

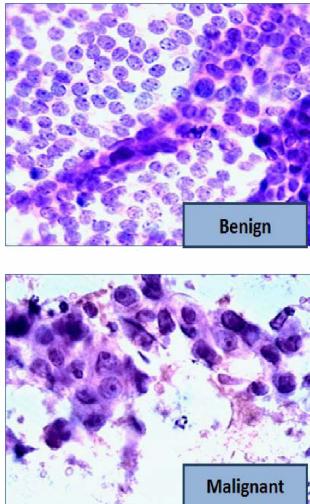
# **Predict breast cancer with machine learning**

**Fan Yang**

# Project summary

---

- Breast cancer is the second main cause of cancer death in women.
- Fine needle aspiration (FNA) is a common method of breast cancer exam that samples cells from a breast lesion with a fine needle.
- Breast cancer is diagnosed as either benign or malignant based on cell nuclei morphology.



- Dataset downloaded from Kaggle website contain cell nuclei features from 569 breast cancer patient.
- **Goal of project is to build classifier model detecting malignant breast cancer.**

Image cited from Ahmad *et al.* 2013

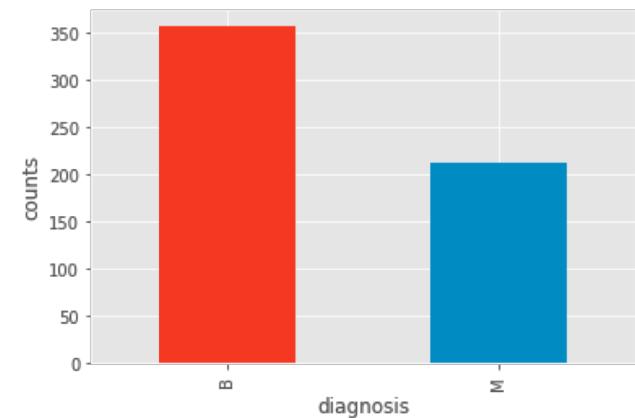
# Dataset overview

df.head()												
	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave_points_mean	...	...
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	...
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	...	...
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	...
3	84348301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	...
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	...	...

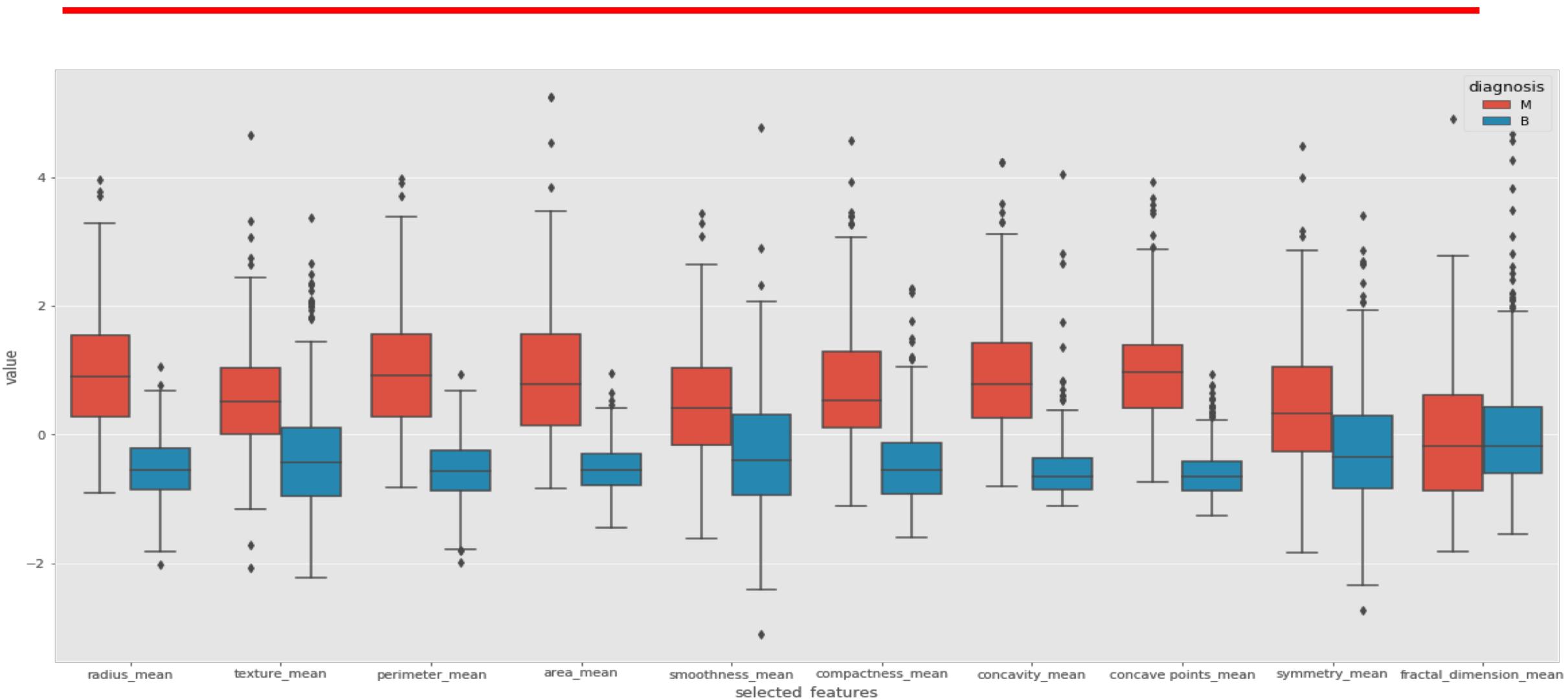
- 1) ID number
- 2) Diagnosis (M = malignant, B = benign)

