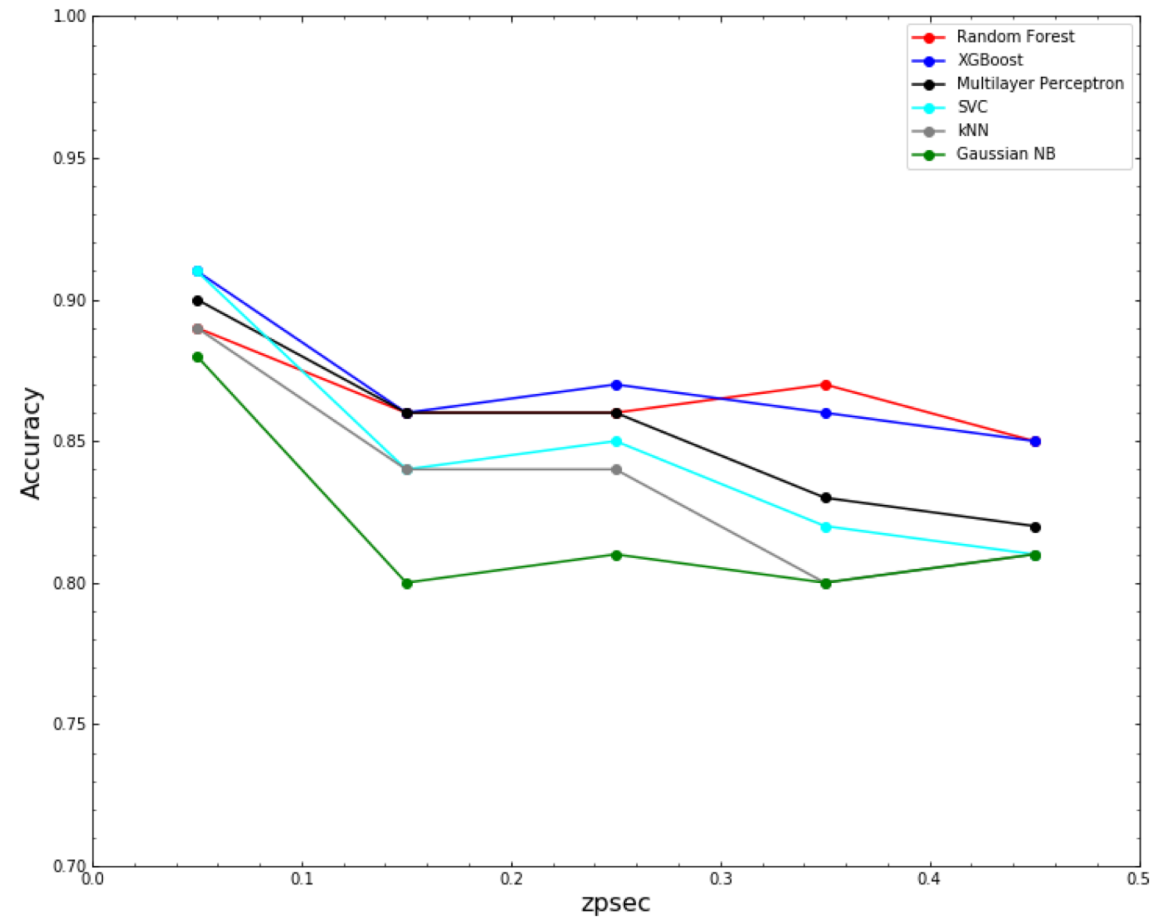
## 4 structural

XGBoost - Parameters	
Feature	Importance
Re	0.484281
ba	0.200457
sersic_r	0.163282
Cl	0.151981