Nuclei morphology is represented by following features:

- a) radius (mean of distances from center to points on the perimeter)
- b) texture (standard deviation of gray-scale values)
- c) perimeter
- d) area
- e) smoothness (local variation in radius lengths)
- f) compactness ( $\text{perimeter}^2 / \text{area} - 1.0$ )
- g) concavity (severity of concave portions of the contour)
- h) concave points (number of concave portions of the contour)
- i) symmetry
- j) fractal dimension ("coastline approximation" - 1)

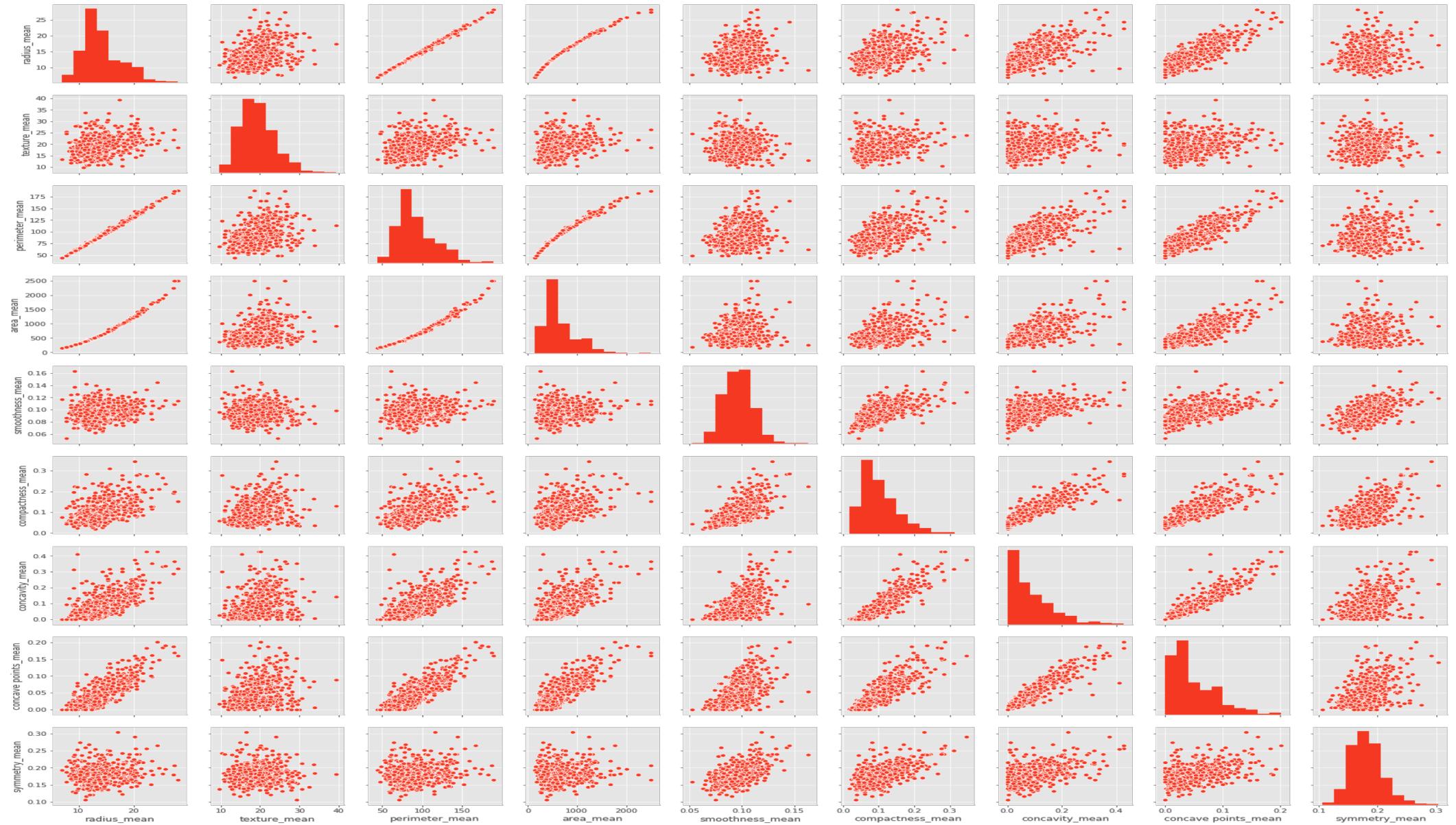


# Differential distribution of cell nuclei features between benign and malignant samples

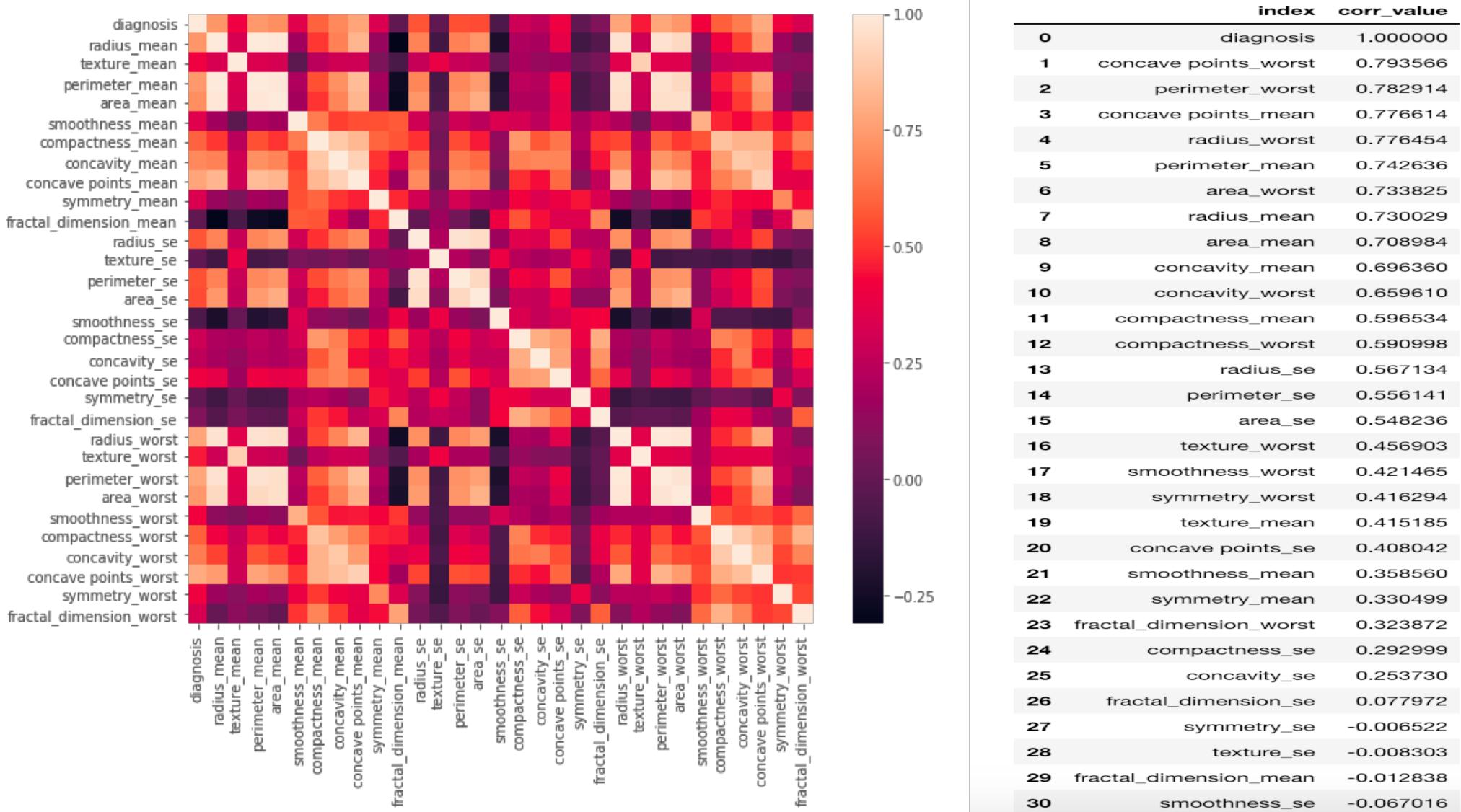


# Relationship among numerical variables representing cell nuclei morphology

---

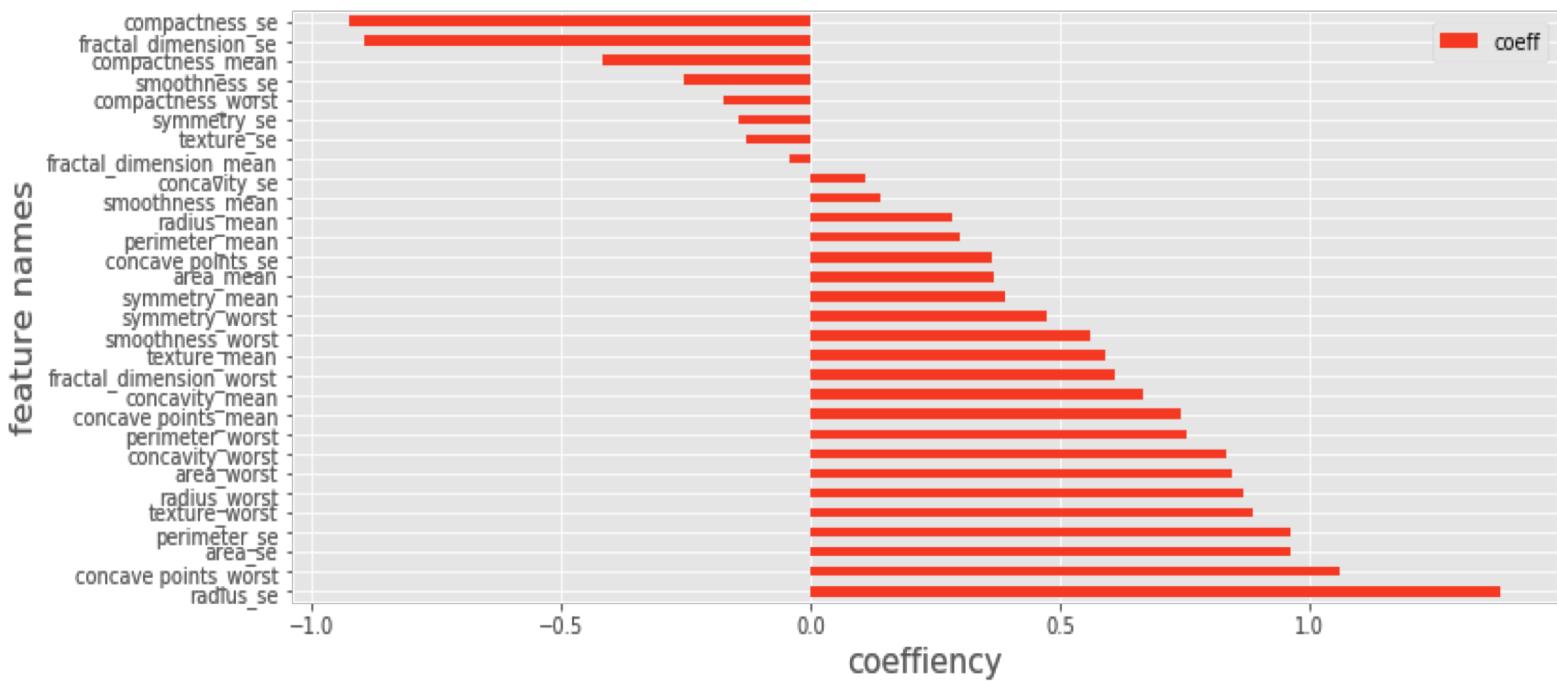
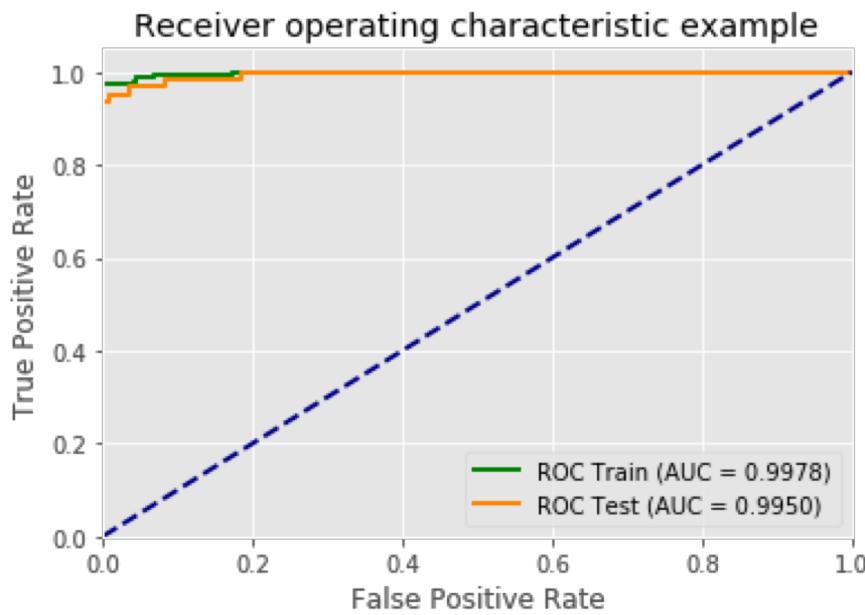