## 4 photometric

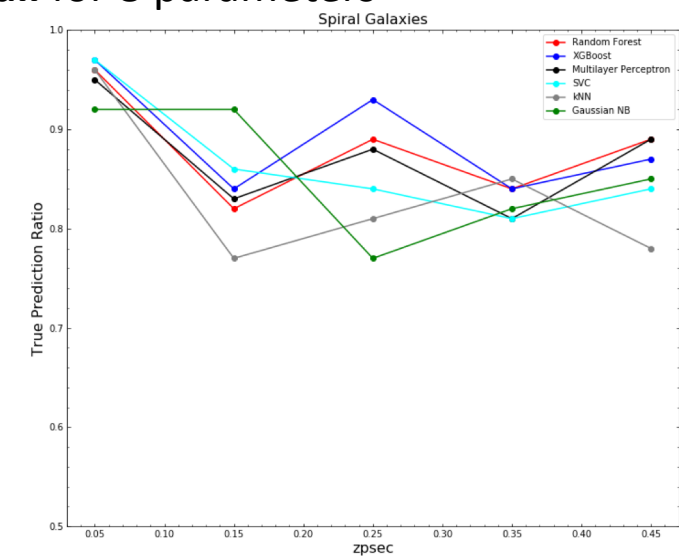
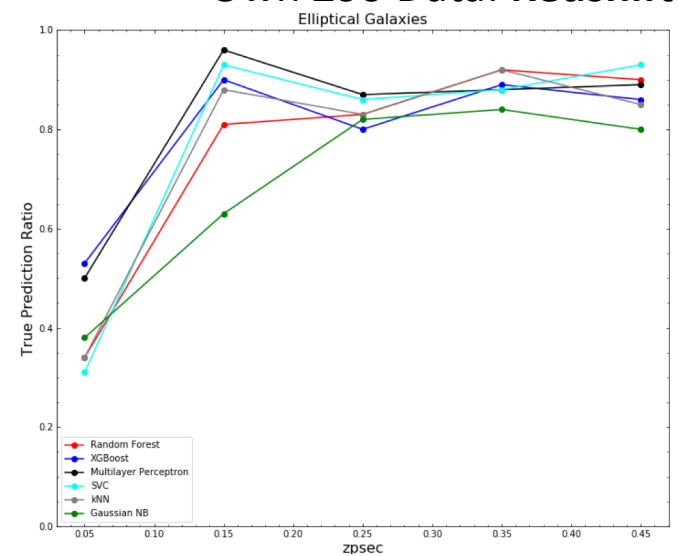
XGBoost - Parameters	
Feature	Importance
ui	0.346192
gr	0.255196
gi	0.226703
absmagR	0.171909

# Own Zoo – Accuracy distribution as a function of redshift

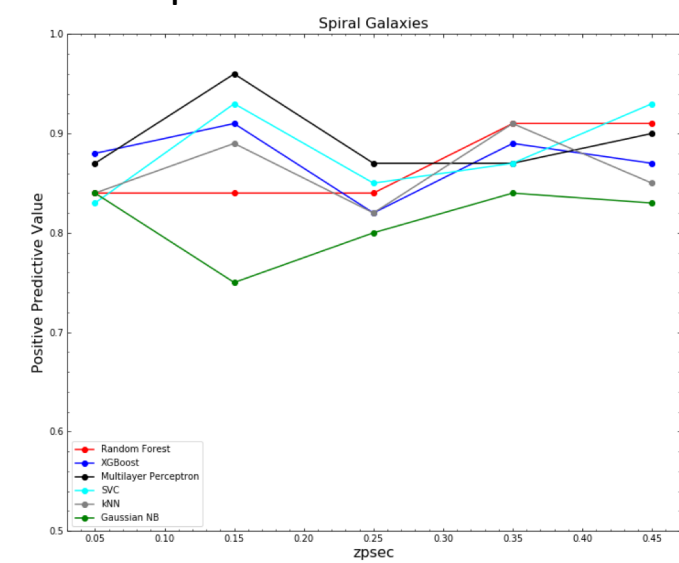
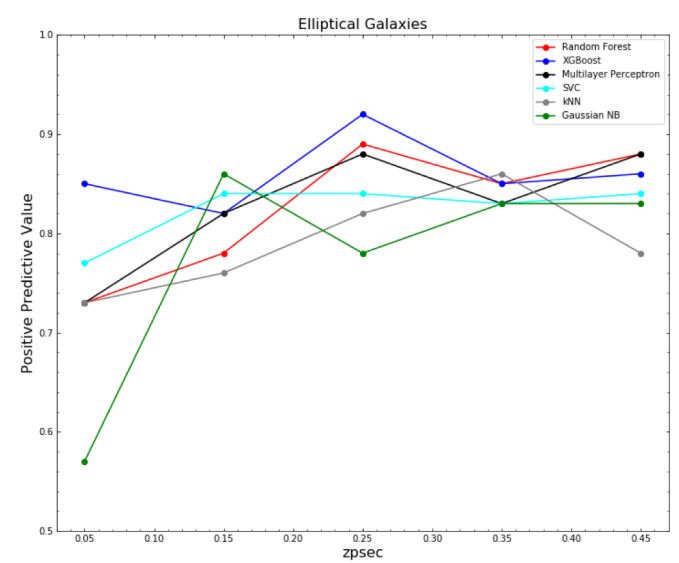


# Own Zoo Data: Recalls and Precisions as a function of redshift

## Own Zoo Data: Redshift Distribution of Recall for 8 parameters

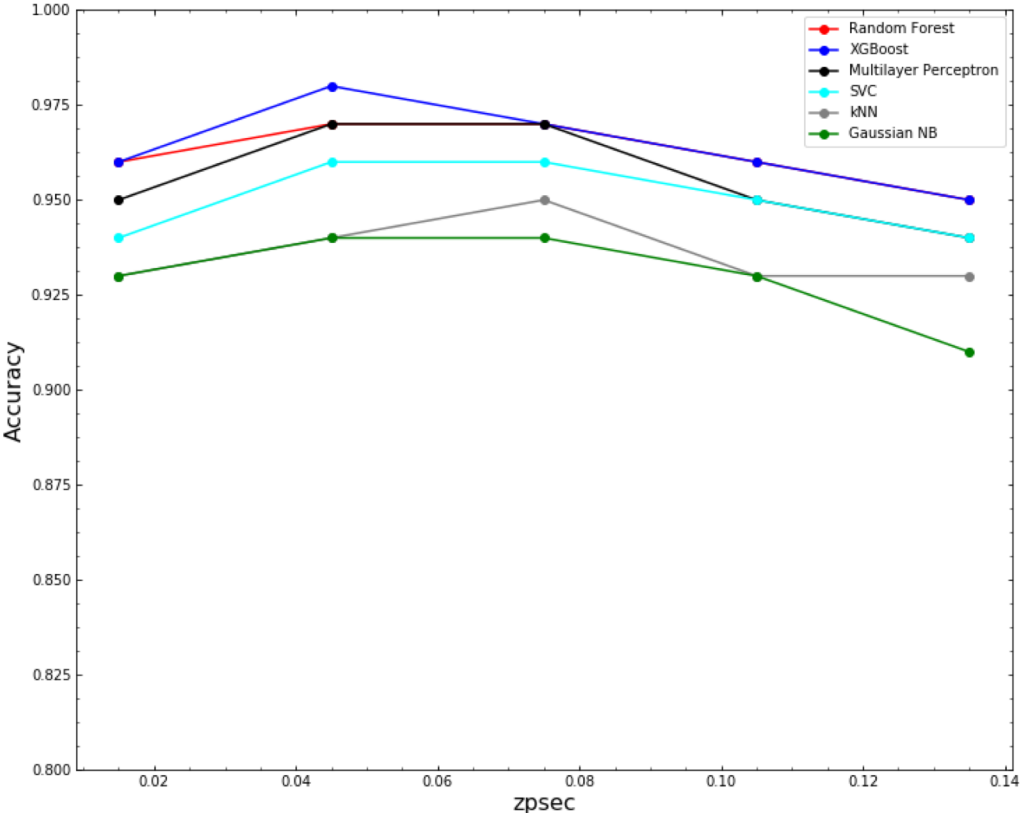


## Own Zoo Data: Redshift Distribution of Precision for 8 parameters

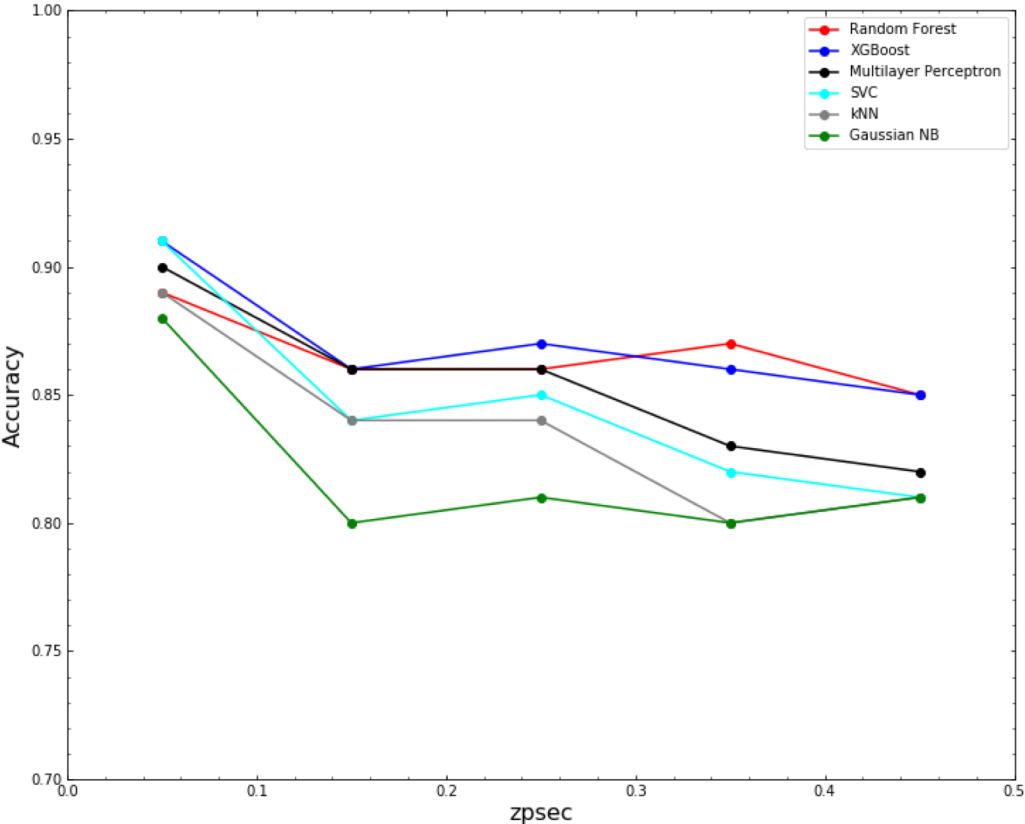


# Accuracy distribution as a function of redshift

### Galaxy Zoo Data $0 < z < 0.14$



### Own Zoo Data $0 < z < 0.5$



# Conclusions

- All 6 ML algorithms give almost the same results with 8 parametric features.
- The best accuracy scores are obtained from RF and XGBoost algorithms in two different data.
- The photometric parameters are less effective than structural parameters.
- To select of the parametric features is crucial than the ML algorithms itself and play very important role in the scores of accuracy.
- The accuracy is slightly decreasing with higher redshifts.
- After a certain redshift human eye won't be able to distinguish galaxy classes.
- More classes of morphological types means less accuracy performances.
- Future Work:
  - To extend visually classified sample and test the algorithms.
  - To choose a robust and effective parameters by using the PCA or features selection algorithms!
  - To apply the algorithms to the higher redshifted CFHLTS-W1 field galaxy sample with sersic indices.