# Correlation between breast cancer diagnosis and cell nuclei features



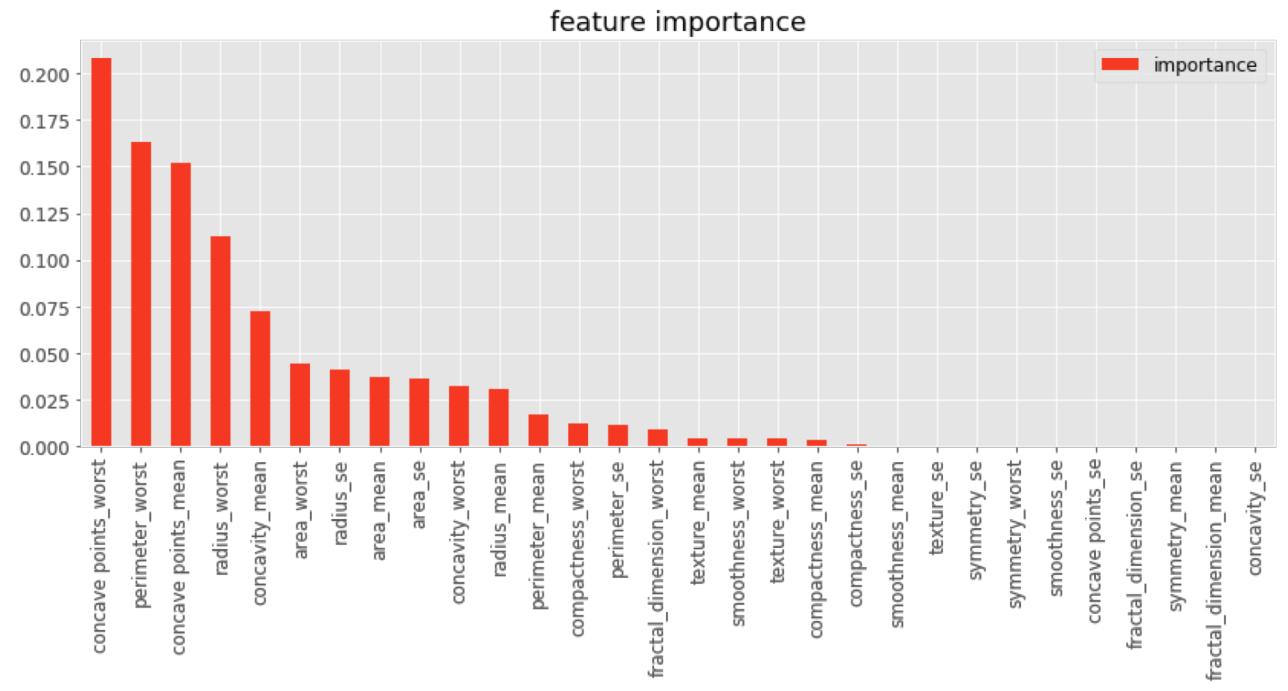
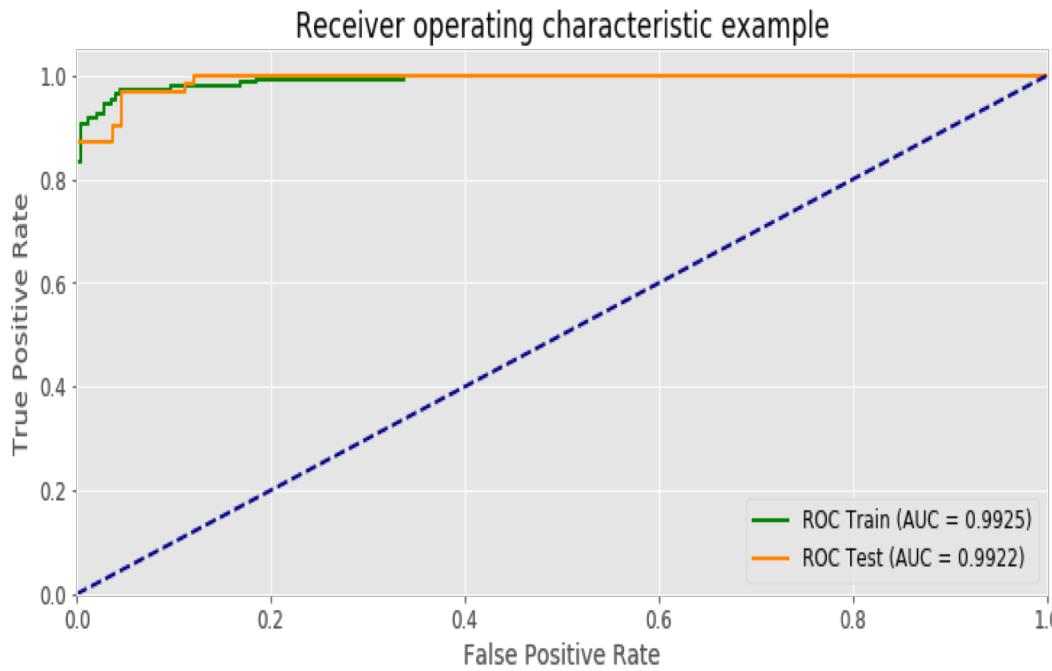
# Build Logistic regression model detecting malignant breast cancer

---



# Build Random forest model detecting malignant breast cancer

---



# Project conclusions

---

1. Morphological features of cell nuclei are criteria for breast cancer diagnosis.
2. This project explores how these features correlate to each other and how they cause a differential distribution between benign or malignant cancer samples
- .
3. Logistic regress and random forest models were built to detecting malignant breast cancer
4. Features importance reveal that concave points and radius are the most useful information for image-based breast cancer diagnosis.